



LEGAL NOTICE:

© Copyright 2008 to 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 to 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
c/o VTM, Inc.
3855 SW 153rd Drive
Beaverton, OR 97003
USA
info@nvmexpress.org

Technical input submitted to the NVM Express® Workgroup is subject to the terms of the NVM Express® Participant's agreement. Copyright © 2008 to 2022 NVM Express, Inc.

NVM Express® Technical Proposal (TP)

Technical Proposal ID	TP8012 Boot Specification
Revision Date	2022.11.15
Builds on Specification(s)	NVM Express Base 2.0b
Affects Specification(s)	NVM Express Boot Specification (New) NVM Express Base Specification NVM Express Command Specification NVM Express PCIe Transport Specification NVM Express RDMA Transport Specification NVM Express TCP Transport Specification NVM Express Management Interface Specification NVM Express Zoned Namespace Command Set Specification NVM Express Key Value Command Set Specification
References	

Technical Proposal Author(s)

Name	Company
Doug Farley	Dell Technologies
Phil Cayton	Intel
Lenny Szubowicz	Red Hat
Rob Davis	NVIDIA
Stephanie Yuan	HPE
Wenhua Liu	VMware
Maciej Rabeda	Intel
Vincent Zimmer	Intel
Curtis Ballard	HPE
Fred Knight	NetApp

Technical Proposal Overview

This proposal for creating a body of NVMe-oF Boot Collateral, which includes both a new standard for NVMe-oF Boot of Namespaces on NVMe-oF Subsystems, and an informative NVMe-oF Boot Document. This collateral leverages standardized network boot mechanisms to describe how NVMe-oF Hosts can discover, attach to, load OS, and boot from NVMe-oF Subsystems utilizing supported transports. An initial approach will focus on enabling NVMe/TCP with other transports to follow.

Document Revisions

Date	Comments
08/17/2021	<p>bit/byte table changes:</p> <p>Many Heap Objects were identified Length and Offset but not all. Normalized all heap format to Offset, Length. This reordered offsets.</p> <p>HFI Descriptor – Added Length Field after ACPI Version Field BEFORE Flags, consumed reserved byte from AFTER flags. This reordered all HFI offsets.</p> <p>HFI TCP Info 1) Added Length Field after Structure Version, 2) moved HFI Transport Type after Length (Needs Discussion). This reordered all HFI offsets.</p> <p>Added length field into Extended Adapter Info table. This reordered EAI offsets.</p> <p>Language around back-propagation from OS on nBFT viability.</p> <p>Added a handful of header clarifications to differentiate handling of multiple nBFT's.</p> <p>Added new heap object for driver image path string. Consumes reserved bytes.</p> <p>Clarifying Text for Route Metric.</p> <p>Optional fields defined with (Optional) tag and mandatory default values when unused.</p>
10/19/2021	<p>Added Unavailable Namespace bit in SSNS</p> <p>bit/byte table change: Added a new sub-structure, the SNSS Extended Info Structure as Heap overflow for SSNS fields. This moved elements out of the SSNS into the Extended Structure and reordered both the SSNS and Flags.</p> <p>bit/byte table change: Discovery Structure converted into sub-structure format like the other fields. NumDisc reordered offsets in part of Control structure.</p>
11/2/2021	<p>Added Definitions section.</p> <p>Added second bit to SSNS Unavailable Namespace Flag to clarify specified vs. unspecified behavior from the pre-OS driver.</p> <p>Changed all Flag wording from “Enabled” to “Selected” when Unselected/Selected is used to indicate Flag state.</p>
11/16/2021	<p>bit/byte table change: Changed signature to NBFT from nBFT</p> <p>Added clarifying verbiage in Discovery structure that Index 0 is reserved for non-discovery cases in the SSNS.</p>
12/17/2021	<p>Add HFI Reference</p> <p>bit/byte table change: Transport Type Value realigned with TRTYPE from NVMe Base DLP Fields. URI Format removed suffix-less variants.</p> <p>bit/byte table change: Transport Transceiver ID Changes</p> <p>bit/byte table change: SSNS got a separate transport specific flags field</p> <p>bit/byte table change: SSNS Offsets reordered.</p> <p>Added HostID/NQN Configured Flags</p> <p>Overall should/shall cleanup round 2.</p>

1/4/2022	Minor notes from 1/4 BWG meeting.
1/12/2022	Added Definitions, cleaned up HostID
1/26/2022	bit/byte table change: Second revision at merging Informative Annex Content Merged in Wenhua's InfoAnnex 7.1.2.1 content Template and Formatting changes Clarified null termination of heap objects needs to be included in object length Added version to HFI transport type Lots of explanatory text from James Added Host Primary flag to host flags
2/14/2022	bit/byte table change: Added Length and Version to Control structure, removed from all primary tables, removed length from all Heap tables.
2/15/2022	Updated document based on feedback/suggestions from Phil.
2/15/2022	Notes from FMDS and BootTG meeting/review.
2/17/2022	Further reformatting to NVMe Style and editorial from PaulS.
2/22/2022	Updated Figure 5 (NBFT Table Header) to NVMe style. All tables may have experienced bit/byte changes. Inserted conventions sections 1.4.1 ... 1.4.3 from Base Spec. Changed all table captions to Figure captions and updated all cross-references. Updated Table of Contents and Table of Figures. Changed double spaces to single. Deleted spaces at the end of paragraphs. Other miscellaneous changes are marked with comments. Fixes from JohnG and updates to the TOC formatting to correct Word Style. Changed previous changelog term 'Breaking' to 'bit/byte table change' for clarity that this line item directly impacted the table data structure. Fixed a lot of broken references.
2/22/2022	TG Notes and Feedback on changes. Extracted BootTG identified topics for future TPARs to Kavi. Fixed stray references Collapsed a number of info-annex sections.
2/23/2022	Fixed stray references Fixed grammar/spelling/capitalization/spacing issues Identified incomplete sections needing attention
2/23/2022	Converted more figures to NVMe format.
2/25/2022	More work on figures. Revised annexes.
2/28/2022	Changed terminology from "structure", "sub-structure", etc. to "descriptor". Made consistent the naming of index fields inside descriptors. Numerous editorial changes. Added comments for items needing review.
3/2/2022	Changed heap object offset/length fields to common Heap Object Reference field format. Padded structures of length 17 bytes or greater to a multiple of 32 bytes. Mandatory fields use "shall" requirements in their definitions. Optional fields use "if" clauses in their definitions. The O/M column is left for checking, but may be removed. Replaced "Driver Device Path Signature" with "Driver Signature".
3/4/2022	Figure 6 (NBFT Layout): Corrected the collection of heap objects. Figure 8 (NBFT Table): Converted the Driver Signature Offset & Length fields to a Heap Object Reference. Removed "heap object" from the descriptions of the offset fields for HFI Descriptor List, SSNS Descriptor List, Security Profile Descriptor List, and Discovery Descriptor List. Wording changes in Transport Types section. Correct Host Descriptor HostID/NQN fields as they are mandatory. Attempted to clean up Route Metric. A.3 intro paragraph and formatting. References to conventions for URI Figure 27 – edk2 flow Figure 28 – reference Removed Extended Auth Flag, waiting on TP8017

	<p>Restored lost administrative override flags for Auth and Security Policy flags; combined into one bit pair rather than separate</p> <p>Reword Controller ID in SSNSEX</p> <p>Restored lost dhcp string heap object into snssex</p>
3/7/2022	<p>SSNS Descriptor: Deleted duplicate SSNS Extended Information Descriptor Heap Object Reference.</p> <p>Changed all instances of “this field specifies” to “this field indicates”</p> <p>Changed instances of “zero” to “0h” as appropriate.</p> <p>Changed “section A.” to “Annex A.”.</p> <p>Changed all remaining instances of SSNS to SSNS.</p> <p>Changed instances of “structure” to “descriptor” when it referred to a named descriptor.</p> <p>Figure 1: revert image to remove the entrance of NVMe-Boot from the NVMe family of specs (for now).</p> <p>Figure 6: Shade/color the figure.</p> <p>Figure 17: rewrite to remove the term “primary boot record” from the Non-bootable Entry Flag.</p> <p>Note to host/OS that a connection in pre-os may have not been formed</p> <p>Probably entirely ‘phase 3’ edits, however.... Noted here in case somebody gets time/bored:</p> <p><global> search/replace “A value of 0” and update to “If cleared to 0h”</p> <p><global> ensure consistency with “pre-os” capitalization</p> <p><global> search and replace “it” as appropriate throughout the document (but might be good for another try later)</p> <p><global> search/replace “Header”, “Control”, etc. capitalization. Also “discovery controller” “I/O controller”, etc. capitalization as appropriate. Probably missed a few out of abundance of caution.</p> <p><global> search for all the “e.g.” instances and remove “etc.” from the end.</p>
3/17/2022	<p>Title change per TWG and Judy’s note – from ‘NVMe Boot Specification’ to ‘Boot Specification’ (implied NVM Express from gfx header)</p> <p>Word choice changes in HostID and HostNQN extended descriptions.</p> <p>Reformatted HFI-Adv-Capabilities to bullet format</p> <p>HostID/NQN Override, remove useless ‘must’ in non-normative statement.</p> <p>Case and typo corrections throughout.</p> <p>Corrected missing Mandatory from 0.52 Port ID</p> <p>Corrected Introduction for Generic NVMe Boot</p> <p>Annex B (UEFI ECR Content) finally removed, separate ECR now</p>
3/29/2022	<p>Narrow NS ID type - Only support Namespace UUIDs with NVMe boot – RFC4122 UUID GUID defined in the NVMe Base Spec (CNS3 NIDT 03H specifies only supports NS UUID; cannot be extended or global unique ID)</p> <p>Figure 16 (NIDT field, mandatory value for NID)</p> <p>Figure 19 (updated NID to delete the words “Globally Unique ID”)</p> <p>Figure 29 (SSNS – remove NSID, modify NIDT, NID fields, remove note)</p> <p>remove ‘A4.1.3 NID Types’ (as we excluded EUID64) based on discussion</p>
4/6/2022	<p>Reverted NGUID and EUI deprecation, added 1.5.x sections for string format including decorators per TWG 3/31</p> <p>Removed Per-HFI HostID and HostNQN overrides for future TPAR per feedback</p> <p>Removed HFI Advanced Capabilities Bits, Descriptors to Reserved Status (future TPAR)</p> <p>Added A4.2.4 NID decoding by drivers or environments</p> <p>Fix byte offsets in Security Profile Descriptor (Cipher Suite Offset 4 vs. 6bytes)</p> <p>Added tail page with Reference to Boot Spec from Base Spec</p> <p>Added TP Header like TP4015 used for new doc + spec change</p>
5/10/2022	<p>Added Annex Doc header styles</p> <p>Correct doc references</p> <p>Started edits to references format from ECN105.</p> <p>Removed Annex Informative A.2.2.x Use of the Adapter Capabilities Bits Section, Normative content was removed in 0.57 and we missed the informative stuff.</p> <p>Fix byte offset math in HFI Info descriptor to 4B aligned</p> <p>Globally case change Pre-OS to pre-OS.</p>
5/11/2022	<p>Converted UEFI Diagram to embedded visio from word art</p> <p>Attempt to fix navigation pane hierarchy</p> <p>Fix missing SSNS descriptor labels in figure 29</p>
05/16/2022	<p>General document editing – changed/rationalized capitalization throughout, fixed grammar issues, updated informative annex text</p> <p>Large must to shall or informative statement overhaul.</p>

05/17/2022	Add in Figure 1 changes for all primary specifications in NVM Express – in TP beyond Annex A content General editing from 5/17 BTG meeting – capitalizations, tense. New TP4126 landing zone identified for Section 3.x before NBFT.
05/19/2022	Example offset values in the HFI, SSNS and DISC descriptor lists in Figures 10, 15 and 24 were not matching the actual length of the descriptors definitions. Grammatical and capitalization changes throughout Term corrections to example calculations in section 3.1.2. Lots of rewording of sentences in the annex. Merged collected review feedback.
05/24/2022	Added/fixed references throughout. Found/commented on Figure 14 (HFI Extended Adapter Information Descriptor) and associated text that has been removed but still is referenced. Need to refactor figures and references to figures for any figure past figure 13. 3.1.2.4.1 HFI – convert “will always” to “shall”. First pass integration of new Linux informative section from Lenny, needs further edits
05/26/2022	Highlighted “Figure” and Section hyperlinks for Editor References to “Heap” that are not Field Names to lower case
06/1/2022	Fix (x26) superscript ® on NVM Express Base Specification Convert stale (x12) NVMe Base Specification to NVM Express® Base Specification Drop incorrect ‘e’ on PCI Firmware Specification reference
07/24/2022	Addressed 45-day review comments from Samsung and Intel. All comments/changes are editorial in nature. SNSS ‘Security Field Flag’ renamed ‘Security Flag’. Added notes/suggestions for the Editor to address.
08/01/2022	Renamed HFI “PCI Address” to “PCI Express Routing ID for the HFI Transport Function”.
08/11/2022	Integrated
09/29/2022	Reviewers completed review – submitted to Editor
10/19/2022	Integrated edits from ratification review
11/07/2022	Fixed the section number styles Changed “a NVMe...”to “an NVMe..” Changed “a ‘M’” to “an ‘M’” Fixed spacing in tables.
11/08/2022	Removing the Table of contents to avoid material from other specifications being included. Changed the Header and Footer to be TP specific as opposed to final specification
11/15/2022	Added a period to a list in section A.4.1.2. Removed the revision column of this table.



NVM Express®

Boot Specification

Revision 1.0

November 8th, 2022

Please send comments to info@nvmexpress.org

NVM Express® Boot Specification, Revision 1.0 is available for download at <https://nvmexpress.org>. The NVM Express Boot Specification, Revision 1.0 incorporates TP8012.

SPECIFICATION DISCLAIMER

LEGAL NOTICE:

© Copyright 2008 to 2022 NVM Express, Inc. ALL RIGHTS RESERVED.

This NVM Express Boot Specification, revision 1.0 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 NVM Express Boot Specification 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 to 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 153rd Drive
Beaverton, OR 97003 USA
info@nvmexpress.org

1 Introduction

1.1 Overview

The NVM Express® (NVMe®) Base Specification defines an interface for host software to communicate with non-volatile memory subsystems over a variety of memory-based transports and message-based transports.

This document defines constructs and guidelines for booting from NVM Express interfaces.

The current revision defines extensions to the NVM Express specification family that enable booting over NVMe over supported fabric transports. This specification includes:

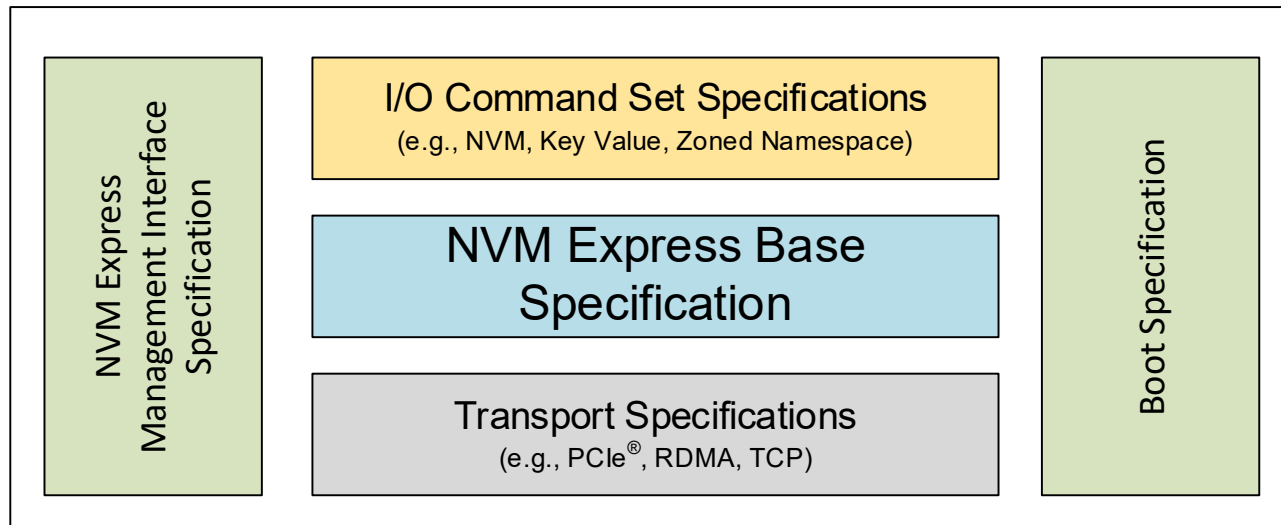
- material specific to the NVMe-oF TCP transport; and
- operating system best practices and their consumption of pre-boot configuration (e.g., UEFI).

Other published specifications referred to in this document, even if required for compliance, are outside the scope of this specification; this includes published specifications for fabrics and other technologies referred to by this document or any NVM Express Transport binding specification.

1.2 Scope

Figure 1 shows the relationship of the NVM Express® Boot Specifications to other specifications within the NVMe® family of specifications.

Figure 1: NVMe Family of Specifications



The NVM Express specification family structure shown in Figure 1 is intended to show the applicability of NVM Express specifications to each other, not a hierarchy, protocol stack, or system architecture.

The NVM Express Base Specification defines a protocol for host software to communicate with non-volatile memory subsystems over a variety of memory-based transports and message-based transports.

The NVM Express Management Interface Specification defines an optional management interface for all NVM Express Subsystems.

NVM Express I/O Command Set specifications define data structures, features, log pages, commands, and status values that extend the NVM Express Base specification.

NVM Express Transport specifications define the binding of the NVMe protocol including controller properties to a specific transport.

The NVM Express Boot Specification (i.e., this specification) defines constructs and guidelines for booting from NVM Express interfaces.

1.3 Definitions

Key Terms and acronyms used throughout this document. Definitions pertinent to domain specific terms are additionally in Annex A.1.

1.3.1 Boot Device

The storage volume which contains the relevant instructions and necessary next step content that a Platform may load an operating system or preform other basic functions beyond that of a pre-OS environment. In UEFI, this is commonly a device with an EFI System Partition.

1.3.2 Discovery Controller (DC)

A controller that exposes capabilities that allow a host to retrieve a Discovery log page. A Discovery controller does not implement I/O Queues or provide access to a non-volatile memory storage medium. Refer to the NVM Express Base Specification.

1.3.3 Diffie-Hellman (DH)

A method of securely exchanging cryptographic keys (refer to RFC 2631).

1.3.4 Discovery Log Page (DLP)

A log page is provided by a Discovery controller contains one or more entries. Each entry indicates information necessary for the host to connect to an NVM subsystem. Refer to the NVM Express Base Specification.

1.3.5 EFI System Partition (ESP)

An EFI defined partition in a UEFI compliant system (refer to the Unified Extensible Firmware Interface Specification) in which all OS loaders and applications will be stored and read by pre-OS Firmware for the purposes of loading or booting an UEFI compatible operating system. On Linux systems this might be represented by a directory structure under /boot/efi for example. For UEFI complaint Platforms, the presence of an ESP makes a namespace bootable.

1.3.6 Heap Space

Refer to section 3.1.1.1.

1.3.7 Host Fabric Interface (HFI)

A Host Fabric Interface is a structure that contains information about various supported fabric interfaces and their transport specific characteristics. These are logical or physical interfaces to local or remote access as defined by the HFI.

1.3.8 iSCSI Boot Firmware Table (iBFT)

A block of information that contains parameters that are useful to the iSCSI Boot process defined in the iSCSI Boot Firmware Table (iBFT) Specification. The iBFT is the mechanism by which iSCSI Boot parameter values are conveyed to the operating system.

1.3.9 Namespace (NS)

The representation of a data device that NVMe-oF hosts attach to in order to locate Bootable and Non-bootable UEFI Block or File Entries.

1.3.10 Namespace ID (NSID)

The namespace ID is used by a controller to provide access to a namespace as defined in the NVM Express Base Specification.

1.3.11 Namespace Identifier (NID)

The Namespace Identifier is a globally unique value defined in the Namespace Identification Descriptor list (i.e., CNS 03h) in the NVM Express Base Specification.

1.3.12 NVM Subsystem

An NVM subsystem includes one or more domains, one or more controllers, zero or more namespaces, and one or more NVM subsystem ports. An NVM subsystem may include a non-volatile memory storage medium and an interface between the controller(s) in the NVM subsystem and non-volatile memory storage medium. Refer to the NVM Express Base Specification.

1.3.13 NVM Subsystem Port

An NVMe over Fabrics protocol interface between an NVM subsystem and a fabric. An NVM subsystem port is a collection of one or more physical fabric interfaces that together act as a single interface. Refer to the NVM Express Base Specification.

1.3.14 NVMe Boot Firmware Table (NBFT)

The NVMe Boot Firmware Table is the data structure defined by this specification as a means of passing context from a pre-OS Boot environment to an OS runtime.

1.3.15 NVMe over Fabrics (NVMe-oF)

An implementation of the NVM Express interface that complies with either the message-only or the message/memory-based implementation of the memory-based transport model definition. Refer to the NVM Express Base Specification.

