



1
2
3
4

Document Identifier: DSP0238

Date: 2021-03-02

Version: 1.2.0

5
6
7

Management Component Transport Protocol (MCTP) PCIe VDM Transport Binding Specification

8
9
10
11

Supersedes: 1.1.0

Document Class: Normative

Document Status: Published

Document Language: en-US

12

13 Copyright Notice

14 Copyright © 2009, 2018, 2021 Distributed Management Task Force, Inc. (DMTF). All rights reserved.

15 DMTF is a not-for-profit association of industry members dedicated to promoting enterprise and systems
16 management and interoperability. Members and non-members may reproduce DMTF specifications and
17 documents, provided that correct attribution is given. As DMTF specifications may be revised from time to
18 time, the particular version and release date should always be noted.

19 Implementation of certain elements of this standard or proposed standard may be subject to third party
20 patent rights, including provisional patent rights (herein "patent rights"). DMTF makes no representations
21 to users of the standard as to the existence of such rights, and is not responsible to recognize, disclose,
22 or identify any or all such third party patent right, owners or claimants, nor for any incomplete or
23 inaccurate identification or disclosure of such rights, owners or claimants. DMTF shall have no liability to
24 any party, in any manner or circumstance, under any legal theory whatsoever, for failure to recognize,
25 disclose, or identify any such third party patent rights, or for such party's reliance on the standard or
26 incorporation thereof in its product, protocols or testing procedures. DMTF shall have no liability to any
27 party implementing such standard, whether such implementation is foreseeable or not, nor to any patent
28 owner or claimant, and shall have no liability or responsibility for costs or losses incurred if a standard is
29 withdrawn or modified after publication, and shall be indemnified and held harmless by any party
30 implementing the standard from any and all claims of infringement by a patent owner for such
31 implementations.

32 For information about patents held by third-parties which have notified the DMTF that, in their opinion,
33 such patent may relate to or impact implementations of DMTF standards, visit
34 <http://www.dmtf.org/about/policies/disclosures.php>.

35 PCI-SIG, PCIe, and the PCI HOT PLUG design mark are registered trademarks or service marks of PCI-
36 SIG.

37 All other marks and brands are the property of their respective owners.

38 This document's normative language is English. Translation into other languages is permitted.

39

40

CONTENTS

41 Foreword 4

42 Introduction..... 5

43 1 Scope 7

44 2 Normative references 7

45 3 Terms and definitions 8

46 4 Symbols and abbreviated terms..... 8

47 5 Conventions 9

48 5.1 Reserved and unassigned values..... 9

49 5.2 Byte ordering..... 9

50 6 MCTP over PCI Express VDM transport..... 9

51 6.1 Packet format..... 9

52 6.2 Supported media..... 12

53 6.3 Physical address format for MCTP control messages..... 12

54 6.4 Message routing 13

55 6.5 Bus owner address 13

56 6.6 Bus address assignment for PCIe 14

57 6.7 Host dependencies 14

58 6.8 Discovery Notify message use for PCIe 14

59 6.9 MCTP over PCIe endpoint discovery..... 15

60 6.10 MCTP messages timing requirements..... 19

61 ANNEX A (informative) Notations and conventions..... 21

62 ANNEX B (informative) Change log..... 22

63

64 Figures

65 Figure 1 – MCTP over PCI Express Vendor Defined Message (VDM) packet format 10

66 Figure 2 – Flow of operations for full MCTP Discovery over PCIe 17

67 Figure 3 – Flow of operations for Partial Endpoint Discovery..... 18

68

69 Tables

70 Table 1 – PCI Express medium-specific MCTP packet fields..... 10

71 Table 2 – Supported media..... 12

72 Table 3 – Physical address format..... 12

73 Table 4 – Timing specifications for MCTP Control messages on PCIe VDM..... 19

74

75

Foreword

76 The *Management Component Transport Protocol (MCTP) PCIe VDM Transport Binding Specification*
77 (DSP0238) was prepared by the PMCI Working Group.

78 DMTF is a not-for-profit association of industry members dedicated to promoting enterprise and systems
79 management and interoperability.

80 **Acknowledgments**

81 The DMTF acknowledges the following individuals for their contributions to this document:

82 **Editors:**

- 83 • Hemal Shah – Broadcom Inc.
- 84 • Tom Slaight – Intel Corporation

85 **Contributors:**

- 86 • Patrick Caporale – Lenovo
- 87 • Yuval Itkin – NVIDIA Corporation
- 88 • Eliel Louzoun – Intel Corporation
- 89 • Patrick Schoeller – Hewlett Packard Enterprise
- 90 • Bob Stevens – Dell Technologies

91

Introduction

92 The Management Component Transport Protocol (MCTP) over PCIe VDM transport binding defines a
93 transport binding for facilitating communication between platform management subsystem components
94 (e.g., management controllers, management devices) over PCIe.

95 The [MCTP Base Specification](#) describes the protocol and commands used for communication within and
96 initialization of an MCTP network. The MCTP over PCIe VDM transport binding definition in this
97 specification includes a packet format, physical address format, message routing, and discovery
98 mechanisms for MCTP over PCIe VDM communications.
99

101 Management Component Transport Protocol (MCTP) PCIe 102 VDM Transport Binding Specification

103 1 Scope

104 This document provides the specifications for the Management Component Transport Protocol (MCTP)
105 transport binding using PCIe Vendor Defined Messages (VDMs).