1.3.16 NVMe Qualified Name (NQN)

Used to uniquely describe a host or NVM subsystem for the purposes of identification and authentication. Refer to the NVM Express Base Specification.

1.3.17 Platform

A computer or hardware device and/or associated operating system, or a virtual environment, on which software is installed or run (refer to ISO19770-2:4.1.17). Examples of platforms include a computer system running Linux®, Microsoft Windows 10®, or Java™.

An NVMe-oF Platform should be a Platform that is capable of instantiating a NVMe-oF association as an NVMe host. Refer to NISTIR 7695 from ISO/IEC 19770-2.

Technical input submitted to the NVM Express® Workgroup is subject to the terms of the NVM Express® Participant's agreement. Copyright © 2008 to 2022 NVM Express, Inc.

1.3.18 Pre-OS Driver

The NVMe-oF environment that is responsible for initializing any host specific NVMe configuration and populating that configuration data into the NBFT.

1.3.19 Storage Subsystem

An NVM subsystem that may include a non-volatile memory storage medium and exposes namespaces, but is not a Discovery Service.

1.3.20 Subsystem and Namespace Descriptor (SSNS)

An NBFT structure that contains information about how to locate a specific namespace by its globally unique CNS 03h NID value.

1.3.21 Unified Extensible Firmware Interface (UEFI)

Firmware interface between the platform and the operating system. Defined by the UEFI Forum (refer to uefi.org).

1.3.22 UEFI Human Interface Infrastructure (UEFI HII)

UEFI core functionality describes the data and APIs used to manage the system's configuration and their human interfaces including basic mechanisms for managing both user input and UEFI protocols. Refer to the Unified Extensible Firmware Interface (UEFI) Specification.

1.3.23 Extended System Description Table (XSDT)

The XSDT is the first table used by an operating systems ACPI subsystem and contains addresses of most other ACPI tables on the system. If a system has an NBFT, its address will be in the system's XSDT. Refer to the Advanced Configuration and Power Interface Specification.

1.4 References

NVM Express® Base Specification, Revision 2.0c. Available from <https://www.nvmexpress.org>.

NVM Express® Command Set Specification, Revision 1.0b. Available from <https://www.nvmexpress.org>.

NVM Express® TCP Transport Specification, Revision 1.0b. Available from <https://www.nvmexpress.org>.

Win 32 API: NL_PREFIX_ORIGIN enumeration, Accessed April 19, 2022. Available from https://docs.microsoft.com/en-us/windows/win32/api/nldef/ne-nldef-nl_prefix_origin.

RFC 7143, M. Chadalapaka, J. Satran, K. Meth, D. Black, "Internet Small Computer System Interface (iSCSI) Protocol (Consolidated)", April 2014. Available from <https://datatracker.ietf.org/doc/html/rfc7143>.

RFC 3720, J. Satran, K. Meth, C. Sapuntzakis, M. Chadalapaka, E. Zeidner, "Internet Small Computer Systems Interface (iSCSI)", April 2004. Available from <https://datatracker.ietf.org/doc/html/rfc3720>.

RFC 4173, P. Sarkar, D. Missimer, C. Sapuntzakis, "Bootstrapping Clients using the Internet Small Computer System Interface (iSCSI) Protocol", April 2005. Available from <https://datatracker.ietf.org/doc/html/rfc4173>.

iSCSI Boot Firmware Table (iBFT) as Defined in ACPI 3.0b Specification, Version 1.03. Available from <https://download.microsoft.com/download/1/e/5/1e5a2287-366c-431a-8e05-b958540230b1/iBFT.docx>.

IEEE Guidelines for Use of Extended Unique Identifiers, August 2017. Available from <https://standards.ieee.org/wp-content/uploads/import/documents/tutorials/eui.pdf>.

Technical input submitted to the NVM Express® Workgroup is subject to the terms of the NVM Express® Participant's agreement. Copyright © 2008 to 2022 NVM Express, Inc.

IEEE 802.1q-2018: IEEE Standard for Local and Metropolitan Area Networks—Bridges and Bridged Networks, August 2018. Available from <https://standards.ieee.org/ieee/802.1Q/6844/>.

IETF RFC 2132, S. Alexander, R. Droms, “DHCP Options and BOOTP Vendor Extensions”, March 1997. Available from <https://rfc-editor.org/rfc/rfc2132.txt>.

RFC 2631, E. Rescorla, “Diffie-Hellman Key Agreement Method”, June 1999. Available from <https://www.ietf.org/rfc.html>.

R. Hinden, B. Carpenter, L. Masinter, “Format for Literal IPv6 Addresses in URL's”, December 1999. Available from <https://rfc-editor.org/rfc/rfc2732.txt>.

RFC 3986, T. Berners-Lee, R. Fielding, L. Masinter, “Uniform Resource Identifier (URI): Generic Syntax”, January 2005. Available from <https://www.ietf.org/rfc.html>.

IETF RFC 4038, M-K. Shin, Ed., Y-G. Hong, J. Hagino, P. Savola, E. M. Castro, “Application Aspects of IPv6 Transition”, March 2005. Available from <https://rfc-editor.org/rfc/rfc4038.txt>.

IETF RFC 4249, B. Haberman, “IP Forwarding Table MIB”, April 2006. Available from <https://www.ietf.org/rfc/rfc4292.txt>.

IETF RFC 4291, R. Hinden, S. Deering, “IP Version 6 Addressing Architecture”, February 2006. Available from <https://rfc-editor.org/rfc/rfc4291.txt>.

DMTF DSP0270, “Redfish Host Interface Specification”, Revision 1.3. Available from https://www.dmtf.org/sites/default/files/standards/documents/DSP0270_1.3.0.pdf.

DMTF DSP8010 “Redfish 2021.4 Schema Bundle”, Version 2021.4, December 2021. Available from https://www.dmtf.org/sites/default/files/standards/documents/DSP8010_2021.4.zip

Tianocore, “EDK2 Build Specification”, Version 1.0.0. Available from <https://edk2-docs.gitbook.io/edk-ii-build-specification/>.

ACPI “Advanced Configuration and Power Interface (ACPI) Specification”, Version 6.4. Available from <https://uefi.org/specs/ACPI/6.4/>.

UEFI, “Unified Extensible Firmware Interface (UEFI) Specification”, Version 2.9. Available from <https://uefi.org/specifications>.

IANA TLS Parameters Registry, TLS Cipher Suites, accessed April 26th, 2022. Available from <https://www.iana.org/assignments/tls-parameters/tls-parameters.xhtml#tls-parameters-4>.

IANA TLS Parameters Registry, TLS Supported Groups, accessed April 26th 2022. Available from <https://www.iana.org/assignments/tls-parameters/tls-parameters.xhtml#tls-parameters-8>.

NISTIR 7695 from ISO/IEC 19770-2. Available from <https://nvlpubs.nist.gov/nistpubs/Legacy/IR/nistir7695.pdf>

PCI Express Base Specification, Revision 6.0. Available from <https://members.pcisig.com/wg/PCI-SIG/document/16609>

PCI Firmware Specification, Revision 3.3. Available from <https://members.pcisig.com/wg/PCI-SIG/document/15350>

1.5 Conventions

1.5.1 Keywords

Several keywords are used to differentiate between different levels of requirements.

Technical input submitted to the NVM Express® Workgroup is subject to the terms of the NVM Express® Participant's agreement. Copyright © 2008 to 2022 NVM Express, Inc.

1.5.1.1 mandatory

A keyword indicating items to be implemented as defined by this specification.

1.5.1.2 may

A keyword that indicates flexibility of choice with no implied preference.

1.5.1.3 obsolete

A keyword indicating functionality that was defined in a previous version of the NVM Express specification and that has been removed from this specification.

1.5.1.4 optional

A keyword that describes features that are not required by this specification. However, if any optional feature defined by the specification is implemented, the feature shall be implemented in the way defined by the specification.

1.5.1.5 R

“R” is used as an abbreviation for “reserved” when the figure or table does not provide sufficient space for the full word “reserved”.

1.5.1.6 reserved

A keyword referring to bits, bytes, words, fields, and opcode values that are set-aside for future standardization. Their use and interpretation may be specified by future extensions to this or other specifications. A reserved bit, byte, word, field, property, or register shall be cleared to 0h, or in accordance with a future extension to this specification. The recipient of a command or a register write is not required to check reserved bits, bytes, words, or fields. Receipt of reserved coded values in defined fields in commands shall be reported as an error. Writing a reserved coded value into a controller property field produces undefined results.

1.5.1.7 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 the specification.

1.5.1.8 should

A keyword indicating flexibility of choice with a strongly preferred alternative. Equivalent to the phrase “it is recommended”.

1.5.2 Numerical Descriptions

A 0's based value is a numbering scheme in which the number 0h represents a value of 1h, 1h represents 2h, 2h represents 3h, etc. In this numbering scheme, there is no method to represent the value of 0h. Values in this specification are 1-based (i.e., the number 1h represents a value of 1h, 2h represents 2h, etc.) unless otherwise specified.

Size values are shown in binary units or decimal units. The symbols used to represent these values are as shown in Figure 2.

Figure 2: Decimal and Binary Units

Decimal		Binary	
Symbol	Power (base-10)	Symbol	Power (base-2)
kilo / k	10 ³	kibi / Ki	2 ¹⁰
mega / M	10 ⁶	mebi / Mi	2 ²⁰
giga / G	10 ⁹	gibi / Gi	2 ³⁰
tera / T	10 ¹²	tebi / Ti	2 ⁴⁰
peta / P	10 ¹⁵	pebi / Pi	2 ⁵⁰
exa / E	10 ¹⁸	exbi / Ei	2 ⁶⁰
zetta / Z	10 ²¹	zebi / Zi	2 ⁷⁰
yotta / Y	10 ²⁴	yobi / Yi	2 ⁸⁰

The ^ operator is used to denote the power to which that number, symbol, or expression is to be raised.

Some parameters are defined as an ASCII string. ASCII strings shall contain only code values 20h through 7Eh. For the string “Copyright”, the character “C” is the first byte, the character “o” is the second byte, etc. The string is left justified and shall be padded with spaces (ASCII character 20h) to the right if necessary. A hexadecimal ASCII string is an ASCII string that uses a subset of the code values: “0” to “9”, “A” to “F” uppercase, and “a” to “f” lowercase.

Hexadecimal (i.e., base 16) numbers are written with a lower case “h” suffix (e.g., 0FFFh, 80h). Hexadecimal numbers larger than eight digits are represented with an underscore character dividing each group of eight digits (e.g., 1E_DEADBEEFh).

Binary (i.e., base 2) numbers are written with a lower case “b” suffix (e.g., 1001b, 10b). Binary numbers larger than four digits are written with an underscore character dividing each group of four digits (e.g., 1000_0101_0010b).

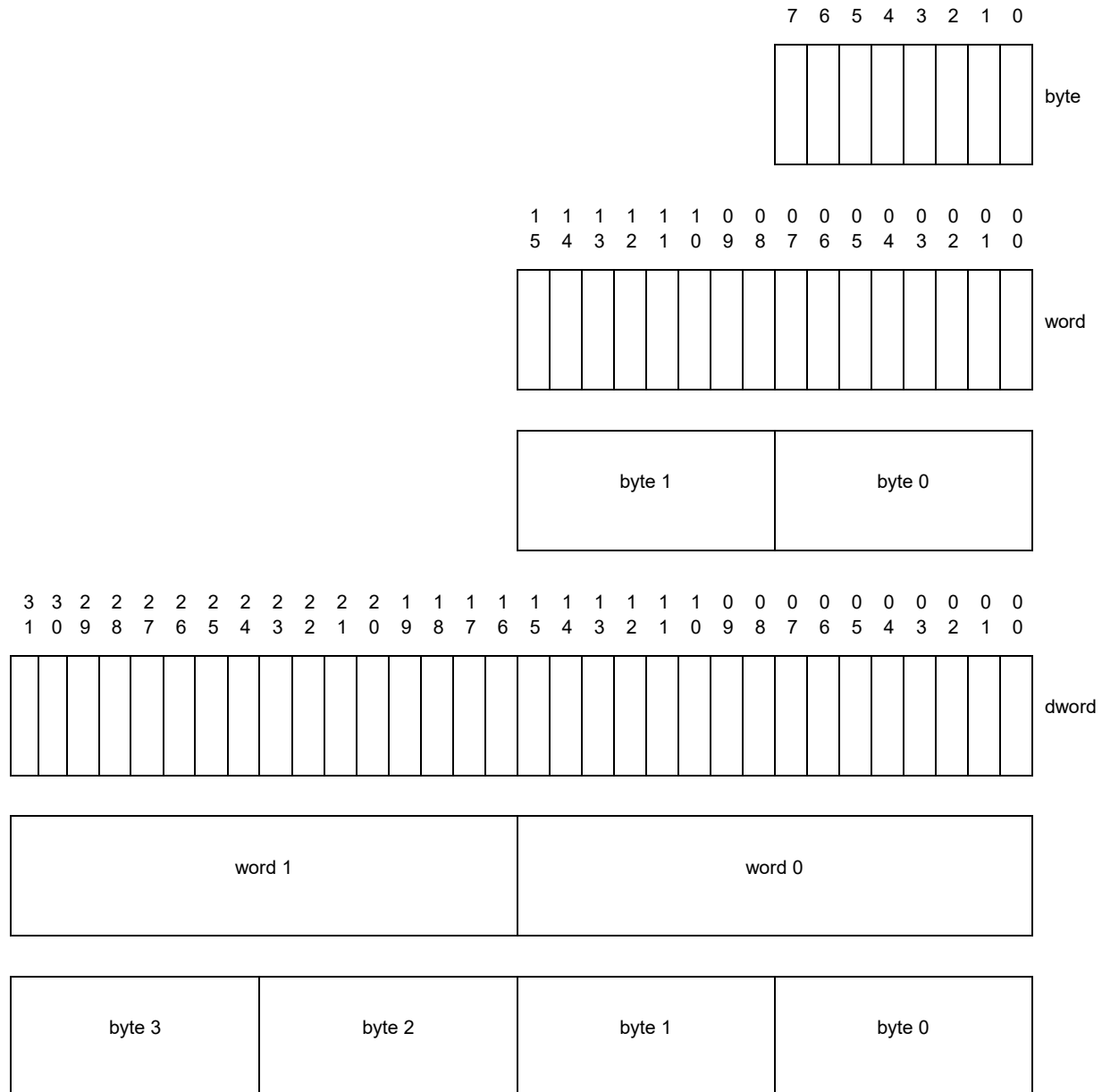
All other numbers are decimal (i.e., base 10). A decimal number is represented in this specification by any sequence of digits consisting of only the Western-Arabic numerals 0 to 9 not immediately followed by a lower-case b or a lower-case h (e.g., 175). This specification uses the following conventions for representing decimal numbers:

- the decimal separator (i.e., separating the integer and fractional portions of the number) is a period;
- the thousands separator (i.e., separating groups of three decimal digits in a portion of the number) is a comma;
- the thousands separator is used in only the integer portion of a number and not the fractional portion of a number; and
- the decimal representation for a year does not include a comma (e.g., 2019 instead of 2,019).

1.5.3 Byte, Word, and Dword Relationships

Figure 3 illustrates the relationship between bytes, words and dwords. A qword (quadruple word) is a unit of data that is four times the size of a word; it is not illustrated due to space constraints. Unless otherwise specified, this specification specifies data in a little-endian format.

Figure 3: Byte, Word, and Dword Relationships



1.5.4 Data Endian-ness

This specification conforms to the Conventions section, Keywords section, and Byte, Word, and Dword Relationships section of the NVM Express Base Specification.

1.5.5 IP Addresses and Formatting

Throughout this document, IP Addresses as identified in this document shall follow the common conventions specified in IETF RFC 4291, RFC 4038, RFC 3986, and RFC 2732. Excerpts and examples included here are done so for readability and do not supersede those RFC's unless explicitly stated.

1.5.5.1 IPv6 and IPv4 Byte Representations

Within this specification, when IP addresses are used in byte fields, then those fields shall consist of reserving 16 bytes to accommodate IPv6 addresses.

Additionally, IPv4 addresses shall be encapsulated in IPv6 address following IETF RFC 4291 and RFC 4038 where IPv4 is modeled inside the IPv6 Scope (which is included in a v4-to-v6 range) in "IPv4-mapped IPv6 address" that looks like:

Figure 4: IPv4 in IPv6 Format

MSB		LSB
80 bits	16 bits	32 bits
0000	FFFF	IPv4 Address

This field has the MSB starting at the left decreasing to right in the data.

1.5.5.2 String Representation

When IP literals are used in string arrays, those IP literals shall follow the IETF RFC 3986 and RFC 2732 grammar. Specifically, IPv6 Addresses shall follow IETF RFC 2732 square bracket encapsulation of IPv6 String Literals with port and prefix outside square bracket encapsulated literals.

Examples of such are:

"[2001:0DB8:7654:3210:FEDC:BA98:7654:3210]:4420"

"[::FFFF:129.144.52.38]:8009"

Within string arrays, IPv4 Addresses may be encapsulated or not. Field requirements will state if a URI is expected to be a string array.

1.5.6 Host Name String Format

Throughout this document, when a well-formed string representation of a Host Name is required, the string shall follow the following common host representation as specified in RFC 3986, including UTF-8 formatting of the Host Name String.

Host name if supplied shall be in one of the following formats:

- <HOST NAME>
- <HOST NAME>.
- <HOST NAME>.<DOMAIN>

1.5.7 NVMe-oF URI Format

Throughout this document, when well-formed string representations of an address URI are required to be specified, the following URI format shall be adhered to.

An NVMe-oF scheme, that conforms to IETF RFC 3986, consists of two parts: the nvme class, and a +<trtype> transport specific identifier. A valid URI shall be formatted in UTF-8.

Valid Transport types string representation are found in Figure 7.

When referencing IP based transports, Host and Port of an Authority statement shall be valid IP literals, decimal formatted addresses, or resolvable names. Allocation, and registration of which are outside the scope of this document.

Examples of well-formed string representations of an address URI using a tcp transport specific identifier include:

```
"nvme+tcp://192.168.1.1:4420/"
```

```
"nvme+tcp://[FE80::1010]:4420/"
```

1.5.8 Optional and Mandatory Fields in Tables

Tables and Structures within this document may contain a 'O/M' column with entries in that column that indicate if the use or content of specific fields in that Table or Structure is Optional or Mandatory. If a record includes an 'M' then that field is mandatory. If a record includes a 'O', it indicates the field is optional.

Optional tables may be absent unless specified otherwise. Optional fields may be cleared to 0h unless otherwise specified. Optional referenced objects, such as heap objects, may be absent and their reference fields may be cleared to 0h unless otherwise specified.

1.5.9 Namespace Identifier String Formats

This section specifies representations that contain lexical decoration and tokenization that may be used when well-known textual representations of Namespace Identifiers are necessary for interoperability. In all cases, the contents of internal and in-memory formats are defined elsewhere, and this section only relates to their textual representations associated with the NVM Express Boot Specification. Further example uses of the Namespace Identifier string formats are described in Section A.3.2.3.

1.5.9.1 NVMe EUI-64 String Format

The EUI-64 textual representation shall follow the IEEE RA EUI-64 hyphenated textual format (refer to the IEEE Guidelines for Use of Extended Unique Identifiers) with the generic prefix decorator "eui:".

This format consists of:

- The string "eui:"; and
- a string representation of the 8 byte IEEE RA formatted EUI-64.

The following is an example of such an NVMe EUI-64 formatted string: "eui:AC-DE-48-23-45-67-01-9F".

1.5.9.2 NVMe NGUID String Format

The NGUID textual representation shall be a prefix of "nvme-nguid:" followed by the NGUID tokens: vendor specific extension identifier, IEEE OUI, and an extension identifier, with a ASCII '-' as a field separator between tokens.

This format consists of:

- The string "nvme-nguid:";
- a string representation of the 8 byte vendor specific extension identifier;
- a ASCII "-";
- a string representation of the 3 byte IEEE OUI;
- a ASCII "-"; and
- a string representation of the 5 byte extension identifier.

The following is an example of such an NVMe NGUID formatted string: “nvme-nguid:FEDCBA9876543210-ABCDEF-0123456789”.

1.5.9.3 NVMe UUID String Format

The UUID textual representation shall follow the IETF RFC 4122 textual format.

The following is an example of such an NVMe UUID formatted string:

“urn:uuid:4eff7f8e-d353-4e9b-a4ec-deea8eab84d7”.

2 Background and General Methodology

2.1 Booting from NVMe

To address the complexities of boot from NVMe, especially NVMe-oF, this document contains normative content that describes generic requirements of information sharing as well as introducing domain specific mechanics such as ACPI tables to provide reliable ways for pre-OS and OS drivers to share configuration context.

2.1.1 Boot from NVMe: PCIe Transport

For the PCIe-based controllers, NVMe storage behaves very similarly to other types of block storage as indicated in Annex A.2.7. Expansion ROM images or boot environment firmware initializes and configures the PCI functions which are the controllers. The controlling firmware enumerates namespaces as block storage devices to the boot environment. Similar to SCSI-based storage, the controlling firmware converts the basic block storage device interface (select/read/write/deselect) used by the boot environment to the NVM Express protocol commands and PCIe transport interface for use with the namespaces. When the OS environment assumes control, its own drivers detects, initializes, and configures the NVMe controllers and enumerate the NVM subsystem's namespaces with the OS.

As the NVM subsystem configuration is predominantly static (i.e., changes only with explicit administrative action taken on the NVM subsystem), the storage devices enumerated to the boot environment will be the same set of storage devices enumerated to the OS. No interaction between the boot environment and OS environment is necessary to ensure device presence.

2.1.2 Boot from NVMe-oF: FC Transport

NVMe-oF Boot from Fibre Channel (FC) storage currently leverages the implementation for Boot from SAN on Fibre Channel. Rather than speaking the SCSI protocol to FC devices, the FC host adapters utilize the NVMe-oF and FC-NVMe protocols with the FC devices. Expansion ROM images for the FC host adapter detects NVMe-oF Subsystems, both discovery and storage subsystems, and creates NVMe-oF associations (including an NVMe controller) to NVM subsystems. Namespaces on the NVM subsystems will then be enumerated to the boot environment.

Note: As a FC endpoint may support both the SCSI and NVMe protocols, an FC host adapter may enumerate storage devices that are both SCSI and NVMe with the boot environment and OS environment.

When creating NVMe-oF FC associations, the FC host adapter typically creates a Host NQN and a HostID value to be used for the associations while in the boot environment. Expansion ROM images typically present administrative interfaces to view and/or set the values to be used for the boot environment. Values for OS-based associations are typically assigned by the OS.

FC supports connectivity to a large number of FC-NVMe devices and NVM subsystems, and the NVM subsystems may enumerate a large number of storage devices. As such, FC host adapters have leveraged the same management methodologies developed for SCSI to allow administrators to view and select the NVMe-oF-based storage devices that are to be enumerated to the boot environment.

FC storage is self-discovering. The host adapter typically sees the same devices regardless of whether it is the boot environment or the OS environment. The host adapter expansion ROM may filter the devices presented to the boot environment. However the devices presented to the OS environment are typically all of the devices, which includes those presented to the boot environment. Therefore, device presentation for boot devices occurs in both the boot environment and the OS environment without requiring the two environments to interact.

2.1.3 Boot from NVMe-oF: TCP Transport

Boot from NVMe-oF storage on TCP, behaves similarly to NVMe PCIe-based storage and NVMe-oF FC-based storage. Once an NVMe-oF association is created and a controller exists, expansion ROM software or boot environment software enumerates namespaces on the NVMe subsystem and convert the basic block storage device interface (select/read/write/deselect) used by the boot environment to the NVM Express protocol commands.

Similar to historical requirements in iSCSI storage (See further detail in Annex A.2.4), NVMe-oF TCP associations are not self-discovering. The host adapter, or NVMe-oF stack utilizing a network adapter, requires information to locate and/or create the NVMe-oF TCP associations. Information includes IP addresses and port numbers for the initial TCP connections, Host NQN and Host ID values, and Subsystem NQN values to be connected to. To aid in the amount of information required to be known and stored, IP services such as DHCP may be used to query for the NVMe-oF TCP boot device session information. Even in those cases, address information and service type information may have to be specified.

To address these issues, the NVMe Boot Firmware Table (NBFT) provides a mechanism to pass NVMe-oF Boot Association information from the boot environment to the OS environment using the ACPI API. The NBFT is defined in such a manner that transports other than TCP, as supported by the NVM Express Transport specifications, may create tables for boot associations such that they are passed to the OS.

2.1.3.1 NVMe Boot Firmware Table

Section 3.1 defines the NVMe Boot Firmware Table (NBFT) which is a multi-level table structure that provides a standardized means of passing configuration and connection context from a pre-OS Boot environment to an OS runtime in regard to the administratively configured run-time as specified by its user. It is formatted to follow conventions of, and to be included as, an ACPI XSDT table. The naming follows the iSCSI Boot Firmware Table (iBFT).

3 Boot Mechanisms

This section specifies different mechanisms and their requirements that may be used for booting over NVMe-oF.

3.1 ACPI NVMe Boot Firmware Table (NBFT)

The NVMe Boot Firmware Table (NBFT) structure is an ACPI structure in host read/write memory that the platform boot firmware reserves for ACPI usage. The NBFT shall be referenced in the ACPI XSDT. This structure is mandatory if the pre-OS environment supports the NVMe Boot Transport mechanism and NBFT. The presence of the table Header, Control, and Host Descriptor structures shall indicate to the OS that the system supports booting from an NVMe boot mechanism. The NBFT uses a structure approach for associating the local and remote aspects of NVMe Boot configurations. If a BIOS supports NVMe Boot with the mechanism defined in this specification, but it is not administratively configured, then a bare table without a valid control structure shall be present.

ACPI NBFT content is a one-way means of communication from the pre-OS environment to the runtime OS environment. Changes that an OS may make to a NBFT in memory are not able to be captured by the pre-boot environment. If changes to the pre-boot configuration are required, another mechanism such as DMTF's Redfish may be required, and its use and configuration are outside the scope of this document.

Note that some systems may provide more than one NBFT. If there are multiple NBFTs on a system, each of them shall be pointed to by the XSDT. Individual NBFTs are independent of each other (i.e., no cross references are possible between any separate NBFT tables on the same system).

The topic of multiple NBFTs is discussed further in section A.3.1.3 of this document.

3.1.1 NBFT Table Overview

The layout of a NBFT is divided up into representative descriptors, some mandatory and some optional that fully describe features as described in Figure 5.

The table Header, Control, and Host Descriptor structures have fixed well-known addresses, all other table structures have starting locations and lengths based on their content defined in the Control structure. Additionally, the Header, Control, and Host Descriptors are the only mandatory structures for a valid minimal NBFT, all other structures are optional unless pre-OS connectivity is established that requires a structure to be conveyed to an OS. Lastly, the Header, Control and Host Descriptors shall be present in the order defined in Figure 8, however all other structures are referenced by offset in the Control structure and may be contiguously added to the NBFT in any order.

Figure 5: NBFT Elements

Element	Descriptor Type Designator	Description	Reference
Header	N/A	An ACPI structure header with some additional NBFT specific info.	3.1.2.1
Control Descriptor	1h	Indicates the location of host, HFI, SSNS, security, and discovery descriptors.	3.1.2.2
Host Descriptor	2h	Host information.	3.1.2.3
HFI Descriptor	3h	An indexable table of HFI Descriptors, one for each fabric interface on the host.	3.1.2.4
Subsystem Namespace Descriptor	4h	An indexable table of SSNS Descriptors.	3.1.2.5
Security Descriptor	5h	An indexable table of Security descriptors.	3.1.2.6

Technical input submitted to the NVM Express® Workgroup is subject to the terms of the NVM Express® Participant's agreement. Copyright © 2008 to 2022 NVM Express, Inc.

Figure 5: NBFT Elements

Element	Descriptor Type Designator	Description	Reference
Discovery Descriptor	6h	An indexable table of Discovery Ddescriptors.	3.1.2.7
Heap Resident Structures			
HFI Transport Descriptor	7h	Indicated by an HFI Descriptor, corresponds to a specific transport for a single HFI..	3.1.2.4.1.1
Reserved	8h		
SSNS Extended Info Descriptor	9h	Indicated by an SSNS Descriptor if required.	3.1.2.5.5

Figure 6 shows graphically the relationships of the elements of the NBFT. In addition to the two elements shown in the Heap, other elements are placed in the Heap as described in subsequent sections.

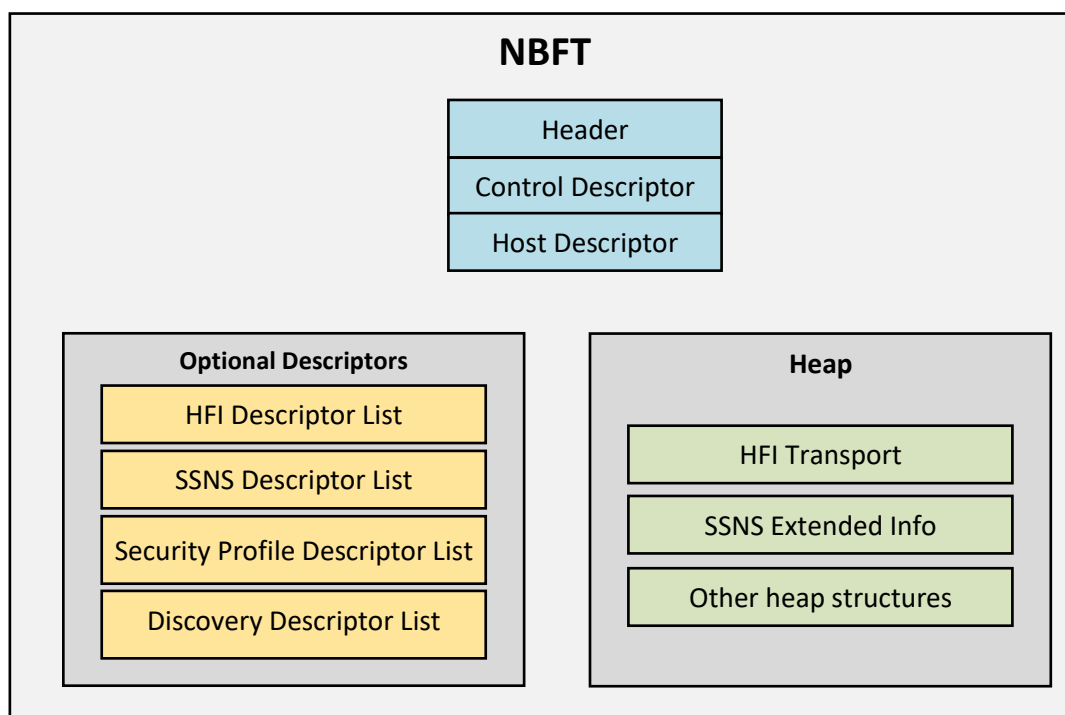
Figure 6: NBFT Layout

Table elements are typically compacted in memory, but this is not required, and the tree shall be traversed by the links (i.e., offset fields).

3.1.1.1 Heap

The Heap is a structure within the NBFT and exists as a contiguous range of addressable variable storage space used by other NBFT structures. A heap contains objects that may be binary data, a null-terminated string, or structured data that may be variable length in size. The null termination of a string shall count towards the length of referenced heap object. The heap may or may not be compact. If a pre-OS NBFT driver removes an entry, the pre-OS NBFT driver may or may not recompact the heap. Order of entries in the heap shall not be assumed. Offsets to heap objects in NBFT descriptors shall be specified as a relative offset from the start of the NBFT. When addressing a valid object in the heap, length and offset address

shall be set to a non-zero value. For unused and optional heap objects, length and offset address shall both be cleared to 0h in NBFT descriptors.

3.1.1.2 Index Values

List elements in NBFT descriptors are specified as Index values. Index values are used as part of the structure offsets for providing the enumerated elements in the system with discrete configuration. As part of resolving the unique access model of a device, HFI, SSNS, and their security profiles shall be paired together. Index values unless otherwise stated, shall start from 1. Index values for NBFT descriptors are not used in any specific order, and index order should not denote priority to consumers of the data.

3.1.1.3 Transport Types

The NVM Express Base Specification defines Numeric Designator values for Transport types. Throughout this document these Numeric Designators are used for transport specific structures. Transport types supported by the NBFT are listed in Figure 7. In addition, for fields that use string data such as URIs, a string representation is also provided in Figure 7.

Figure 7: NBFT Interface Transport Types

Numeric Designator	Transport	String Designator	Default Transport Service Identifier
3h	NVMe/TCP (802.3 + TCP/IP)	"tcp"	34343230h (i.e., "4420") with a length of 4 Bytes

3.1.1.4 Driver Signature

The Driver Signature is a NBFT Header structure that shall contain a per-system unique identification of the driver source of the NBFT structure data. In the NBFT Table Header, the Driver Signature is a heap object consisting of a UTF-8 string for the unique identity of the driver source of NBFT structure data in that system. On UEFI-based systems, the Driver Signature shall consist of UEFI Driver Image Device Path string. Examples of this are found in Annex A.3.1.3.4.

3.1.2 NBFT Structure

3.1.2.1 NBFT Header Structure

Structure of the NBFT Header Table (refer to Figure 8) follows the normal conventions of an ACPI system description table.

The host or OS software that is compliant with this specification should ignore any new identifiers and fields indicating larger minor revision numbers.

Figure 8: NBFT Table

Bytes	Description
Header	
03:00	Signature: An ASCII string representation of the table identifier. This field shall be set to the value 4E424654h (i.e., "NBFT") for the NVMe-oF Boot Firmware Table.
07:04	Length: The length of the table, in bytes, including the header, starting from offset 0h. This field is used to record the size of the entire table.
08	Major Revision: The major revision of the structure corresponding to the Signature field. Larger major revision numbers should not be assumed backward compatible to lower major revision numbers with the same signature. This field shall be set to '1'.

Figure 8: NBFT Table

Bytes	Description						
09	Checksum: The entire table, including the Checksum field, shall sum to 0h to be considered valid. This checksum shall follow ACPI table checksum requirements. Refer to the Advanced Configuration and Power Interface (ACPI) Specification.						
15:10	OEMID: OEMID shall be populated by the NBFT driver writer by an OEM-supplied string that identifies the OEM. All trailing bytes shall be null.						
23:16	OEM Table ID: This field shall be populated by the NBFT driver writer with an OEM-supplied string that the OEM uses to identify the particular data table. This field is particularly useful when defining a definition block to distinguish definition block functions. The OEM assigns each dissimilar table a new OEM Table ID.						
27:24	OEM Revision: An OEM-supplied revision number. Larger numbers are assumed to be newer revisions.						
31:28	Creator ID: Vendor ID of utility that created the table. For instance, this may be the ID for the ASL Compiler.						
35:32	Creator Revision: Revision of utility that created the table. For instance, this may be the ID for the ASL Compiler.						
39:36	Heap Offset (HO): This field indicates the offset in bytes of the heap, if any, from byte offset 0h of the NBFT Table Header.						
43:40	Heap Length (HL): The length of the heap, if any.						
49:44	<p>Driver Signature Heap Object Reference: This field indicates the offset in bytes of a heap object containing the Driver Signature, if any, from byte offset 0h of the NBFT Table Header. For more detail on Driver Signature, refer to section A.3.1.3.4.</p> <table> <tr> <th>Bytes</th><th>Description</th></tr> <tr> <td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr> <tr> <td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr> </table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description						
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.						
05:04	Length: Length in bytes of the heap object, if any.						
50	<p>Minor Revision: The minor revision of the structure corresponding to the Signature field. If the major revision numbers are the same, any minor revision number differences shall be backwards compatible with the same signature.</p> <p>This field shall be cleared to '0'.</p>						
63:51	Reserved						
Control Descriptor							
64	Structure ID: This field specifies the element (refer to Figure 5). This field shall be set to 1h (i.e., Control).						
65	<p>Major Revision: The major revision of the structure corresponding to the Signature field. Larger major revision numbers should not be assumed backward compatible to lower major revision numbers with the same signature.</p> <p>This field shall be set to '1'.</p>						
66	<p>Minor Revision: The minor revision of the structure corresponding to the signature field. If the major revision numbers are the same, any minor revision number differences shall be backwards compatible with the same signature.</p> <p>This field shall be cleared to '0'.</p>						
67	Reserved						
69:68	Control Structure Length (CSL): This field indicates the length in bytes of the Control Descriptor.						
70	<p>Flags:</p> <table> <tr> <th>Bits</th><th>Feature</th></tr> <tr> <td>07:01</td><td>Reserved</td></tr> <tr> <td>00</td><td>Block Valid: If set to '1', then this structure is valid. If cleared to '0', then this structure is not valid.</td></tr> </table>	Bits	Feature	07:01	Reserved	00	Block Valid: If set to '1', then this structure is valid. If cleared to '0', then this structure is not valid.
Bits	Feature						
07:01	Reserved						
00	Block Valid: If set to '1', then this structure is valid. If cleared to '0', then this structure is not valid.						
71	Reserved						

Figure 8: NBFT Table

Bytes	Description						
77:72	Host Descriptor (HDESC): This field indicates the location and length of the Host Descriptor (refer to Figure 9).						
	<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the Host Descriptor from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the Host Descriptor. The Length field shall be set to a non-zero value.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the Host Descriptor from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.	05:04	Length: Length in bytes of the Host Descriptor. The Length field shall be set to a non-zero value.
	Bytes	Description					
03:00	Offset: Offset in bytes of the Host Descriptor from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.						
05:04	Length: Length in bytes of the Host Descriptor. The Length field shall be set to a non-zero value.						
78	Host Descriptor Version (HSV): This field indicates the version of the Host Descriptor. This field shall be set to '1'.						
79	Reserved						
Host Fabric Interface Descriptor Information							
83:80	HFI Descriptor List Offset (HFIO): If this field is set to a non-zero value, then this field indicates the offset in bytes of the HFI Descriptor List (refer to Figure 10), if any, from byte offset 0h of the NBFT Table Header. If the NumHFI field is cleared to 0h, then this field is reserved.						
85:84	HFI Descriptor Length (HFIL): This field indicates the length in bytes of each HFI Descriptor, if any. If the NumHFI field is cleared to 0h, then this field is reserved.						
86	HFI Descriptor Version (HFIV): This field indicates the version of each HFI Descriptor. This field shall be set to '1'.						
87	Number of Host Fabric Interface Descriptors (NumHFI): This field indicates the number of HFI Descriptors (refer to Figure 11) in the HFI Descriptor List (refer to Figure 10), if any. If no interfaces have been configured, then this field shall be cleared to 0h.						
Subsystem Namespace Descriptor Information							
91:88	SSNS Descriptor List Offset (SSNSO): This field indicates the offset in bytes of the SSNS Descriptor List (refer to Figure 14), if any, from byte offset 0h of the NBFT Table Header. If the NumSSNS field is cleared to 0h, then this field is reserved.						
93:92	SSNS Descriptor Length (SSNSL): This field indicates the length in bytes of each SSNS Descriptor, if any. If the NumSSNS field is cleared to 0h, then this field is reserved.						
94	SSNS Descriptor Version (SSNSV): This field indicates the version of the SSNS Descriptor. This field shall be set to '1'.						
95	Number of Subsystem and Namespace Descriptors (NumSSNS): This field indicates the number of Subsystem Namespace (SSNS) Descriptors (refer to Figure 15) in the SSNS Descriptor List (refer to Figure 14), if any.						
Security Profile Information							
99:96	Security Profile Descriptor List Offset (SECO): This field indicates the offset in bytes of the Security Profile Descriptor List (refer to Figure 20), if any, from byte offset 0h of the NBFT Table Header. If the NumSec field is cleared to 0h, then this field is reserved.						
101:100	Security Profile Descriptor Length (SECL): This field indicates the length in bytes of each Security Profile Descriptor, if any. If the NumSec field is cleared to 0h, then this field is reserved.						
102	Security Profile Descriptor Version (SECV): This field indicates the version of the Security Profile Descriptor. This field shall be set to '1'.						
103	Number of Security Profile Descriptors (NumSec): This field indicates the number of Security Profile Descriptors (refer to Figure 21), if any, in the Security Profile Descriptor List (refer to Figure 20).						
Discovery Descriptor Information							
107:104	Discovery Descriptor Offset (DISCO): This field indicates the offset in bytes of the Discovery Descriptor List (refer to Figure 23), if any, from byte offset 0h of the NBFT Table Header. If the NumDisc field is cleared to 0h, then this field is reserved.						
109:108	Discovery Descriptor Length (DISCL): This field indicates the length in bytes of each Discovery Descriptor, if any. If the NumDisc field is cleared to 0h, then this field is reserved.						
110	Discovery Descriptor Version (DISCV): This field indicates the version of the Discovery Descriptor. This field shall be set to '1'.						
111	Number of Discovery Descriptors (NumDisc): This field indicates the number of Discovery Descriptors (Figure 24), if any, in the Discovery Descriptor List (refer to Figure 23), if any.						
127:112	Reserved						

3.1.2.1.1 Optional: HFI Descriptor Offsets

If the Number of HFI Descriptors (NumHFI) field is greater than 1h, the offset for each additional HFI Descriptor is calculated from the descriptor information. That calculation is based on the two-byte location value offset by the number of HFI's specified in the NumHFI field.

An example calculation for the offset of a specific HFI Descriptor Index is:

$$\text{IF(NumHFI >1) } \{ \text{HFIO} + (\text{HFIL} * \text{NumHFI}) \}.$$

3.1.2.1.2 Optional: Subsystem Namespace Descriptor Offsets

If the number of Subsystem and Namespace Descriptors (NumSSNS) field specified in the NumSSNS field, is greater than 1h, then additional Namespace records are present and are located at an offset after the first Namespace record. The offset calculation is based on the two-byte location value of Namespaces specified in NumSSNS.

An example calculation for the Offset is:

$$\text{IF(NumSSNS >1) } \{ \text{SSNS_1_Offset} + (\text{SSNSL} * \text{NumSSNS}) \}.$$

3.1.2.1.3 Optional: Security Profile Descriptor Offsets

If the Number of Security Descriptors (NumSec) field referenced in the NumSec field, is greater than 0h, then the user has specified security profiles and this optional structure shall be populated. The offset for the number of security profiles is a calculated value based on the number of Security Profiles specified in the NumSec field.

An example calculation for the Offset is:

$$\text{IF(NumSec > 1) } \{ \text{SECO} + (\text{SECL} * (\text{NumSec} - 1)) \}.$$

3.1.2.1.4 Optional: Discovery Descriptor Offsets

If the Number of Discovery Descriptors (NumDisc) field referenced in the NumDisc field, is greater than 0h, then the user has specified discovery entries and the corresponding optional structure shall be populated. The offset for the number of Discovery Descriptor entries is a calculated value based on the number of Discovery Entries specified in the NumDisc field.

An example calculation for the Offset is:

$$\text{IF(NumDisc > 1) } \{ \text{DISCO} + (\text{DISCL} * (\text{NumDisc} - 1)) \}.$$

3.1.2.2 Control Descriptor

The Control Descriptor (refer to Figure 8) presents global configuration options that are meaningful to the system as a whole. This includes offset reference designators for the configured system elements. Byte Offset in the Control Descriptor is specified as a relative offset from the start of the NBFT.

The Control Descriptor provides lookup features into the offsets of each additional Table structure via structure indices. The Host Descriptor defines the identity of the host known by the pre-OS environment. Lastly, the heap space, an area for string and configuration attributes of a non-fixed length resides at the end of the NBFT and is defined in the NBFT Table Header.

3.1.2.3 Host Descriptor

The Host Descriptor (refer to Figure 9) specifically defines details that are scoped globally to the host. Only one instance of this descriptor shall exist per NBFT for a physical or logical system that implements boot

from NVMe-oF. The byte offset in the Host Descriptor Table is specified as a relative offset from the start of this Descriptor rather in reference from the start of the NBFT.

Figure 9: Host Descriptor

Bytes	Description									
00	Structure ID: This field shall be set to 2h (i.e., Host Descriptor; refer to Figure 5).									
01	Host Flags:									
	<table><tr><th>Bits</th><th>Description</th></tr><tr><td>00</td><td>Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.</td></tr><tr><td>01</td><td>HostID Configured: If set to '1', then the Host ID field contains an administratively-configured value. If cleared to 0, then the Host ID field contains a driver default value.</td></tr><tr><td>02</td><td>Host NQN Configured: If set to '1', then the Host NQN indicated by the Host NQN Heap Object Reference field contains an administratively-configured value. If cleared to 0, then the Host NQN indicated by the Host NQN Offset field contains a driver default value.</td></tr></table>	Bits	Description	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.	01	HostID Configured: If set to '1', then the Host ID field contains an administratively-configured value. If cleared to 0, then the Host ID field contains a driver default value.	02	Host NQN Configured: If set to '1', then the Host NQN indicated by the Host NQN Heap Object Reference field contains an administratively-configured value. If cleared to 0, then the Host NQN indicated by the Host NQN Offset field contains a driver default value.	
	Bits	Description								
	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.								
	01	HostID Configured: If set to '1', then the Host ID field contains an administratively-configured value. If cleared to 0, then the Host ID field contains a driver default value.								
	02	Host NQN Configured: If set to '1', then the Host NQN indicated by the Host NQN Heap Object Reference field contains an administratively-configured value. If cleared to 0, then the Host NQN indicated by the Host NQN Offset field contains a driver default value.								
	Primary Administrative Host Descriptor: This field indicates whether the Host Descriptor in this NBFT was selected as the primary NBFT for administrative purposes of platform identity as a hint to the OS.									
	If multiple NBFT tables are present, only one NBFT should be administratively selected. There is no enforcement mechanism for this to be coordinated between multiple NBFT tables, but this field should be set to Selected (i.e., 10b) if more than one NBFT is present.									
	04:03									
	<table><tr><th>Value</th><th>Description</th></tr><tr><td>00b</td><td>Not Indicated by Driver: The driver that created this NBFT provided no administrative priority hint for this NBFT.</td></tr><tr><td>01b</td><td>Unselected: The driver that created this NBFT explicitly indicated that this NBFT should not be prioritized over any other NBFT.</td></tr><tr><td>10b</td><td>Selected: The driver that created this NBFT explicitly indicated that this NBFT should be prioritized over any other NBFT.</td></tr><tr><td>11b</td><td>Reserved</td></tr></table>	Value	Description	00b	Not Indicated by Driver: The driver that created this NBFT provided no administrative priority hint for this NBFT.	01b	Unselected: The driver that created this NBFT explicitly indicated that this NBFT should not be prioritized over any other NBFT.	10b	Selected: The driver that created this NBFT explicitly indicated that this NBFT should be prioritized over any other NBFT.	11b
Value	Description									
00b	Not Indicated by Driver: The driver that created this NBFT provided no administrative priority hint for this NBFT.									
01b	Unselected: The driver that created this NBFT explicitly indicated that this NBFT should not be prioritized over any other NBFT.									
10b	Selected: The driver that created this NBFT explicitly indicated that this NBFT should be prioritized over any other NBFT.									
11b	Reserved									
05:07										
Reserved										
17:02	Host ID: This field shall be set to the Host Identifier (refer to the NVM Express Base Specification). This field shall not be empty if the NBFT and NVMe Boot are supported by the Platform.									
23:18	Host NQN Heap Object Reference: this field indicates a heap object containing a Host NQN (refer to the NVM Express Base Specification) This object shall not be empty if the NBFT and NVMe Boot are supported by the Platform.									
	<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. This field shall be set to a non-zero value.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any. This field shall be set to a non-zero value.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. This field shall be set to a non-zero value.	05:04	Length: Length in bytes of the heap object, if any. This field shall be set to a non-zero value.			
	Bytes	Description								
	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. This field shall be set to a non-zero value.								
05:04	Length: Length in bytes of the heap object, if any. This field shall be set to a non-zero value.									
31:24	Reserved									

The value of the Host ID field in the NBFT Host Descriptor shall be set to the value used by the pre-OS driver that populates the NBFT and was used as part of any Connect command by the pre-OS driver. Upon reaching run-time, an OS or application may use a different value for the Host ID, the value in the Host ID field is to represent the Host ID used as part of the boot lifecycle.

If the driver value used was administratively configured, the HostID Configured bit shall be set to '1'. If the driver assumed a value or used a non-administrative authoritative source (e.g., default values with respect to a platform or device), then the HostID Configured bit shall be cleared to '0'. If no authoritative sources are present to supply a value for the HostID field, the HostID field shall be cleared to 0h by the driver.

The Host NQN is used by the Connect command to identify the host in this transaction. Upon reaching run-time, an OS or application may use a different Host NQN than the Host NQN which was used as part of the Boot lifecycle. If the value of Host NQN was not administratively configured, the Host NQN Configured bit shall be cleared to '0' and the value of Host NQN in the Host Descriptor shall be populated with the value used by the driver.

3.1.2.4 Host Fabric Interface (HFI) Descriptor

This Host Fabric Interface Descriptor List (refer to Figure 10) consists of an array of contiguous Host Fabric Interface Descriptors (refer to Figure 11). Each Host Fabric Interface Descriptor indicates the transport specific parameters associated with one fabric interface. Fabric interfaces may be a representation of specific host adapters with specific transports as defined by the HFI Transport Type flag. The same fabric interface or host adapter may support multiple transport types with their own unique configuration, which results in multiple HFI Descriptors in sequence within the NumHFI defined in the Control Descriptor.

Transport specific details are addressed in section 3.1.2.4.1.

Figure 10: Host Fabric Interface (HFI) Descriptor List

Bytes	Description
31:00	Host Fabric Interface Descriptor 1: The first Host Fabric Interface Descriptor (refer to Figure 11), if any.
64:32	Host Fabric Interface Descriptor 2: The second Host Fabric Interface Descriptor (refer to Figure 11), if any.
...	
(n*32)-1:(n-1)*32	Host Fabric Interface Descriptor n: The last Host Fabric Interface Descriptor (refer to Figure 11), if any.

Figure 11: Host Fabric Interface (HFI) Descriptor

Bytes	Description						
00	Structure ID: This field shall be set to 3h (i.e., Host Fabric Interface Descriptor; refer to Figure 5).						
01	HFI Descriptor Index: This field indicates the number of this HFI Descriptor in the Host Fabric Interface Descriptor List (i.e., the number 1h represents HFI Descriptor 1, 2h represents HFI Descriptor 2, etc.).						
02	HFI Descriptor Flags: <table border="1"> <tr> <th>Bits</th><th>Description</th></tr> <tr> <td>07:01</td><td>Reserved</td></tr> <tr> <td>00</td><td>Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.</td></tr> </table>	Bits	Description	07:01	Reserved	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.
Bits	Description						
07:01	Reserved						
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.						
03	HFI Transport Type: Refer to Figure 7.						
15:04	Reserved						
21:16	HFI Transport Info Descriptor Heap Object Reference: If this field is set to a non-zero value, then: <ul style="list-style-type: none"> a) this field indicates the location and size of a heap object containing a HFI Transport Info Descriptor (refer to section 3.1.2.4.1); and b) the Offset field and the Length field shall be set to non-zero values. If this field is cleared to 0h, then this field is reserved. <table border="1"> <tr> <th>Bytes</th><th>Description</th></tr> <tr> <td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr> <tr> <td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr> </table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description						
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.						
05:04	Length: Length in bytes of the heap object, if any.						
31:22	Reserved						

3.1.2.4.1 HFI Transport Info Descriptor

This section contains all the HFI Transports that are supported. All HFI Transport Info Descriptors consist of a common header format with NBFT structure ID of 7h with a HFI Transport Type specific payload that corresponds to a Transport type described in the HFI Transport Type field (refer to Figure 7). Figure 12 shows the common representation of what a Transport Specific HFI Transport Info Descriptor header shall contain.

Figure 12: HFI Transport Info Descriptor Header Template

Bytes	Description
00	Structure ID: This field shall be set to 7h (i.e., HFI Transport Info; refer to Figure 5).
01	Version: This field shall be set to 1h.
02	HFI Transport Type: Refer to Figure 7.
03	Transport Info Version: Implementations compliant to this specification shall set this field to 1h.
05:04	HFI Descriptor Index: The value of the HFI Descriptor Index field of the HFI Descriptor (refer to Figure 11) whose HFI Transport Info Descriptor Heap Object Reference field indicates this HFI Transport Info Descriptor.

3.1.2.4.1.1 HFI Transport Info Descriptor – NVMe/TCP

This section defines the HFI Transport details specific to the NVMe/TCP transport defined in The NVM Express TCP Transport Specification in Figure 13.

Figure 13: HFI Transport Info Descriptor – NVMe/TCP

Bytes	O/M ¹	Description													
00	M	Structure ID: This field shall be set to 7h (i.e., HFI Transport Info; refer to Figure 5).													
01	M	Version: This field shall be set to 1h.													
02	M	HFI Transport Type: This field shall be set to 03h (i.e., NVMe/TCP; refer to Figure 7).													
03	M	Transport Info Version: Implementations compliant to this specification shall set this field to 1h.													
05:04	M	HFI Descriptor Index: The value of the HFI Descriptor Index field of the HFI Descriptor (refer to Figure 11) whose HFI Transport Info Descriptor Heap Object Reference field indicates this HFI Transport Info Descriptor.													
06	M	HFI Transport Flags:													
		<table><tr><th>Bits</th><th>Description</th></tr><tr><td>07:03</td><td>Reserved</td></tr><tr><td>02</td><td>DHCP Override: If set to '1', then HFI information was populated by consuming the DHCP on this interface. If cleared to '0', then the HFI information was set administratively by a configuration interface to the driver and pre-OS environment.</td></tr><tr><td>01</td><td>Global Route vs. Link Local Override Flag: If set to '1', then the BIOS utilized this interface described by HFI to be the default route with highest priority. If cleared to 0, then routes are local to their own scope.</td></tr><tr><td>00</td><td>Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.</td></tr></table>	Bits	Description	07:03	Reserved	02	DHCP Override: If set to '1', then HFI information was populated by consuming the DHCP on this interface. If cleared to '0', then the HFI information was set administratively by a configuration interface to the driver and pre-OS environment.	01	Global Route vs. Link Local Override Flag: If set to '1', then the BIOS utilized this interface described by HFI to be the default route with highest priority. If cleared to 0, then routes are local to their own scope.	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.			
		Bits	Description												
		07:03	Reserved												
		02	DHCP Override: If set to '1', then HFI information was populated by consuming the DHCP on this interface. If cleared to '0', then the HFI information was set administratively by a configuration interface to the driver and pre-OS environment.												
01	Global Route vs. Link Local Override Flag: If set to '1', then the BIOS utilized this interface described by HFI to be the default route with highest priority. If cleared to 0, then routes are local to their own scope.														
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to 0, then this descriptor is reserved.														
10:07	M	PCI Express Routing ID for the HFI Transport Function: This field indicates the PCI Express Routing ID as specified in the PCI Express Base Specification.													
		Note: If the platform supports PCIe ARI, then the Device bits and the Function bits may be concatenated for that purpose.													
		<table><tr><th>Bits</th><th>Description</th><th>Refer to</th></tr><tr><td>31:16</td><td>Segment Group Number</td><td>PCI Firmware Specification</td></tr><tr><td>15:08</td><td>Bus</td><td rowspan="3">PCI Express Base Specification</td></tr><tr><td>07:03</td><td>Device</td></tr><tr><td>02:00</td><td>Function</td></tr></table>	Bits	Description	Refer to	31:16	Segment Group Number	PCI Firmware Specification	15:08	Bus	PCI Express Base Specification	07:03	Device	02:00	Function
		Bits	Description	Refer to											
		31:16	Segment Group Number	PCI Firmware Specification											
15:08	Bus	PCI Express Base Specification													
07:03	Device														
02:00	Function														

Figure 13: HFI Transport Info Descriptor – NVMe/TCP

Bytes	O/M ¹	Description						
16:11	M	MAC Address: The MAC address of this HFI, in EUI-48™ format, as defined in the IEEE Guidelines for Use of Extended Unique Identifiers. This field shall be set to a non-zero value.						
18:17	O	VLAN: If this field is set to a non-zero value, then this field contains the VLAN identifier if the VLAN associated with this HFI, as defined in IEEE 802.1q-2018. If no VLAN is associated with this HFI, then this field shall be cleared to 0h.						
19	O	IP Origin: If this field is set to a non-zero value, then this field indicates the source of Ethernet L3 configuration information used by the driver for this interface. Valid values are defined in the Win 32 API: NL_PREFIX_ORIGIN enumeration specification. This field should be cleared to '0' if the IP Origin field is unused by driver.						
35:20	M	IP Address: This field indicates the IPv4 or IPv6 address of this HFI. This field shall be set to a non-zero value. The format of this field is defined by section 1.5.5.1.						
36	M	Subnet Mask Prefix: This field indicates the IPv4 or IPv6 subnet mask in CIDR routing prefix notation.						
52:37	O	IP Gateway If this field is set to a non-zero value, this field indicates the IPv4 or IPv6 address of the IP gateway for this HFI. If this field is cleared to 0h, then no IP gateway is specified.						
53		Reserved						
55:54	O	Route Metric: If this field is set to a non-zero value, this field indicates the cost value for the route indicated by this HF. This field contains the value utilized by the pre-OS driver when chosing among all available routes. Lower values relate to higher priority. Refer to IETF RFC 4249. If the pre-OS driver supports routing and did not configure a specific route metric for this interface, then the pre-OS driver should set this value to 500. If the pre-OS driver does not support routing, then this field should be cleared to 0h.						
71:56	O	Primary DNS: If this field is set to a non-zero value, this field indicates the IPv4 or IPv6 address of the Primary DNS server for this HFI, if any, from byte offset 0h of the NBFT Table Header. If this field is cleared to 0h, then no Primary DNS is specified.						
87:72	O	Secondary DNS: If this field is set to a non-zero value, this field indicates the IPv4 or IPv6 address of the Secondary DNS server for this HFI, if any, from byte offset 0h of the NBFT Table Header. If this field is cleared to '0', then no Secondary DNS is specified.						
103:88	O	DHCP Server: If the DHCP Override bit is set to '1', then this field indicates the IPv4 or IPv6 address of the DHCP server used to assign this HFI address. If that bit is cleared to '0', then this field is reserved.						
109:104	O	Host Name Heap Object Reference: If this field is set to a non-zero value, then: a) this field indicates the location and size of a heap object containing a Host Name string (refer to section 1.5.6); and b) the Offset field and the Length field shall be set to non-zero values. If this field is cleared to 0h, then this field is reserved. <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object, if any.							
127:110		Reserved						

Notes:
1. O/M definition: O = Optional, M = Mandatory.

3.1.2.4.1.2 HFI Info Payload Host Name

In a DHCP scenario the HFI Info Payload Host Name may be the name stored as the Option 12 host name following the Host Name String Format (refer to section 1.5.6).

3.1.2.5 Subsystem Namespace (SSNS) Descriptor

The Subsystem Namespace Descriptor List (refer to Figure 14) consists of an array of Subsystem and Namespace (SSNS) descriptors, defined in Figure 15, with NumSSNS contiguous entries. Each SSNS descriptor defines the namespace specific parameters required to uniquely identify the attached namespace that was available during system initialization and the boot process. If multiple transports correspond to the same SSNS descriptor there may be multiple SSNS descriptors for the same NID with at least one per transport type.

Figure 14: Subsystem Namespace Descriptor List

Bytes	Description
127:00	Subsystem Namespace Descriptor 1 The first Subsystem Namespace Descriptor (refer to Figure 15), if any.
255:128	Subsystem Namespace Descriptor 2: The second Subsystem Namespace Descriptor (refer to Figure 15), if any.
...	
(n*128)-1:(n-1)*128	Subsystem Namespace Descriptor n: The last Subsystem Namespace Descriptor (refer to Figure 15), if any.

Figure 15: Subsystem Namespace (SSNS) Descriptor

Bytes	O/M ¹	Description						
00	M	Structure ID: This field shall be set to 4h (i.e., SSNS; refer to Figure 5).						
02:01	M	SSNS Descriptor Index: This field indicates the number of this Subsystem Namespace Descriptor in the Subsystem Namespace Descriptor List (i.e., the number 1h represents Subsystem Namespace Descriptor 1, the number 2h represents Subsystem Namespace Descriptor 2, etc.).						
04:03	M	SSNS Flags: Refer to Figure 16.						
05	M	Transport Type: Refer to Figure 7.						
07:06	M	Transport Specific Flags: Refer to Figure 17.						
08	O	<p>Primary Discovery Controller Index: The Discovery Descriptor Index field of the Discovery Descriptor (refer to Figure 24) that is associated with this SSNS Descriptor.</p> <p>If a Discovery controller was used to establish this record this value shall be set to a non-zero value.</p> <p>If this namespace was associated with multiple Discovery controllers, those Discovery controllers shall have records in the Discovery Descriptor to facilitate multi-path rediscovery as required.</p> <p>If no Discovery controller was utilized to inform this namespace record, this field shall be cleared to 0h.</p>						
09		Reserved						
15:10	M	<p>Subsystem Transport Address Heap Object Reference: This field indicates the location and size of a heap object containing the Subsystem Transport Address.</p> <p>For IP based transports types, shall be an IP Address, following the format described in section 1.5.5.1.</p> <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.	05:04	Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.							
05:04	Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.							
21:16	M	<p>Subsystem Transport Service Identifier Heap Object Reference: This field indicates the location and size of a heap object containing an array of bytes indicating the Subsystem Transport Service Identifier.</p> <p>Valid values are found in Figure 7.</p> <table><tr><th>Bytes</th><th>Description</th></tr></table>	Bytes	Description				
Bytes	Description							

Figure 15: Subsystem Namespace (SSNS) Descriptor

Bytes	O/M ¹	Description						
		<table><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.</td></tr></table>	03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.	05:04	Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.		
03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.							
05:04	Length: Length in bytes of the heap object. The Length field shall be set to a non-zero value.							
23:22	M	Subsystem Port ID: Port in the NVM subsystem associated with this transport address used by the pre-OS driver (refer to the NVM Express Base Specification).						
27:24	O	Namespace ID: This field indicates the namespace identifier (NSID) of the namespace indicated by this descriptor. This field shall be cleared to 0h if not specified by the user. If this value is cleared to 0h, then consumers of the NBFT shall rely on the NID.						
28	M	Namespace Identifier Type (NIDT): This field contains the value of the Namespace Identifier Type (NIDT) field in the Namespace Identification Descriptor (refer to the NVM Express Base Specification) for the namespace indicated by this descriptor. If a namespace supports multiple NIDT entries for uniqueness, the order of preference is NIDT field value of 3h (i.e., UUID) before 2h (i.e., NSGUID), and 2h before 1h (i.e., EUI-64).						
44:29	M	Namespace Identifier (NID): This field contains the value of the Namespace Identifier (NID) field in the Namespace Identification Descriptor (refer to the NVM Express Base Specification) for the namespace indicated by this descriptor.						
45	O	Security Profile Descriptor Index: If the Use Security Flag bit in the SSNS Flags field is set to '1', then this field indicates the value of the Security Profile Descriptor Index field of the Security Profile Descriptor (refer to Figure 21) associated with this namespace. If the Use Security Flag bit is cleared to '0', then no Security Profile Descriptor is associated with this namespace and this field is reserved.						
46	M	Primary HFI Descriptor Index: This field indicates the value of the HFI Descriptor Index field of the HFI Descriptor (refer to Figure 11) for the interface associated with this namespace. If multiple HFIs are associated with this record, subsequent interfaces should be populated in the Secondary HFI Associations field.						
47		Reserved						
53:48	O	Secondary HFI Associations Heap Object Reference: If this field is set to a non-zero value, then: a) this field indicates an array of bytes, in which each byte contains the value of the HFI Descriptor Index field of an HFI Descriptor in the HFI Descriptor List (refer to Figure 10); and b) the Offset field and the Length field shall be set to non-zero values. If this field is cleared to 0h, then no secondary HFI associations are specified. <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object, if any.							
59:54	M	Subsystem and Namespace NQN Heap Object Reference: This field indicates the location and size of a heap object containing the Subsystem and Namespace NQN (refer to the NVM Express Base Specification). <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any. The Length field shall be set to a non-zero value.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.	05:04	Length: Length in bytes of the heap object, if any. The Length field shall be set to a non-zero value.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header. The Offset field shall be set to a non-zero value.							
05:04	Length: Length in bytes of the heap object, if any. The Length field shall be set to a non-zero value.							
65:60	O	SSNS Extended Information Descriptor Heap Object Reference: If the SSNS Extended Info In-use Flag bit is set to '1', then:						

Figure 15: Subsystem Namespace (SSNS) Descriptor

Bytes	O/M ¹	Description						
		<div><div><div>a) this field indicates the offset in bytes of a heap object containing an SSNS Extended Information Descriptor (refer to Figure 19) heap object from byte offset 0h of the NBFT Table Header; and</div><div>b) the Offset field and the Length field shall be set to non-zero values.</div></div><div>If the SSNS Extended Info In-use Flag bit is cleared to '0', then this field is reserved.</div><table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table></div>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object, if any.							
127:66		Reserved						

Notes:

1. O/M definition: O = Optional, M = Mandatory.

3.1.2.5.1 Subsystem and Namespace Descriptor Flags

Flags for the Subsystem and Namespace Descriptor (refer to Figure 17).

Figure 16: Subsystem and Namespace Specific Flags Field

Bits	Description
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is not valid. A host that supports NVMe-oF Boot, but does not currently have a remote Subsystem and Namespace assigned may clear this bit to '0'.
01	Non-bootable Entry Flag: If set to '1', this flag indicates that this SSNS Descriptor contains a namespace of administrative purpose to the boot process, but the pre-OS may not have established connectivity to or evaluated the contents of this Descriptor. Such namespaces may contain supplemental data deemed relevant by the Administrator as part of the pre-OS to OS hand off. This may include properties such as a UEFI device path that may not have been created for this namespace. This means an OS runtime may still require the contents of such a namespace to complete later stages of boot. If cleared to '0', then this namespace did not have any special administrative intent.
02	Use Security Flag: If set to '1', then there is a Security Profile Descriptor associated with this SSNS record and the Security Profile Descriptor Index field is valid. If cleared to '0', then there is no Security Profile Descriptor associated with this SSNS record and the Security Profile Descriptor Index field is not valid.
03	DHCP Root-Path Override Flag: If set to '1', then this SSNS descriptor was populated by consuming the DHCP Root-Path (refer to section 3.1.2.5.3) on this interface. If cleared to '0', then the DHCP Root-Path was not used in populating the SSNS descriptor.
04	SSNS Extended Info In-use Flag: If set to '1', then the SSNS Extended Information Offset field and the SSNS Extended Information Length field are valid. This flag, if set to '1', indicates that a Subsystem and Namespace Extended Information Descriptor corresponding to this descriptor is present.
05	Separate Discovery Controller Flag: If set to '1', then the Discovery controller associated with this volume is on a different transport address than the specified in the Subsystem Transport Address Heap Object Reference. If cleared to '0', then the Discovery controller is the same as the Subsystem Transport Address Heap Object Reference.
06	Discovered Namespace Flag: If set to '1', then this namespace was acquired through discovery. If cleared to '0', then this namespace was explicitly configured in the system.

Figure 16: Subsystem and Namespace Specific Flags Field

Bits	Description	
08:07	Unavailable Namespace Flag: This field indicates the availability of the namespace at a specific point in time. Such use is only a hint and its use does not guarantee the availability of that referenced namespace at any future point in time.	
	Value	Description
	00b	Not Indicated by Driver: No information is provided.
	01b	Available: A referenced namespace described by this flag was previously accessible by the pre-OS driver.
	10b	Unavailable: This namespace was administratively configured but unattempted, unavailable or inaccessible when establishing connectivity by the pre-OS driver.
	11b	Reserved
15:09	Reserved	

3.1.2.5.2 Transport Specific Flags

The Transport Specific Flags (refer to Figure 17) are flags which have specific meaning depending on the Transport Type of the SSNS Descriptor they belong to. Aside from the Transport Specific Flags In Use bit (i.e., Bit 0), which indicates the use of this field, all other bits shall rely on their association with the Transport column of Figure 17.

Figure 17: SSNS Transport Specific Flags Field

Transport	Bits	Description
Any	0	Transport Specific Flags in Use: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is not valid.
3h	01	PDU Header Digest (HDGST) Flag: If set to '1', then the host or administrator required the connection described by this Subsystem and Namespace Descriptor to use the NVM Header Digest Enabled. A consumer of this information should attempt to use NVM Header Digest when recreating this connection if enabled. If cleared to '0', then the host or administrator did not require the connection described by this Subsystem and Namespace Descriptor to use the NVM Header Digest Enabled.
	02	Data Digest (DDGST) Flag: If set to '1', then the host or administrator required the connection described by this Subsystem and Namespace Descriptor to use the NVM Data Digest Enabled. If cleared to '0', then the host or administrator did not require the connection described by this Subsystem and Namespace Descriptor to use the NVM Data Digest Enabled. A consumer of this field should attempt to use NVM Data Digest when recreating this connection if enabled.
	15:03	Reserved

3.1.2.5.3 DHCP Root-Path

The format of this field shall follow the format defined by IETF RFC 2132, DHCP Option 17, which specifies the path name that contains the the root disk for the host. The path is formatted with a hierarchical NVMe-oF URI Format (refer to Figure 7), as defined in section 4.2.4, with UTF-8 characters and formatting otherwise as specified in RFC 3986. A definition of the valid URI tokens and their mandatory or optional requirements are specified in Figure 18.

Valid representations are:

- “NVME<+PROTOCOL>://<SERVERNAME/IP>[:TRANSPORT PORT]/<SUBSYS NQN>/<NID>”;
or
- “NVME<+PROTOCOL>://<DISCOVERY CONTROLLER ADDRESS>[:DISCOVERY CONTROLLER PORT]/NQN.2014-08.ORG.NVMEXPRESS.DISCOVERY/<NID>”

Figure 18: DHCP Root-Path Elements

Token	O/M ¹	Description
<PROTOCOL>	M	Value shall correspond to valid Transport Type string designator as defined in section 2.3.1.1.1.3.
<SERVER NAME/IP>	M	The transport address of the Subsystem Port. IPv6 address shall be square brace enclosed per IETF RFC 2732.
<TRANSPORT PORT>	O	If Transport Port is blank, assume 4420 for IP Based Transports.
<DISCOVERY CONTROLLER ADDRESS>	O	If Discovery Controller Address is blank, assume integrated with transport address.
<DISCOVERY CONTROLLER PORT>	O	For protocol types: tcp If Discovery Controller Port is unused, but well known Discovery Subsystem NQN Address is populated a driver shall set this value to 8009. If unused but there is no well-known NQN, then this value may be set to the default as defined by the Transport Type in Figure 7.
<SUBSYSTEM NQN>	M	NQN of the Subsystem being indicated by this record.
<NID>	M	Desired namespace (i.e., NVMe NID) that the host shall attach to from available Namespaces
Notes: 1. O/M definition: O = Optional, M = Mandatory.		

3.1.2.5.4 Primary Discovery Controller Index Field

The Primary Discovery Controller Index field indicates the primary or first Discovery controller associated with a namespace as registered in the NBFT. In the cases of hierarchical Discovery controllers, the last controller in a topology that is encountered that refers connectivity to the owning NVM subsystem I/O controller should be specified by this field. In the cases of multipath Discovery controllers, only one Discovery controller per Namespace path shall be tracked by this field but all Discovery controllers in the chain should be registered in the Discovery Descriptor.

If this Discovery controller requires Security Parameters for connectivity, those values shall be listed in the Discovery Descriptor associated with this Index record.

3.1.2.5.5 Subsystem and Namespace Extended Information Descriptor

The Subsystem and Namespace Extended Information Descriptor, defined in Figure 19, when utilized, contains additional transport specific detail relevant to a host reestablishing connectivity to the referenced namespaces that may otherwise only be available through administrative configuration and without which some OS runtime environments may operate deficiently. This information is otherwise available and provided to the pre-OS environment through configuration means beyond the scope of this document. An OS driver should treat this information as informational and may require further advanced functionality to properly leverage.

The Subsystem and Namespace Extended Information Descriptor is structured as a standard table located in heap space as required.

Figure 19: Subsystem and Namespace Extended Information Descriptor

Bytes	O/M ¹	Description								
00	M	Structure ID: This field shall be set to 9h (i.e., SSNS Extended Info; refer to Figure 5).								
01	M	Version: This field shall be set to 1h.								
03:02	M	SSNS Descriptor Index: This field indicates the value of the SSNS Descriptor Index field of the Subsystem and Namespace Descriptor (refer to Figure 15) whose SSNS Extended Information Descriptor Heap Object Reference field indicates this descriptor.								
07:04	M	Flags:								
		<table><tr><th>Bits</th><th>Feature</th></tr><tr><td>31:02</td><td>Reserved</td></tr><tr><td>01</td><td>Administrative ASQSZ: If set to '1', then the value of the ASQSZ field was provided by administrative configuration for this SSNS record. If cleared to 0h, then the value of the ASQSZ field was either obtained by discovery or assumed by the driver.</td></tr><tr><td>00</td><td>Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.</td></tr></table>	Bits	Feature	31:02	Reserved	01	Administrative ASQSZ: If set to '1', then the value of the ASQSZ field was provided by administrative configuration for this SSNS record. If cleared to 0h, then the value of the ASQSZ field was either obtained by discovery or assumed by the driver.	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.
		Bits	Feature							
		31:02	Reserved							
		01	Administrative ASQSZ: If set to '1', then the value of the ASQSZ field was provided by administrative configuration for this SSNS record. If cleared to 0h, then the value of the ASQSZ field was either obtained by discovery or assumed by the driver.							
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.									
09:08	O	Controller ID: The controller identifier of the first controller associated with the Admin Queue by the driver. If a controller identifier is not administratively specified or direct configuration is not supported by the driver, then this field shall be cleared to 0h.								
11:10	M	Admin Submission Queue Size (ASQSZ): The Admin Submission Queue Size utilized for the respective SSNS by the driver.								
17:12	O	DHCP Root Path String Heap Object Reference: If the SSNS DHCP Root Path Override flag bit is set to '1', then: a) this field indicates the offset in bytes of a heap object containing an DHCP Root Path String as defined in section 3.1.2.5.3 used by the driver. b) the Offset field and the Length field shall be set to non-zero values. If the SNSS DHCP Root Path Overrid flag bit is cleared to '0', then this field is reserved.								
		<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.		
		Bytes	Description							
		03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
		05:04	Length: Length in bytes of the heap object, if any.							

Notes:
1. O/M definition: O = Optional. M = Mandatory.

3.1.2.6 Security Descriptor

The Security Descriptor List (refer to Figure 20) consists of an array of Security descriptors, as defined in Figure 21, with NumSec contiguous entries. Each Security descriptor defines the various connectivity policy features of authentication and transport security related to NVMe connections to SSNS or Discovery descriptors. Security Profile Descriptors may be referenced in other descriptors, (i.e., the SSNS or discovery) using their list indices.

Figure 20: Security Profile Descriptor List

Bytes	Description
63:00	Security Profile Descriptor 1: The first Security Profile Descriptor (refer to Figure 21), if any.
127:64	Security Profile Descriptor 2: The second Security Profile Descriptor (refer to Figure 21), if any.
...	
(n*64)-1:(n-1)*64	Security Profile Descriptor n: The last Security Profile Descriptor (refer to Figure 21), if any.

Figure 21: Security Profile Descriptor

Bytes	O/M ¹	Description								
00	M	Structure ID: This field shall be set to 5h (i.e., Security; refer to Figure 5).								
01	M	Security Profile Descriptor Index: This field indicates the number of this Security Profile Descriptor in the Security Profile Descriptor List (i.e., the number 1h represents Security Profile Descriptor 1, the number 2h represents Security Profile Descriptor 2, etc.).								
03:02	M	Security Profile Descriptor Flags: Refer to Figure 22.								
04	M	Secret Type:								
		<table><tr><th>Bits</th><th>Description</th></tr><tr><td>07:02</td><td>Reserved</td></tr><tr><td>01</td><td>Redfish Host Interface URI: If set to '1', then the Secret Keypath Object Reference is a URI pointing to a Redfish Key Collection Object that contains the PSK. Relies on DMTF DSP 0270 v1.3 or later and DMTF DSP8010 2021.4 or later. If cleared to '0', then the Secret String is not a Redfish Key Collection Object.</td></tr><tr><td>00</td><td>Reserved</td></tr></table>	Bits	Description	07:02	Reserved	01	Redfish Host Interface URI: If set to '1', then the Secret Keypath Object Reference is a URI pointing to a Redfish Key Collection Object that contains the PSK. Relies on DMTF DSP 0270 v1.3 or later and DMTF DSP8010 2021.4 or later. If cleared to '0', then the Secret String is not a Redfish Key Collection Object.	00	Reserved
		Bits	Description							
		07:02	Reserved							
		01	Redfish Host Interface URI: If set to '1', then the Secret Keypath Object Reference is a URI pointing to a Redfish Key Collection Object that contains the PSK. Relies on DMTF DSP 0270 v1.3 or later and DMTF DSP8010 2021.4 or later. If cleared to '0', then the Secret String is not a Redfish Key Collection Object.							
00	Reserved									
05		Reserved								
11:06	O	Secure Channel Algorithm Heap Object Reference: If the Security Policy List field is set to 01b, then: a) this field indicates the location and size of a heap object containing a list of secure channel algorithms from byte offset 0h of the NBFT Table Header; and b) the Offset field and the Length field shall be set to non-zero values. The list is an array of bytes and the values are defined in the Security Type (SECTYPE) field in the Transport Specific Address Subtype Definition in the NVMe TCP Transport Specification. If the Security Policy List field is cleared to 0h, then this field is reserved.								
		<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.		
		Bytes	Description							
		03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
		05:04	Length: Length in bytes of the heap object, if any.							
17:12	O	Authentication Protocols Heap Object Reference: If the Authentication Policy List field is set to 01b, then: a) this field indicates the location and size of a heap object containing a list of authentication protocol identifiers (refer to the NVM Express Base Specification) from byte offset 0h of the NBFT Table Header; and b) the Offset field and the Length field shall be set to non-zero values. If the Authentication Policy List field is cleared to 0h, then this field is reserved.								
		<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.		
		Bytes	Description							
		03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
		05:04	Length: Length in bytes of the heap object, if any.							
23:18	O	Cipher Suite Offset Heap Object Reference: If the Cipher Suites Restricted by Policy bit is set to '1', then: a) this field indicates the location and size of a heap object containing a list of cipher suite identifiers from byte offset 0h of the NBFT Table Header; and b) the Offset field and the Length field shall be set to non-zero values. The list, if any, is an array of bytes and the values are defined in the IANA TLS Parameters Registry. If the Cipher Suites Restricted by Policy bit is cleared to '0', then this field is reserved.								
		<table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.				
		Bytes	Description							
		03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							

Figure 21: Security Profile Descriptor

Bytes	O/M ¹	Description						
		05:04 Length: Length in bytes of the heap object, if any.						
29:24	O	<p>DH Groups Heap Object Reference: If the Authentication DH Groups Restricted by Policy List bit is set to '1', then:</p> <ul style="list-style-type: none">a) this field indicates the location and size of a heap object containing a list of DH-HMAC-CHAP Diffie-Hellman (DH) group identifiers (refer to the NVM Express Base Specification) from byte offset 0h of the NBFT Table Header; andb) the Offset field and the Length field shall be set to non-zero values. <p>If the Authentication DH Groups Restricted by Policy List bit is cleared to '0', then this field is reserved.</p> <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object.							
35:30	O	<p>Secure Hash Functions Offset Heap Object Reference: If the Secure Hash Functions Policy List bit is set to '1', then:</p> <ul style="list-style-type: none">a) this field indicates the offset in bytes of a heap object containing a list of DH-HMAC-CHAP hash function identifiers from byte offset 0h of the NBFT Table Header; andb) the Offset field and the Length field shall be set to non-zero values. <p>The list is an array of bytes and the values are defined in the NVM Express Base Specification.</p> <p>If the Secure Hash Functions Policy List bit is cleared to '0', then this field is reserved.</p> <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object, if any.							
41:36	O	<p>Secret Keypath Offset Heap Object Reference: if this field is set to a non-zero value, then:</p> <ul style="list-style-type: none">a) this field indicates the location and size of a heap object containing a URI; andb) the Offset field and the Length field shall be set to non-zero values. <p>The type of the URI is specified in the Secret Type field.</p> <p>If this field is cleared to 0h, then this field is reserved.</p> <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object, if any.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object, if any.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object, if any, from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object, if any.							
63:42		Reserved						

Notes:
1. O/M definition: O = Optional, M = Mandatory.

Figure 22: Security Profile Descriptor Flags

Bits	Description
15:12	Reserved
11	Secure Hash Functions Policy List: If set to '1', then connections shall use one of the secure hash functions in the Secure Hash Functions Policy List is required. If cleared to '0', then no Secure Hash Functions Policy List is indicated and use of a secure hash function is not required.
10	Authentication DH Groups Restricted by Policy List: If set to '1', then connections shall use one of the authentication DH groups in the Authentication DH Groups List is required. If cleared to '0', then no Authentication DH Groups List is indicated and use of an authentication DH Group is not required.

Figure 22: Security Profile Descriptor Flags

Bits	Description										
09	Cipher Suites Restricted by Policy: If set to '1', then the Cipher Suite Offset field and the Cipher Suite Length field indicate a list of supported cipher suites by the driver. If cleared to '0', then the Cipher Suite Offset field and the Cipher Suite Length field are reserved.										
08:07	Security Policy List: This field indicates whether secure channel protocols were indicated by policy from driver defaults or administrative configuration.										
	<table><tr><th>Value</th><th>Description</th></tr><tr><td>00b</td><td>The Offset field and Length field in the Secure Channel Algorithm Heap Object Reference field are reserved.</td></tr><tr><td>01b</td><td>The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols used by the driver.</td></tr><tr><td>10b</td><td>The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols that were administratively set and used by the driver.</td></tr><tr><td>11b</td><td>Reserved</td></tr></table>	Value	Description	00b	The Offset field and Length field in the Secure Channel Algorithm Heap Object Reference field are reserved.	01b	The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols used by the driver.	10b	The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols that were administratively set and used by the driver.	11b	Reserved
	Value	Description									
	00b	The Offset field and Length field in the Secure Channel Algorithm Heap Object Reference field are reserved.									
	01b	The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols used by the driver.									
10b	The Heap Object specified by the Secure Channel Algorithm Heap Object Reference field indicates a list of authentication protocols that were administratively set and used by the driver.										
11b	Reserved										
06:05	Secure Channel Negotiation Required: This field indicates whether secure channel negotiation (e.g., TLS) is required.										
	<table><tr><th>Value</th><th>Description</th></tr><tr><td>00b</td><td>Secure channel negotiation is not supported by the NVM subsystem.</td></tr><tr><td>01b</td><td>Secure channel negotiation is supported by the NVM subsystem and is not required.</td></tr><tr><td>10b</td><td>Secure channel negotiation is supported by the NVM subsystem and is required.</td></tr><tr><td>11b</td><td>Reserved</td></tr></table>	Value	Description	00b	Secure channel negotiation is not supported by the NVM subsystem.	01b	Secure channel negotiation is supported by the NVM subsystem and is not required.	10b	Secure channel negotiation is supported by the NVM subsystem and is required.	11b	Reserved
	Value	Description									
	00b	Secure channel negotiation is not supported by the NVM subsystem.									
	01b	Secure channel negotiation is supported by the NVM subsystem and is not required.									
10b	Secure channel negotiation is supported by the NVM subsystem and is required.										
11b	Reserved										
04:03	Authentication Policy List: This field indicates whether authentication protocols were indicated by policy from driver defaults or administrative configuration.										
	<table><tr><th>Value</th><th>Description</th></tr><tr><td>00b</td><td>Authentication Protocols Heap Object Reference field Offset and Length are reserved.</td></tr><tr><td>01b</td><td>Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols used by the driver.</td></tr><tr><td>10b</td><td>Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols that were administratively set and used by the driver.</td></tr><tr><td>11b</td><td>Reserved</td></tr></table>	Value	Description	00b	Authentication Protocols Heap Object Reference field Offset and Length are reserved.	01b	Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols used by the driver.	10b	Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols that were administratively set and used by the driver.	11b	Reserved
	Value	Description									
	00b	Authentication Protocols Heap Object Reference field Offset and Length are reserved.									
	01b	Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols used by the driver.									
10b	Authentication Protocols Offset field and the Authentication Protocols Length field indicate a list of authentication protocols that were administratively set and used by the driver.										
11b	Reserved										
02:01	In-Band Authentication Required:										
	<table><tr><th>Value</th><th>Description</th></tr><tr><td>00b</td><td>In-band authentication is not supported by the NVM subsystem.</td></tr><tr><td>01b</td><td>In-band authentication is supported by the NVM subsystem and is not required.</td></tr><tr><td>10b</td><td>In-band authentication is supported by the NVM subsystem and is required.</td></tr><tr><td>11b</td><td>Reserved</td></tr></table>	Value	Description	00b	In-band authentication is not supported by the NVM subsystem.	01b	In-band authentication is supported by the NVM subsystem and is not required.	10b	In-band authentication is supported by the NVM subsystem and is required.	11b	Reserved
	Value	Description									
	00b	In-band authentication is not supported by the NVM subsystem.									
	01b	In-band authentication is supported by the NVM subsystem and is not required.									
10b	In-band authentication is supported by the NVM subsystem and is required.										
11b	Reserved										
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is not valid.										

The Secure Channel Algorithm List is an array of bytes, each byte of which contains a value indicating a supported security type. The length of the list is indicated by the Secure Channel Algorithm Length field. Refer to the Security Type (SECTYPE) field in the Transport Specific Address Subtype Definition in the NVMe TCP Transport Specification.

3.1.2.7 Discovery Descriptor

The Discovery Descriptor is an indexed structure that contains all Discovery controllers that have been collected by a pre-OS driver, such as UEFI BIOS, prior to OS execution within Discovery Records. The intent of this table is to allow an OS client to both see pre-OS environment context and to serve OS clients who may not otherwise dynamically discover subsystems but are able to perform basic connections to Discovery controllers as outlined in the NVMe Express Base Specification to acquire additional namespaces.

The NumDisc field in the NBFT Header Structure (refer to Figure 8) indicates how many Discovery Records are in the Discovery Descriptor List (refer to Figure 23).

Figure 23: Discovery Descriptor List

Bytes	Description
31:00	Discovery Descriptor 1: The first Discovery Descriptor (refer to Figure 24), if any.
63:32	Discovery Descriptor 2: The second Discovery Descriptor (refer to Figure 24), if any.
...	
(n*32)-1:(n-1)*32	Discovery Descriptor n: The last Discovery Descriptor (refer to Figure 24), if any.

Figure 24: Discovery Descriptor

Bytes	O/M ¹	Description						
00	M	Structure ID: This field shall be set to 6h (i.e., Discovery Descriptor; refer to Figure 5).						
01	M	Discovery Descriptor Flags: <table><tr><th>Bits</th><th>Feature</th></tr><tr><td>07:01</td><td>Reserved</td></tr><tr><td>00</td><td>Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.</td></tr></table>	Bits	Feature	07:01	Reserved	00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.
Bits	Feature							
07:01	Reserved							
00	Descriptor Valid: If set to '1', then this descriptor is valid. If cleared to '0', then this descriptor is reserved.							
02	M	Discovery Descriptor Index: This field indicates the number of this Discovery Descriptor in the Discovery Descriptor List (i.e., the number 1h represents Discovery Descriptor 1, the number 2h represents Discovery Descriptor 2, etc.).						
03	M	HFI Descriptor Index: This field indicates the value of the HFI Descriptor Index field of the HFI Descriptor associated with this Discovery Descriptor. If multiple HFIs share a common Discovery controller, there shall be multiple Discovery Descriptor entries with one per HFI.						
04	O	Security Profile Descriptor Index: This field indicates the value of the Security Profile Descriptor Index field of the Security Descriptor associated with this Discovery Descriptor.						
05		Reserved						
11:06	M	Discovery Controller Address Heap Object Reference: This field indicates the location and size of a heap object containing a URI (refer to section 1.5.7) which indicates an NVMe Discovery controller associated with this Discovery Descriptor. If this field is cleared to 0h, then no URI is specified. <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object.							
17:12	O	Discovery Controller NQN Heap Object Reference: If set to a non-zero value, this field indicates the location and size of a heap object containing an NVMe Discovery controller NQN. If the NVMe Discovery controller referenced by this record requires secure authentication with a well known Subsystem NQN, this field indicates the unique NQN for that NVMe Discovery controller. This record is involved formatted as an NQN string as defined in the NVM Express Base Specification. If this field is cleared to 0h, then: a) this field is reserved; and b) the OS shall use the well known discovery NQN for this record. <table><tr><th>Bytes</th><th>Description</th></tr><tr><td>03:00</td><td>Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.</td></tr><tr><td>05:04</td><td>Length: Length in bytes of the heap object.</td></tr></table>	Bytes	Description	03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.	05:04	Length: Length in bytes of the heap object.
Bytes	Description							
03:00	Offset: Offset in bytes of the heap object from byte offset 0h of the NBFT Table Header.							
05:04	Length: Length in bytes of the heap object.							
31:18		Reserved						

Notes:
1. O/M definition: O = Optional, M = Mandatory.

Annex A. Booting (Informative)

This section introduces booting over NVMe-oF concepts; it also references other boot-from-SAN technologies (e.g., iSCSI and FCOE) as examples, however booting over NVMe-oF architecture is designed to be fabric-agnostic. Section A.1 provides the summary of general concepts related to NVMe-oF boot.

A.1 Domain Specific Terms

This domain specific terms section contains terminology that is specific to certain domain topics, such as UEFI that are not generally relevant to the normative sections but are still important for informative content.

A.1.1 SCSI Logical Unit Number (LUN)

A LUN is a SCSI device address of the logical unit on a target device.

A.1.2 UEFI Boot Device

The device handle that corresponds to the device from which the currently executing image was loaded as defined in Unified Extensible Firmware Interface (UEFI) Specification.

A.1.3 UEFI Boot Services

Refer to Unified Extensible Firmware Interface (UEFI) Specification. The collection of interfaces and protocols that are present in the boot environment. The services minimally provide an OS loader with access to platform capabilities required to complete OS boot. Services are also available to drivers and applications that require access to platform capability. Boot services are terminated once the operating system takes control of the platform.

A.1.4 UEFI Device Handle

A handle points to a list of one or more UEFI protocols that are able to respond to requests for services for a given device referred to by the handle. Refer to the Unified Extensible Firmware Interface (UEFI) Specification.

A.1.5 UEFI Device Path

A variable-length binary data structure that is composed of variable-length generic device path nodes and is used to define the programmatic path to a logical or physical device. Refer to the Unified Extensible Firmware Interface (UEFI) Specification.

A.1.6 UEFI Device Path Protocol

A protocol used during boot services to provide the information required to construct and manage Device Paths. Also called the EFI Device Path Protocol. Refer to the Unified Extensible Firmware Interface (UEFI) Specification.

A.1.7 UEFI Protocol

A UEFI Protocol is a firmware interface between UEFI applications and the set of services offered by the UEFI firmware. Protocols serve to enable communication between separately built modules including drivers. A UEFI Protocol provides a block of function pointers and data, named by a GUID and associated with a UEFI handle. Refer to the EDK II UEFI Driver Writer's Guide.

A.2 Boot from SAN Ontology and Concepts

“Boot from SAN” is a term typically used to refer to the capability to boot from SCSI-based Block Storage where the technology of the “SCSI bus” is a network (commonly referred to as a “Storage Area Network” or “SAN”). The storage area network is typically an FC or IP-based network. The storage devices are attached to the network and not to the booting system (e.g., local storage with IDE/SATA). SANs also support multiple access paths via the fabric, through multiple SCSI host adapters and multiple SCSI targets, to the same SCSI block storage devices. This allows the booting system to support redundant and parallel connections to storage devices.

A.2.1 Boot from Block Storage Device

The boot environment traditionally begins by the system legacy BIOS or system firmware (e.g., UEFI) enumerating host adapter hardware and invoking PCIe Expansion ROM images which initialize and configure the host adapter. The Expansion ROM images on storage adapters detect and enumerate supported block storage devices as part of the Expansion ROM's initialization.

A system boot device list is configured with a small number of devices that the system attempts to boot from. The boot firmware selects the first boot device in the list and attempt to read from the device to load the OS boot image into memory. Upon a failure, the boot firmware transitions to the next boot device in the list and repeats the OS load attempt. The process repeats until the OS is successfully loaded via a device in the list. At such time, any remaining devices in the list are ignored. A failure to successfully load a device typically results in the administratively configured behavior for any boot device failure (i.e., a reboot) or a halt at prompt.

In most cases, an administratively managed list of boot devices is used to set the boot order of the system. For example, a prioritized list may start with CD/DVD device, with second USB device, followed by hard disk drive, where the transition between devices is expected to occur based on non-presence of leading devices.

In cases where the same block storage entity may be accessible via multiple host adapters, there are multiple block storage devices enumerated to the boot environment which all map to the same device. In this scenario, the boot device list may be used as a multi-path selection sequence to allow system booting, even if a host adapter fails or loses connectivity to storage.

Simplistic host adapters such as IDE or SATA interface typically enumerate 1 or 2 storage entities based on their physical connectivity to a storage device. More complex host adapters may use the adapter's Expansion ROM image to configure the host adapter with options and settings that allow it to abstract physical connectivity. For example, an IDE or SATA host adapter may support multiple physically attached drives to be abstracted into a single block storage device presented to the boot environment. The host adapter may utilize the physical devices in one of the following manners:

- concatenated to appear as a higher capacity device;
- perform RAID 0 for striping for capacity and performance;
- perform RAID 1 for mirroring for fault-tolerance; or
- some other RAID type.

In general, each time a system boots, simplistic host adapters enumerate the same block storage devices to the boot environment as the storage devices are physically connected to the adapter. Only physical removal of the devices results in a change to the enumerated list. Given the fixed physical connectivity, the devices detected and enumerated in the boot environment are the same devices automatically detected and enumerated in the OS environment.

A.2.2 Boot from SCSI-based Block Storage

SCSI host storage adapters follow the basic behaviors described for Boot from Block Storage Devices. SCSI host adapters are typically attached to “a bus” which is able to connect to multiple SCSI storage devices. The technology of the “bus” may vary e.g., parallel SCSI, SAS, Fibre Channel Arbitrated Loop, Fibre Channel fabric, or Ethernet/IP. The storage adapter uses methods specific to the technology of the bus to detect SCSI targets, and the SCSI command set to scan for SCSI Logical Unit Numbers (LUN), detect block storage devices, and enumerate the block storage devices with the boot environment. The storage adapter converts the basic block storage device interface (i.e., select/read/write/deselect) used by the boot environment to the SCSI protocol commands for use with the SCSI LUNs.

The introduction of bus topologies loosens the physical connectivity of devices to the host adapter. A SCSI backplane with bus connections for SCSI target devices allow for dynamic device removal and insertion, resulting in possible dynamic device enumeration and changes in device enumeration order. Network topologies using switches and routers allow the storage to become even more “remote” and possibly much more dynamic in terms of addition, removal, and detection order.

Many SCSI bus topologies support connectivity to large numbers of storage devices. A single SCSI target may support many SCSI LUNs, each LUN being a block storage device. Enumerating all devices may overwhelm the boot environment. Over time SCSI adapters have expanded their Expansion ROM images to support administrative configuration which selectively identifies the SCSI target and SCSI LUN's that are to be enumerated with the boot environment. The adapter typically saves this configuration in non-volatile store on the adapter. This change has resulted in faster boot enumeration (i.e., fewer targets/LUNs connect to and/or enumerated) and better manageability (i.e., fewer devices enumerated to the boot environment and with more consistent order).

As most SCSI topologies are self-discovering, when the OS environment re-initializes and rescans for SCSI devices, the OS environment rediscovers the same storage devices enumerated with the boot environment as those storage devices are a subset of the overall devices that are discovered by the OS (barring dynamic changes in connectivity). There is no required interaction between the boot environment and the OS to automatically re-establish access to the boot device.

A.2.3 Boot from SAN: Fibre Channel Storage

Boot from Fibre Channel (FC) storage is an example of Boot from SCSI-based Block Storage. Most references to Boot from SAN refer to the administrative settings made per host adapter to identify the FC-based SCSI target devices and the SCSI LUN numbers to enumerate with the boot environment. The FC SCSI host adapters typically perform device and LUN discovery, allow the administrator to select the block storage devices to be presented to the boot environment, and save the configuration on the adapter.

FC storage is self-discovering. The host adapter typically sees the same devices regardless of whether it is operating in the boot environment or the OS environment. The host adapter expansion ROM may filter those presented to the boot environment, however the devices presented to the OS environment are typically all devices, including those presented to the boot environment. Thus, device presentation for boot devices occurs in both the boot environment and the OS environment without requiring the two environments to interact.

A.2.4 Boot from SAN: iSCSI Storage

Boot from iSCSI storage is another example of booting from SCSI-based block storage. Most references to Boot from SAN refer to the administrative settings made per host adapter to specify iSCSI session information for the iSCSI target(s) and SCSI LUN number for the block storage devices to enumerate with

the boot environment. The settings may also include IP addresses for discovery services to help locate the iSCSI target device.

In contrast to FC storage, iSCSI storage is not self-discovering. The host adapter, or iSCSI stack utilizing a network adapter, requires information to locate and/or create the iSCSI session. Information includes IP addresses and port numbers for the initial TCP connection, iSCSI names to use for the host initiator, and iSCSI names for the iSCSI target to connect to. To aid in the amount of information that is required, IP services such as DHCP, iSNS, or SLP may be used to query for the iSCSI boot device session information. Even in those cases, address information and service type information may have to be specified. Information on such services is found in IETF RFC 4173.

As iSCSI storage is not self-discovering, when the OS environment is transitioned to, without additional aid, the OS environment does not know that it should establish an iSCSI session nor what the configuration information is for that session. As such, OS environments have to be prepared by an administrator with the same connectivity information that was in the boot environment. Not only does this require a priori knowledge and configuration of both environments, but it also presents situations where the two environments may become out of sync.

To resolve these issues, the industry created the iSCSI Boot Firmware Table (iBFT). The iBFT is a data structure that is created by the boot environment and passed to the OS environment using the ACPI API. The table contains the iSCSI boot session information that was used by the boot environment so that the OS environment realizes that it is to connect to iSCSI storage and to use the session information in the table for the connectivity.

A.2.5 Boot from NVMe-oF

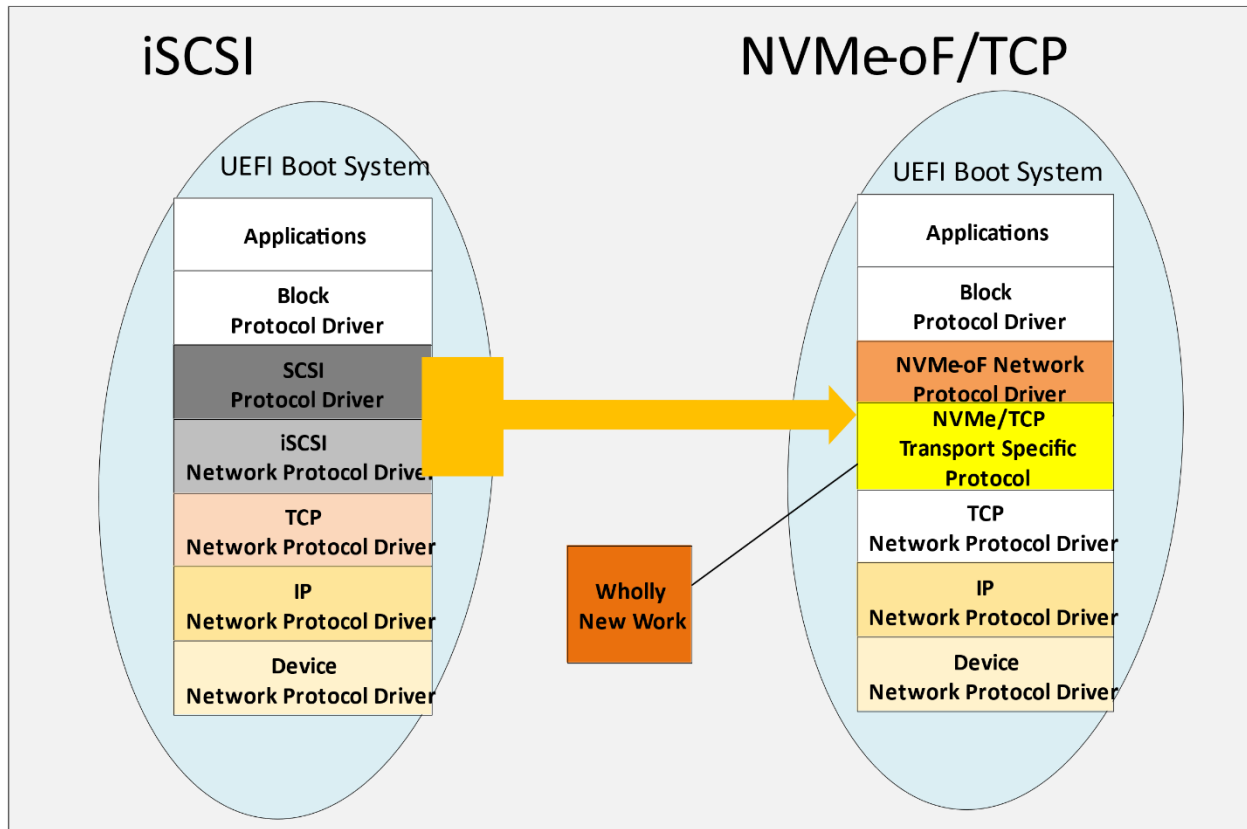
NVMe-oF boot addresses the technology requirements of remote booting from NVMe namespaces. It is similar to other remote block storage boot technologies (i.e., main concepts such as boot flow and handover mechanism are similar).

NVMe-oF provides:

- fabric agnostic capability supporting all NVMe-oF fabric transports;
- lower latency due to elimination of SCSI command translation via direct NVMe command execution; and
- increased scalability and flexibility through the NVMe-oF defined discovery mechanism.

Figure 25 outlines a hypothetical example of UEFI driver stack differences between iSCSI and NVMe-oF (using TCP as transport).

Figure 25: UEFI Network stack differences between booting from iSCSI and booting from NVMe/TCP



A.2.6 Configure Boot from NVMe-oF

Before booting from NVMe-oF, administrative configuration of the UEFI pre-OS driver is required. As an example, an administrator would configure a UEFI system through UEFI HII (Human Interface Infrastructure) to include the required information such as target Subsystem NQN, target IP address, target port number, target namespace, host NQN, security related information. How this information is stored by UEFI firmware is vendor specific, but the information is presented to the OS using an ACPI table named NBFT (NVMe Boot Firmware Table) when OS boot is taking place.

A.2.7 Boot Process

The boot process is separated as two parts: pre-OS boot and OS boot.

A.2.7.1 Pre-OS Boot

Continuing with an UEFI based example of a pre-OS environment; The seven phases in a UEFI boot sequence (refer to Figure 26 and the Tianocore: EDK2 Build Specification) are:

1. Security (SEC)
2. Pre-EFI Initialization (PEI)
3. Drive Execution Environment (DXE)
4. Boot Device Selection (BDS)
5. Transient System Load (TSL)
6. Runtime (RT)

7. After Life (AL)

The DXE driver supporting NVMe-oF boot is loaded and executed in the DXE phase. Namespaces used as boot devices are discovered using the information stored during configuration, the driver connects to the NVMe-oF subsystems, discovers the namespaces (i.e., boot devices) required for OS boot. Note that the pre-OS environment may not require persistent connections to Discovery controllers as it not intended to fully hand over connections to an OS. These namespaces are then be selected at the BDS phase as the final boot device for OS boot.

In the TSL phase, the operating system boot loader is loaded from the boot device and the UEFI firmware hands over the execution to that specific boot loader which continues the operating system boot. By this point, the NBFT has been generated either by UEFI firmware, optional drivers, or an intermediate TSL application; and stored in the main memory which is accessed by the OS as an ACPI table after UEFI ExitBootServices has been called, ending the TSL phase.

From a NVMe-oF host to controller connection perspective, a pre-OS environment may not require or utilize persistent connections to all controllers. Some pre-OS environments are not multi-threaded and may see a high burden of maintaining multiple concurrent connections.

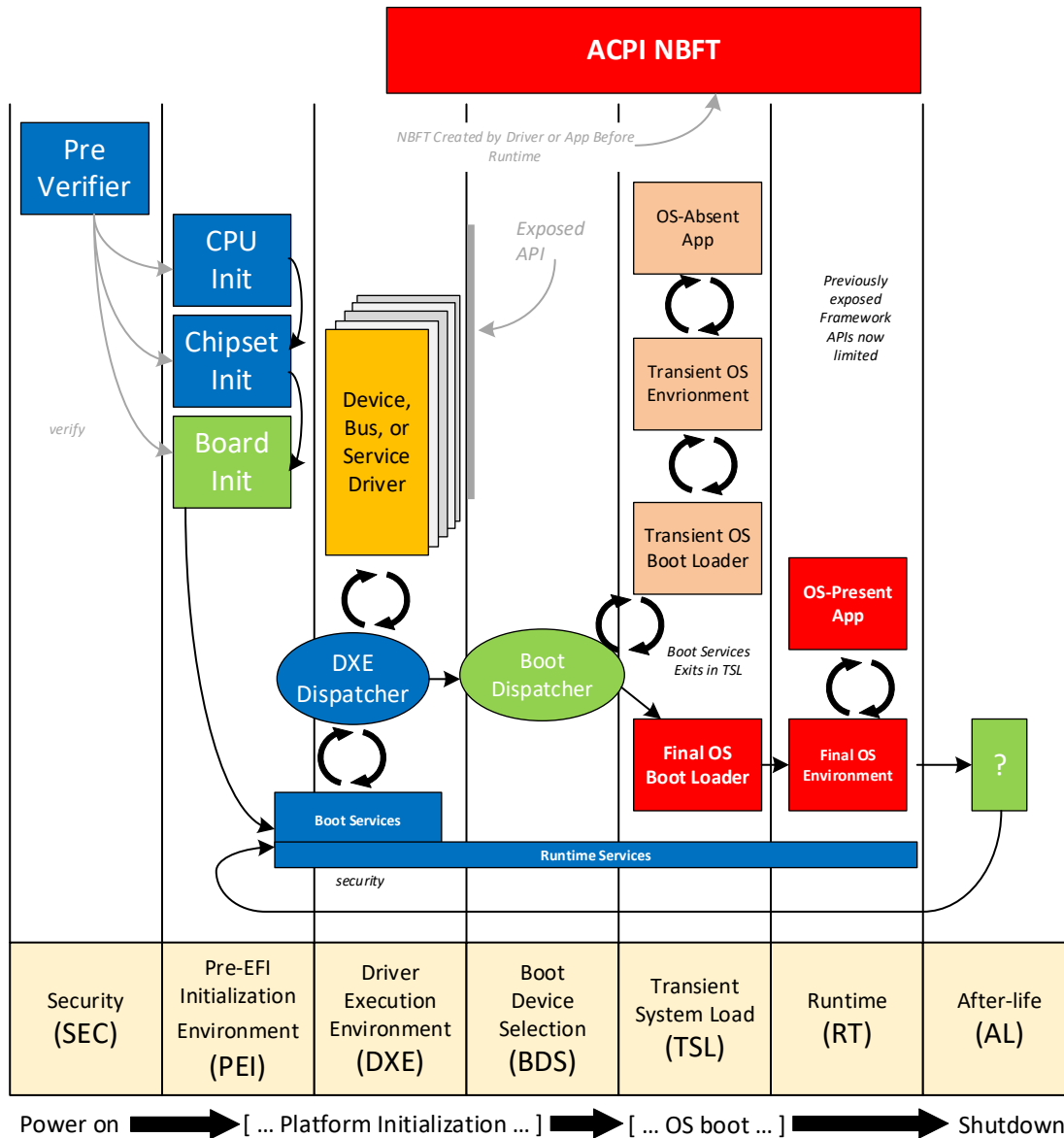
A.2.7.2 OS Boot

Typically, an operating system does not expect to rely on or reuse the pre-OS connections and connection specific data structures. Because of this, at the operating system boot stage, the OS may then use the information stored in NBFT to connect to the NVM subsystem without user intervention.

There are two scenarios in the OS boot stage: operating system installation or normal operating system boot.

1. During OS installation, user may either:
 - a) use the NBFT provided host NQN as its own host NQN; or
 - b) set a separate host NQN if the NVM subsystem has been configured to support multiple host NQNs.
2. In normal operating system boot, as the OS may have its own mechanism to persist the required information for restore the NVMe-oF connections. The OS may either continue using the NBFT or the OS specific mechanism.

Figure 26: Modern EDK2 UEFI Boot Flow



A.3 Transport Specific and interface specific semantics

This section is intended to provide informative guidance on transport and boot mechanism specific semantics. This section includes discussions of the ACPI NBFT and UEFI semantics required to adopt UEFI or UEFI and ACPI specific boot workflows. Because every platform is different, these sections provide generalizations about how to interpret the handoff between pre-OS to OS environments.

A.3.1 ACPI NBFT Discussion and Background

The NBFT contains parameters used to connect to an NVMe-oF storage device, specifically NVMe/TCP. The NBFT allows the pre-OS boot environment and OS environment to use the same parameters. This is

important when booting from an NVMe Storage Device. The NBFT parameters are typically set in the pre-boot environment or an external management tool.

The NBFT parameters can be divided into three areas: host, fabric and storage subsystem. Host parameters are used to uniquely identify the host to the storage. These parameters are typically used for access control lists and reservations. These parameters contain the Host NQN and Host ID.

Host Fabric Interface (HFI) parameters describe the underlying connection between the host and the storage. These parameters contain the transport type (e.g., Ethernet, FC, InfiniBand), transport address, keep alive settings, etc.

Storage subsystem parameters uniquely identify the NVMe Storage Device. These parameters contain the storage NQN, controller ID, NID/Namespace ID. There may be several NVMe Storage Devices connected to the fabric. The storage subsystem parameters are used to pick a specific storage device.

The minimal set of NBFT parameters required to connect to an NVMe namespace are listed below.

A.3.1.1 Required and Information for a NVMe-oF Host

Under NVMe/TCP the following fields are required:

Figure 27: Required NVMe/TCP host NBFT information with example data

Descriptor	Field	Example
Host Descriptor	Host ID	00h 11h 22h 33h 44h 55h 66h 77h 88h 99h aah bbh cch ddh eeh ffh
Host Descriptor	Host NQN	"nqn.2014-08.com.example:nvme.host.sys.xyz"
HFI Descriptor	Transport type	03h
HFI Descriptor	Host Transport Address	192.168.1.1
SSNS	Subsystem Transport Address	192.168.1.2
SSNS	Transport Service ID	"4420"
SSNS	Subsystem NQN	"nqn.2014-08.com.example:nvme.storage.xyz"
SSNS	Controller ID	FFFFh
SSNS	Namespace Identifier Type (NIDT)	03h
SSNS	Namespace Identifier (NID)	00h 11h 22h 33h 44h 55h 16h 77h 88h 99h aah bbh cch ddh eeh ffh

A.3.1.2 Invalid HostID Values

Administrative care should be taken when choosing HostID values due to interoperability concerns with host pre-OS and OS driver versions which may be separately compliant with any revision of the NVM Express Base Specification. Host OS drivers should pay attention to the flag for primary administrative identity as defined in the Primary Administrative Host Descriptor field of host flags in Figure 9 to determine if an administratively assigned value or a default value assumed by pre-OS.

A.3.1.3 Multiple NBFTs

A.3.1.3.1 Coexistence of multiple NBFT instances

NVMe-oF boot requires that the pre-OS environment creates one NBFT ACPI table per driver, which contains information on network devices that have successfully established a connection to NVMe-oF controllers and information on the connections and NVMe-oF controllers reached. Those devices may provide different level of support for NVMe-oF network stack and may be managed by different SW drivers. A situation is possible where those drivers create two unrelated NVMe-oF stacks that co-exist on a single platform, each providing a piece of data necessary to build a complete NBFT. The design choice to allow

for the coexistence of multiple NBFTs is due to the complexity for disparate ecosystem components to share such a structure in the pre-OS environment. Shared use of a single NBFT requires considerations such as:

- simultaneous read/write access to NBFT from NVMe-oF stacks;
- support for different versions of NBFT substructures by NVMe-oF stacks; and
- UEFI context: exchange and synchronization of the TableKey parameter used for adding and removing ACPI tables via EFI_ACPI_TABLE_PROTOCOL.

To prevent issues such as these from occurring, each NVMe-oF stack implementation is advised to build and install its own NBFT table. Operating systems additionally should consider analysis of multiple NBFTs and extract necessary information from them.

For infrastructure or configuration specifics that may lead multiple drivers to establish connection to the same namespace over the same HFI (e.g., two NVMe/TCP drivers), it is possible for two or more NBFT instances to contain identical entries in NBFT substructures. While such redundant information may be used or ignored by OS, such handling is out of scope for this document.

A.3.1.3.2 OS and Post Boot handling of multiple NBFTs

If multiple NBFT tables are present, only one NBFT should be administratively selected. There is no enforcement mechanism for this to be coordinated between multiple NBFT tables, however, it is not recommended to be selected in more than one NBFT. Because an OS may encounter multiple NBFTs that assert primary administrative identity as defined in the Primary Administrative Host Descriptor field of host flags in Figure 9, the OS may correlate the Boot partition to a NBFT in order to adopt a Host NQN and HostID for its own purposes. Usually this is handled by comparing the UEFI boot variable 'BootCurrent' with the NBFT it is associated with, and then use the Host NQN and HostID from the host descriptor.

A.3.1.3.3 Multiple Drivers with Different Versions creating SSNS or HFI structures

Due to the complexities of a pre-OS environment that results from multiple pre-OS drivers of different levels to simultaneously exist, the possibility exists that an OS environment may encounter scenarios with multiple NBFT structures referencing the same HFI devices or SSNS endpoint. Therefore, an OS environment that interacts with multiple NBFT's requires careful handling of potentially common or conflicting entries.

A.3.1.3.4 NBFT Header Driver Signature

The Header NBFT Driver Signature provides a mechanic for an OS driver to understand the identity of the driver that filled in a specific NBFT. These device path signatures are string representations of the issuing drivers device path.

Examples of such are:

NvmExpressDxe driver integrated into UEFI FW image on Flash:

"Fv(7CB8BDC9-F8EB-4F34-AAEA-3EE4AF6516A1)/FvFile(5BE3BDF4-53CF-46A3-A6A9-73C34A6E5EE3)"

Driver coming from OptionROM:

"PciRoot(0x0)/Pci(0x10,0x0)/Pci(0x0,0x1)/Offset(0x13800,0x31DFF)"

A.3.1.3.5 Impact of NBFT updates

In a pre-OS environment, as a driver is assembling configuration data and context for the creation of a NBFT, the driver may choose to build the structures of the NBFT in stages. As there are cross references, namely between index values for Security Profile Descriptors or Discovery Descriptors, care should be taken to ensure any creation, deletion or update of an index is done atomically.

As an example, pre-generating SSNS records before discovery is complete may result in a namespace that is associated with namespaces discovered across multiple HFIs. To handle this, a driver may determine it is required to update an existing SSNS heap object to amend a new HFI. If the driver has pre-allocated the SSNS heap object, this may result in fragmentation of the heap. An approach for drivers that want to avoid fragmentation may consider delaying allocation of the heap object until all HFI and subsequent controller discovery is completed and then come back and re-fill indices.

A.3.2 Implementation considerations

This section contains high level overview of common use cases or problems in consistent implementation and outlines steps to answer specific questions.

A.3.2.1 Host Connection Information in the absence of a Discovery Controller

If a user specifies Discovery controller(s) as the connection target, the Discovery Controller Structure should contain the Discovery controller(s) data input from the HII interface.

If a user selects the auto discovery option in the HII interface, when implemented in the HII interface, the Discovery Controller Structure should contain the Discovery controller(s) discovered.

However, if a user specifies an NVM subsystem as direct connect in the HII interface, the value in Control Structure Discovery Descriptor Offset field is required to be cleared to 0h.

If a host implementation supports a user configuration interface (i.e., UEFI HII), to allow administrator's input of NVM subsystem information (i.e., NVM subsystem qualified name, transport service identifier, port id, controller id, address family, transport type, transport address, maximum size of an admin submission queue), a host is not required to query Discovery Log Pages from the Discovery controller to obtain NVM subsystem information. For this implementation, a Discovery Descriptor is not present with

- the Discovery Offset in the Control Structure cleared to 0h; and
- the Primary Discovery Ctrl Offset field of the Subsystem and Namespace Descriptor cleared to 0h,

to indicate no Discovery controller is utilized. An example of this is shown in Figure 28.

Figure 28: Example Subsystem and Namespace descriptor

Bytes	Description	Example
03:00	Structure ID: 4 = Namespace	4h
05:04	Index: Index = 1 indicates namespace 1 ... Index = n indicates namespace n	
07:06	SSNS Flags: 1 = Namespace Structure Valid	1h
08	Transport Type: 1 = NVMe/TCP	1h
09	Primary Discovery Ctrl Offset: 0 = No Discovery Descriptor reference.	0h
127:10	Reserved	

A.3.2.2 Redfish Keypath Object

The Security Profile Descriptor Keypath URI Object (refer to Figure 21) allows the NBFT to reference the relative URI of a well-formed Redfish™ Key Collection Object, defined in DMTF DSP8010 or later. It is assumed that the base-address of the Redfish URI is obtained through available methods to the platform or OS environment, such as DMTF DSP0270, and is outside the scope of this document.

The following is an example of such a formatted Redfish relative URI:

```
/redfish/v1/KeyService/NVMeoFSecrets/0
```

Where '0' in the Redfish® URI is the Id value of an NVMeoF Key Service Secrets Collection.

The object referenced by this Keypath includes the following important pieces of information:

Id, KeyString, NQN, and HostKeyld

Where:

- Id is again the value of the current record;
- KeyString is the key material for the belonging to Id;
- NQN is the host or Subsystem NQN depending on both HostKeyld value being present and the NQN. If the current record is not a Host NQN, Additionally the SubsystemNQN of the related Discovery or SSNS record should match this value; and
- HostKeyld is the Id value of the record that corresponds to the Host NQN associated with the current record. Its presence indicates the current record is not a Host NQN, but the record should be addressable by a relative local URI with the value of HostKeyld as the Id.

For detailed examples, please see DMTF DSP8010 mockups of NVMeoFSecrets objects.

A.3.2.3 NID decoding by drivers or environments

In some boot environments, it may be possible for administrative and configuration interfaces to digest either binary or textual representations of a NID. Without outside context an undecorated NID of type 2 (i.e., NGUID) or type 3 (i.e., UUID) may be indistinguishable for many valid values of both NID types. Because of this, a driver or environment should avoid explicitly attempting to determine if an undecorated value is one or the other without providing positive selection by a user to ensure proper NID Type attribution.

For example, the following undecorated value is both a valid NGUID and a valid UUID:

- The string "4eff7f8ed3534e9ba4ecdeea8eab84d7"

However, the following decorated string is well understood to be a UUID type NID:

- The string "urn:uuid:4eff7f8e-d353-4e9b-a4ec-deea8eab84d7"

When at all possible, it is recommended to include mechanisms to select NID type in administrative interfaces when decoration is unable to be enforced.

A.4 OS Specific Details and Application Considerations

This OS specific section details how individual OS handle namespace information handed over from pre-boot environment (e.g., virtual machines).

A.4.1 Generic OS Usage of the NBFT

A.4.1.1 Accessing NBFT in operating system

During system boot process, when the execution is handed over to the operating system, all the information, with the exception of security credentials, required for NVMe-oF boot has been presented in the form of a NBFT as an ACPI table, hence the NBFT may be accessed by the OS using OS provided ACPI API during OS boot. Security credentials may be accessed via the Security Descriptors Secret Keypath Offset Heap Object Reference (refer to Figure 21).

For further accessing of a NBFT, an OS exports an NBFT in the OS specific method (e.g., sys interface in Linux, VSI interface in VMware ESXi). However, if an NBFT is accessed as a generic ACPI table using an ACPI API, the NBFT content may be just a stream of bytes without any specific meaning. To make NBFT access more meaningful from an NVMe-oF boot point of view, the OS may implement an NBFT specific interface (e.g., API, UI, tools) to make NBFT access more convenient and user-friendly. Subsequently, other components in the OS may access NBFT without using ACPI API as the ACPI API may not be generally available to all the components.

A.4.1.2 Consuming NBFT

The information contained in an NBFT may be consumed for multiple uses, including:

- using the connect parameters to connect to NVMe subsystem on which the OS is installed, including security;
- using the connect parameters to connect/discover other NVMe subsystems; or
- using the parameters to configure the OS itself (e.g., host NQN, host ID).

The NBFT may be consumed at different stages of the OS lifecycle, including, but not limited to:

- a. during system boot: The OS is able to detect if there is NBFT available and use the parameters in NBFT to connect to the NVMe-oF subsystem and present NVMe namespaces to the OS for further use; or
- b. during OS installation: The OS installer is able to detect the presence of an NBFT and determine that there is at least one NVMe namespace present in the system that supports NVMe-oF boot. The administrator may then choose installing the OS on an NVMe namespace that supports NVMe-oF boot.

The OS should assume the NBFT is recreated from pre-OS context at every boot and that the NBFT is not guaranteed to remain static. Pre-OS context may change between reboots due to platform configuration or administrative changes. Therefore the OS should be tolerant of changes between reboots.

A.4.2 Linux

A.4.2.1 Linux Boot Loader

On entry, a Linux boot loader (e.g., shimx64.efi, grubx64.efi) may assume an association exists to the namespace from which it was loaded. But the boot loader is not required to be aware that it is booting from an NVMe-oF namespace or that an NBFT exists. That is, no NVMe-oF specific changes are required to an EFI boot loader, such as grubx64.efi. The boot loader may continue to use the EFI block device protocol to:

- use NVMe-oF associations made by the UEFI firmware;
- access the ESP;

- find any other partitions;
- locate systems files that contain boot loader configuration;
- locate and load the Linux kernel; and
- load the initrd/initramfs files.

The files that are used by a Linux EFI boot loader are typically located in “/boot” or the ESP. If “/boot” is located in a different NVMe-oF namespace than the ESP, that namespace must have been previously administratively configured such that required association is created during the UEFI DXE phase.

Note that there is no EFI protocol or mechanism by which an EFI boot loader creates additional NVMe-oF associations for the current boot.

A.4.2.2 Linux Kernel Boot

The EFI ExitBootServices() call is the transition from the TSL phase (e.g., EFI grub2) to the OS environment. Typically, this call is made very early in the boot flow of the Linux kernel. On return from the EFI ExitBootServices() call, any NVMe-oF associations created during the UEFI pre-OS environment are closed. Therefore, if continuation of the OS boot requires access to any NVMe-oF namespaces, the booting Linux system must recreate these associations using the Linux kernel drivers and utilities. The information in the NBFTs provided by the UEFI pre-OS environment may be used by the Linux OS to create the required boot-time NVMe-oF associations (e.g., “/”, “/boot”) that are required to continue the Linux OS boot.

Note that the presence of the NBFT does not preclude the use of OS-specific configuration files to recreate or to create new associations to any additional NVMe-oF namespaces that the OS requires. If those configuration files are somehow passed as kernel boot parameters or are located in “/boot”.

Figure 29 summarizes common Linux pre-OS associations that might utilize NBFT supported transports for NVMe-oF namespaces and possible expectations by native Linux OS boot software. When such a pre-OS association is created, parameters used to form that association are present in the NBFT.

Figure 29: Common Linux boot namespaces and expectations on NBFT content

Use of NVMe-oF Namespaces	Use of pre-OS association	Use of parameters in NBFT
ESP	Required by the boot loader ¹	May be used by OS
“/boot”	Required by the boot loader ²	Typically used by OS
“/”	Optional for the boot loader ²	Typically used by OS
other file system mounts, if any	OS configuration dependent ²	OS configuration dependent
Notes:		
1. Established and required by UEFI firmware.		
2. If different namespace than ESP, must be administratively configured.		

A.4.2.3 NBFT access and availability

Linux makes ACPI tables available to userspace as binary files via sysfs. Since the NBFT is a standard ACPI table, if present, the sysfs representation of the NBFT is named:

```
/sys/firmware/acpi/tables/NBFT
```

If multiple NBFTs are present, then the first one has the instance number “1” appended to its name. Any additional NBFTs have ascending instance numbers, e.g.,:

```
/sys/firmware/acpi/tables/NBFT1
```

```
/sys/firmware/acpi/tables/NBFT2
```

/sys/firmware/acpi/tables/NBFT3

...

Software that reads the NBFT should check for and handle the possibility that multiple ACPI NBFT tables are present.

A.4.2.4 Requirements for EFI boot variables for NVMe-oF namespaces

The creation or modification of an EFI boot variable with a value that includes an NVMe-oF device path for an NVMe-oF namespace does not automatically trigger the update the NBFT for the next boot.

Such a boot variable is usable on the next boot, if and only if either:

1. Corresponding information about that NVMe-oF namespace is already in the NBFT; or
2. Prior to the next boot, corresponding information about that NVMe-oF namespace is administratively configured (i.e., via a RedFish API, or manually configured via EFI HII).

This has implications for an OS installer which boots from a non NVMe-oF device, e.g., network boot, and is installing the OS to an NVMe-oF namespace that was configured using OS installer provided inputs or OS-based configuration files. A Linux OS installer typically uses the efibootmgr utility to create the required EFI boot variable for the new OS installation.

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.
<i>Green Bracketed:</i>	Notes to editor

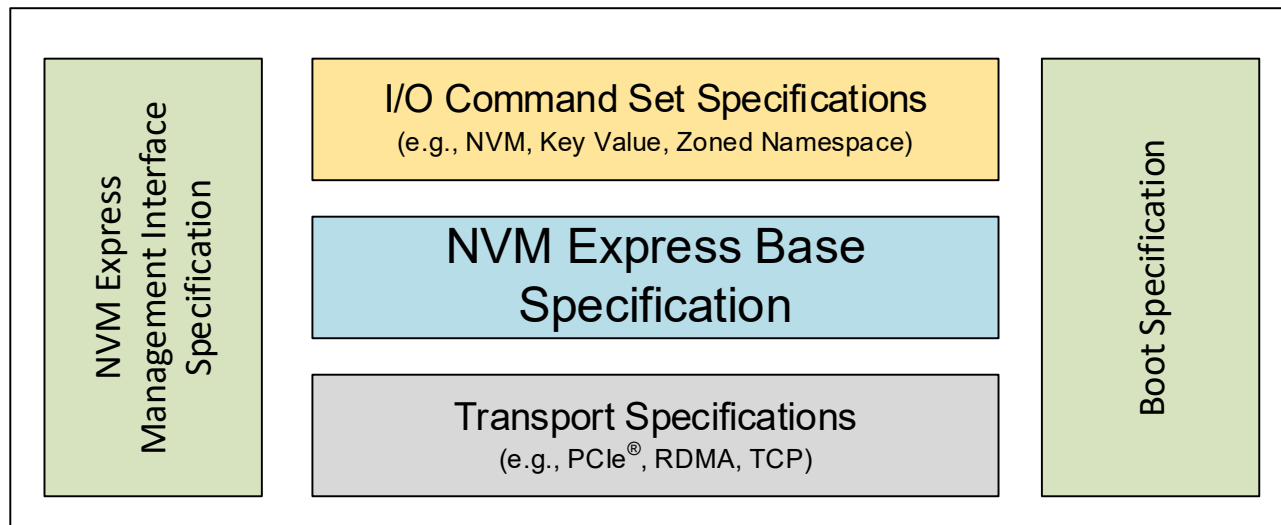
Description for Changes Document for NVM Express Base Specification 2.0b

Editor Update to Section 1.1.1 (NVM Express Specification Family)

1.1 Overview

1.1.1 NVM Express® Specification Family

Change Figure 1 to:



Add the following sentence to the end of section 1.1.1(NVM Express Specification Family)

[...]

...

The NVM Express Boot Specification defines constructs and guidelines for booting from NVM Express interfaces.

1.1.8 References

Editor Add to Section 1.8 (References)

[...]

...

NVM Express Boot Specification, Revision 1.0. Available from <https://www.nvmexpress.org>.

Defines new reference to NVMe Boot Specification regarding textual representations related to Boot of NIDs, References NVMe-Base Specification 2.0B section <5.17.2.3>

Description of Specification Changes

Modify section 5.17.2.3 as modified by ECN110 to add reference Textual representations related to Boot of NIDs in the NVMe Boot Spec

5.17.2.3 Namespace Identification Descriptor list (CNS 03h)

A list of Namespace Identification Descriptor structures (refer to Figure 277) is returned to the host for the namespace specified in the Namespace Identifier (NSID) field if it is an active NSID. Namespace Identification Descriptor structures consist of one or more Namespace Identifiers (NID) of various types as indicated by the Namespace Identifier Type (NIDT) field in each descriptor. Each NID is assigned to a namespace at namespace creation and remains fixed throughout the life of that namespace. If the NSID field does not specify an active NSID, then refer to section 3.2.1.5 for the status code to return. [Textual representations related to Boot of NIDs may be found in the NVM Express Boot Specification.](#)

The contents of the Namespace Identification Descriptor list is preserved across namespace and controller operations (e.g., Controller Level Reset, namespace format, etc.).

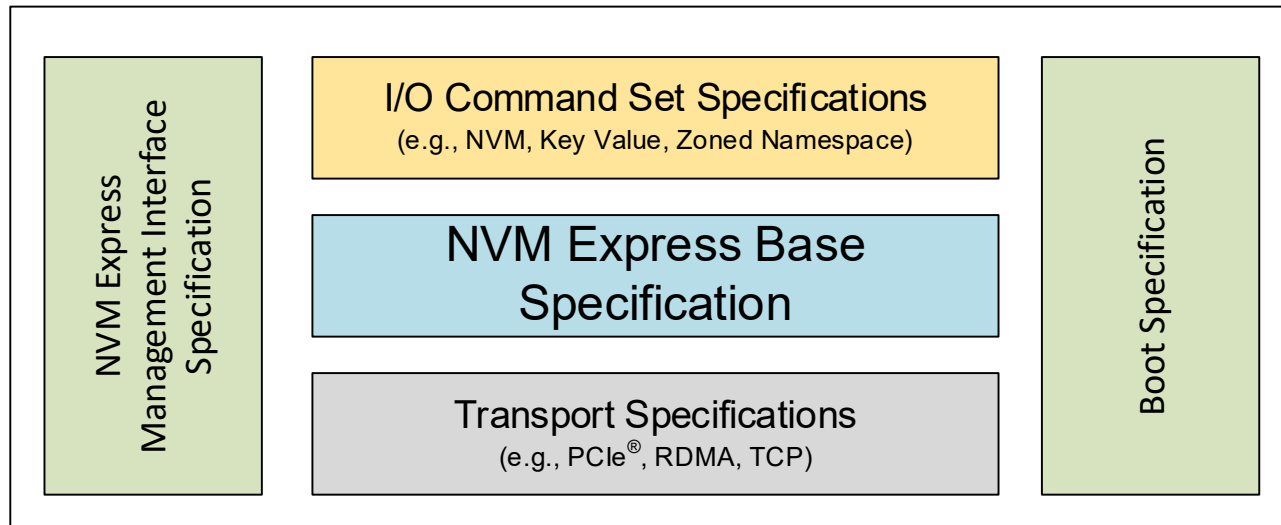
Description for Changes Document for NVM Express NVM Command Set Specification 1.0b

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope) list.

[...]

...

6. NVM Express® Boot Specification

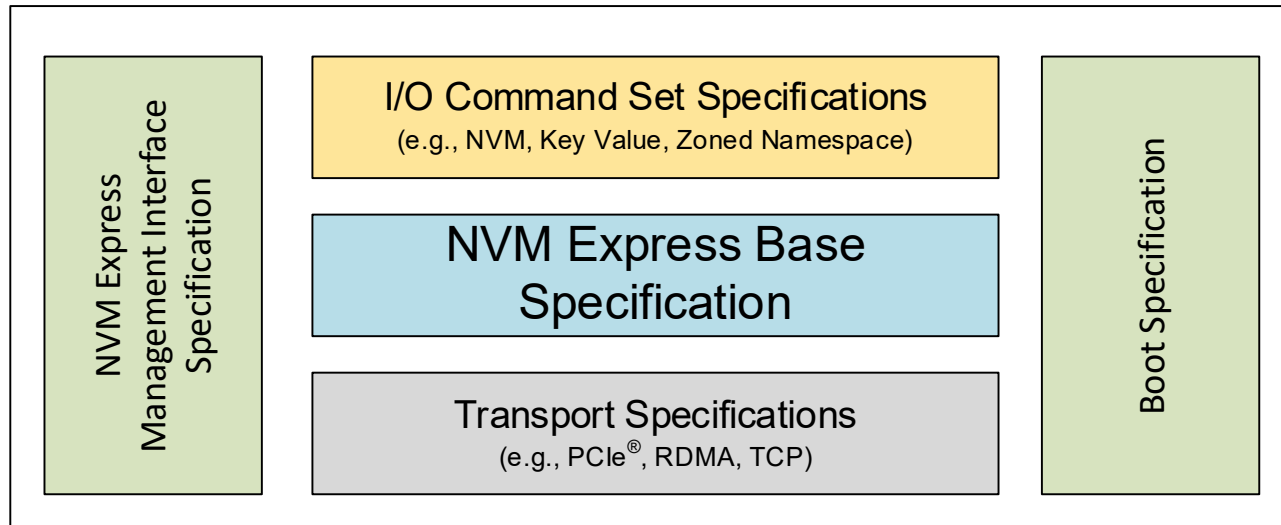
Description for Changes Document for NVMe over PCIe Transport Specification 1.0b

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope)list.

[...]

...

6. NVMe Express® Boot Specification

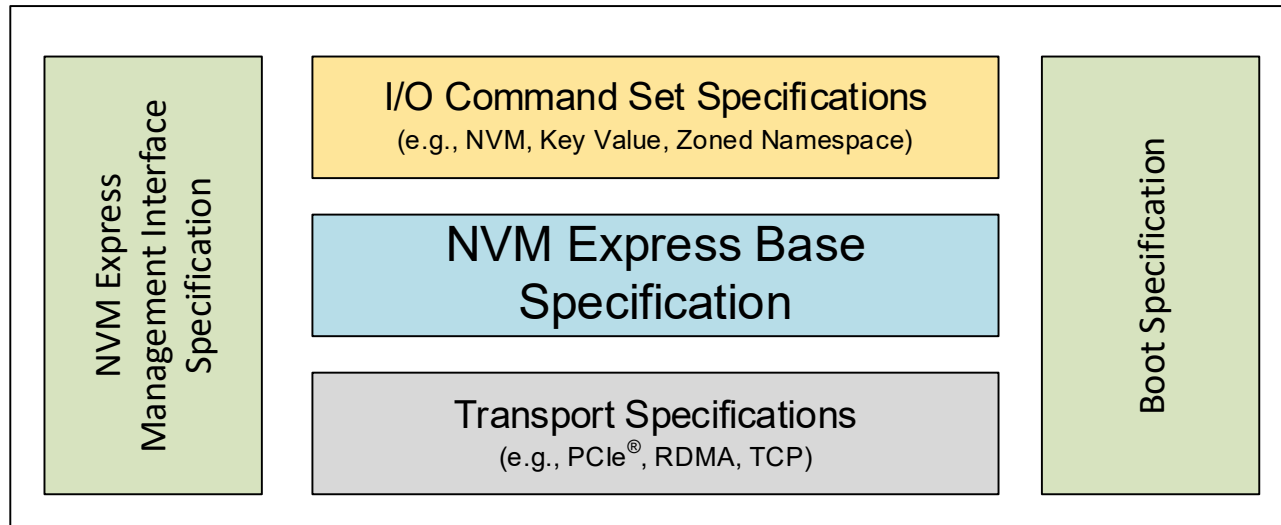
Description for Changes Document for NVM Express Management Interface Specification 1.2b

Editor Update to Section 1.1 (Overview)

1.1 Overview

...

Change Figure 1 to:



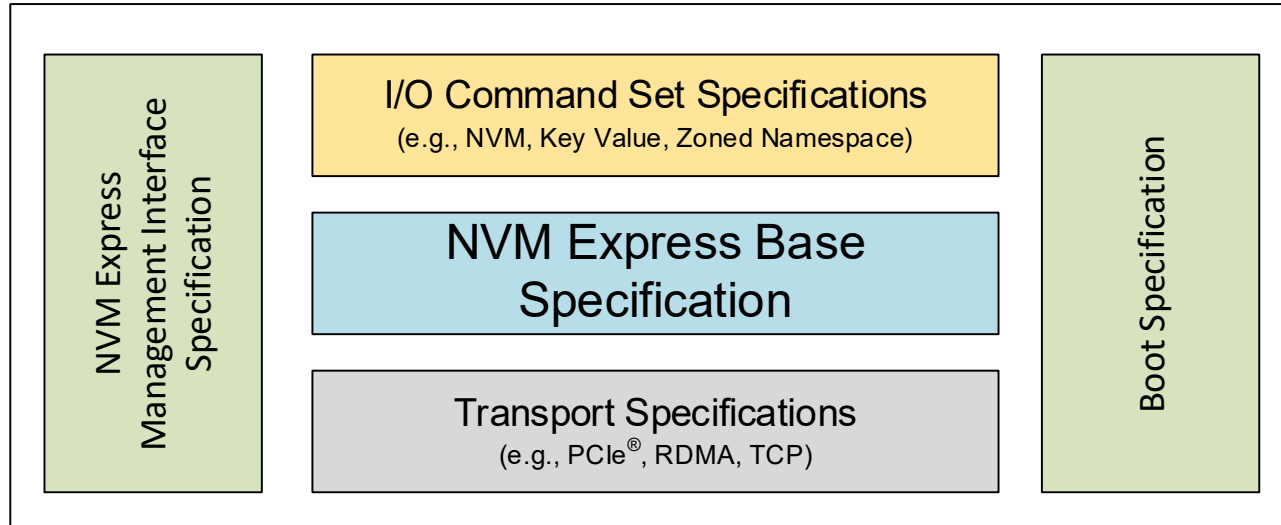
Description for Changes Document for NVM Express RDMA Transport Specification 1.0a

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope) list.

[...]

...

6. NVM Express® Boot Specification

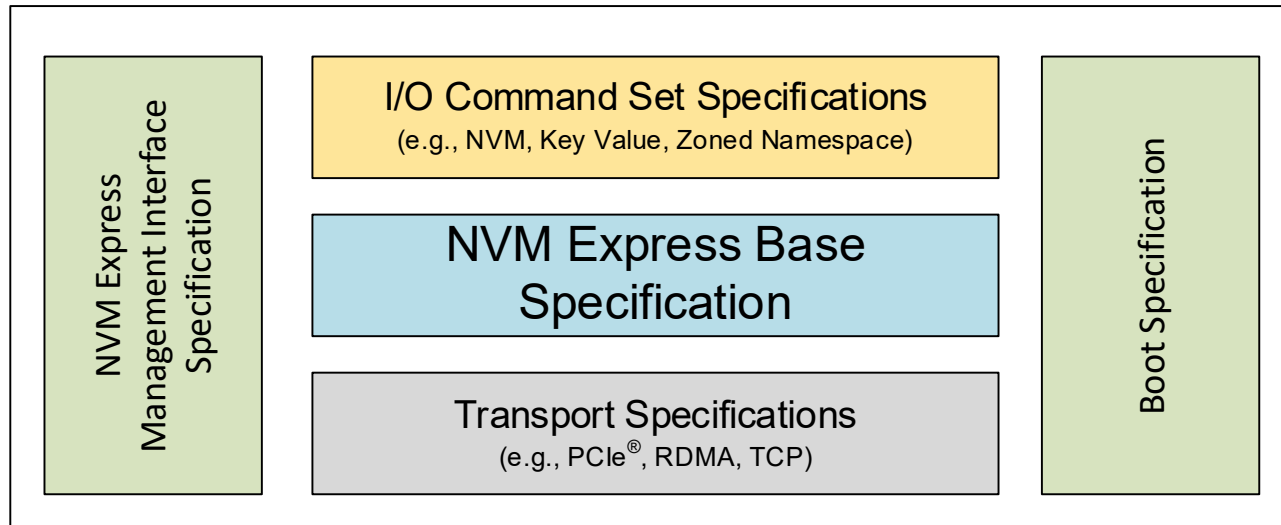
Description for Changes Document for NVM Express TCP Transport Specification 1.0b

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope) list.

[...]

...

6. NVM Express® Boot Specification

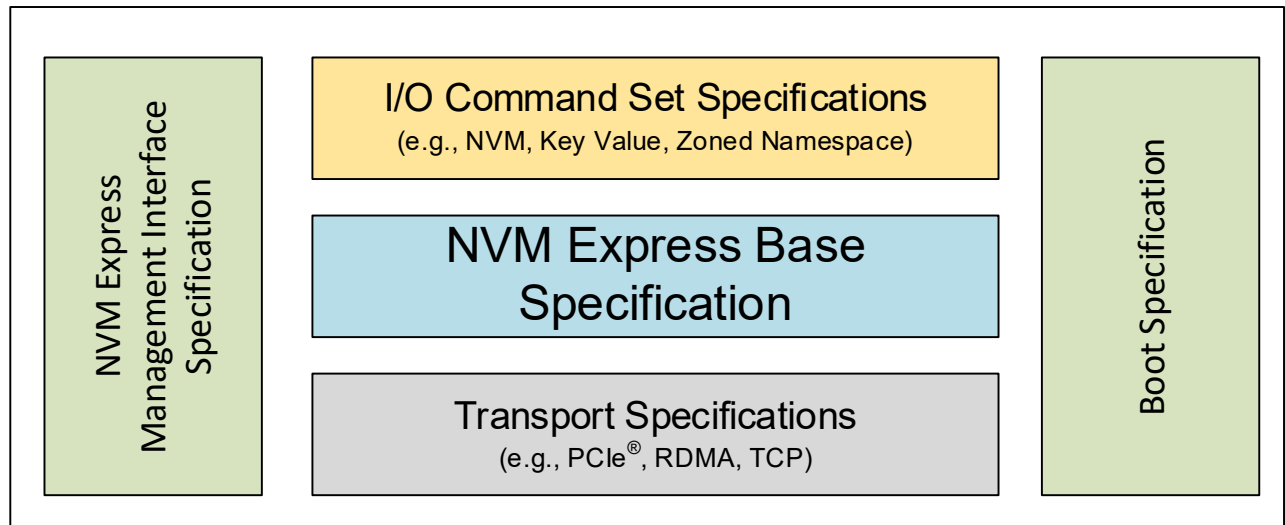
Description for Changes Document for NVM Express Zoned Namespace Command Set Specification 1.1b

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope) list.

[...]

...

6. NVM Express® Boot Specification

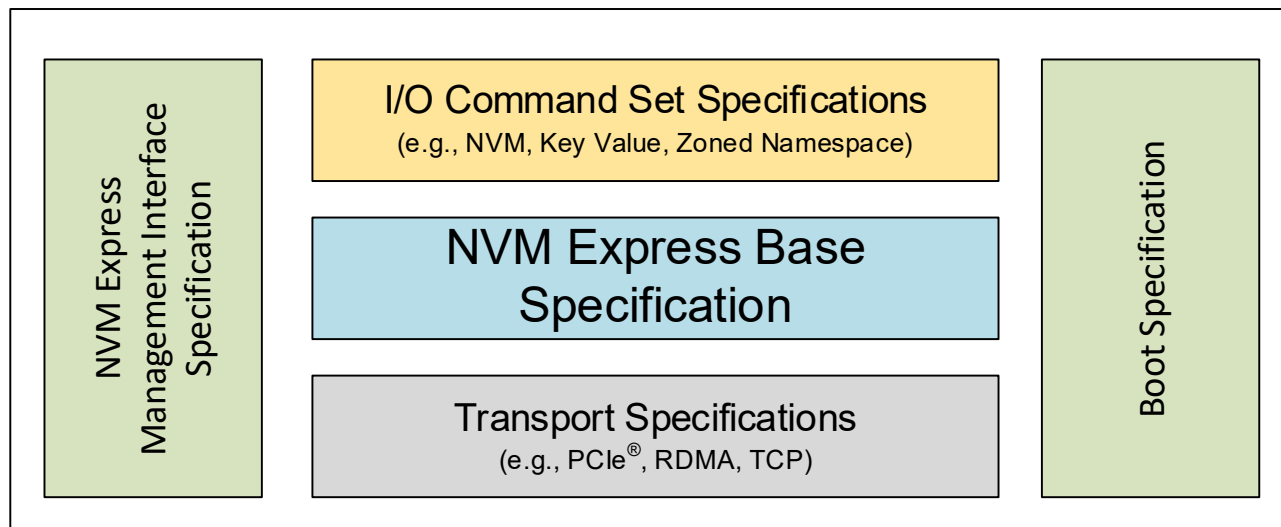
Description for Changes Document for NVM Express Key Value Command Set Specification 1.0b

Editor Update to Section 1.2 (Scope)

1.2 Scope

...

Change Figure 1 to:



Add the following Numbered Bullet to the end of Section 1.2 (Scope) list.

[...]

...

6. NVM Express® Boot Specification