106 2 Normative references

107 The following referenced documents are indispensable for the application of this document. For dated
108 references, only the edition cited applies. For undated references, the latest edition of the referenced
109 document (including any amendments) applies.

110 CXL Consortium, *Compute Express Link™ (CXL™) Specification Revision 1.0*,
111 <https://www.computeexpresslink.org>

112 CXL Consortium, *Compute Express Link™ (CXL™) Specification Revision 1.1*,
113 <https://www.computeexpresslink.org>

114 CXL Consortium, *Compute Express Link™ (CXL™) Specification Revision 2.0*,
115 <https://www.computeexpresslink.org>

116 DMTF DSP0236, *Management Component Transport Protocol (MCTP) Base Specification 1.0*
117 https://www.dmtf.org/sites/default/files/standards/documents/DSP0236_1.0.pdf

118 DMTF DSP0236, *Management Component Transport Protocol (MCTP) Base Specification 1.3*
119 https://www.dmtf.org/sites/default/files/standards/documents/DSP0236_1.3.pdf

120 DMTF DSP0239, *Management Component Transport Protocol (MCTP) IDs and Codes 1.0*
121 https://www.dmtf.org/sites/default/files/standards/documents/DSP0239_1.0.pdf

122 DMTF DSP0239, *Management Component Transport Protocol (MCTP) IDs and Codes 1.8*
123 https://www.dmtf.org/sites/default/files/standards/documents/DSP0239_1.8.pdf

124 ISO/IEC Directives, Part 2, *Rules for the structure and drafting of International Standards*,
125 <http://isotc.iso.org/livelink/livelink?func=ll&objId=4230456&objAction=browse&sort=subtype>

126 PCI-SIG, *PCI Express® Base Specification Revision 1.1*, March 8, 2005,
127 <http://www.pcisig.com/specifications/>

128 PCI-SIG, *PCI Express® Base Specification Revision 2.0*, December 20, 2006,
129 <http://www.pcisig.com/specifications/>

130 PCI-SIG, *PCI Express® Base Specification Revision 2.1*, March 4, 2009,
131 <http://www.pcisig.com/specifications/>

132 PCI-SIG, *PCI Express® Base Specification Revision 3.0*, November 10, 2010,
133 <http://www.pcisig.com/specifications/>

134 PCI-SIG, *PCI Express® Base Specification Revision 3.1a*, December 7, 2015,
135 <http://www.pcisig.com/specifications/>

136 PCI-SIG, *PCI Express® Base Specification Revision 4.0*, October 5, 2017,
137 <http://www.pcisig.com/specifications/>

138 PCI-SIG, *PCI Express® Base Specification Revision 5.0*, May 28, 2019,
139 <http://www.pcisig.com/specifications/>

140 3 Terms and definitions

141 In this document, some terms have a specific meaning beyond the normal English meaning. Those terms
142 are defined in this clause.

143 The terms "shall" ("required"), "shall not", "should" ("recommended"), "should not" ("not recommended"),
144 "may", "need not" ("not required"), "can" and "cannot" in this document are to be interpreted as described
145 in [ISO/IEC Directives, Part 2](#), Clause 7. The terms in parentheses are alternatives for the preceding term,
146 for use in exceptional cases when the preceding term cannot be used for linguistic reasons. Note that
147 [ISO/IEC Directives, Part 2](#), Clause 7 specifies additional alternatives. Occurrences of such additional
148 alternatives shall be interpreted in their normal English meaning.

149 The terms "clause", "subclause", "paragraph", and "annex" in this document are to be interpreted as
150 described in [ISO/IEC Directives, Part 2](#), Clause 6.

151 The terms "normative" and "informative" in this document are to be interpreted as described in [ISO/IEC](#)
152 [Directives, Part 2](#), Clause 3. In this document, clauses, subclauses, or annexes labeled "(informative)" do
153 not contain normative content. Notes and examples are always informative elements.

154 Refer to [DSP0236](#) for terms and definitions that are used across the MCTP specifications. For the
155 purposes of this document, the following additional terms and definitions apply.

156 3.1

157 MCTP PCIe Endpoint

158 a PCIe endpoint on which MCTP PCIe VDM communication is supported

159 4 Symbols and abbreviated terms

160 Refer to [DSP0236](#) for symbols and abbreviated terms that are used across the MCTP specifications. The
161 following symbols and abbreviations are used in this document.

162 4.1

163 PCIe®

164 PCI Express™

165 4.2

166 VDM

167 Vendor Defined Message

168 4.3

169 CXL™

170 Compute Express Link™

171 **5 Conventions**

172 The conventions described in the following clauses apply to this specification.

173 **5.1 Reserved and unassigned values**

174 Unless otherwise specified, any reserved, unspecified, or unassigned values in enumerations or other
175 numeric ranges are reserved for future definition by the DMTF.

176 Unless otherwise specified, numeric or bit fields that are designated as reserved shall be written as 0
177 (zero) and ignored when read.

178 **5.2 Byte ordering**

179 Unless otherwise specified, byte ordering of multi-byte numeric fields or bit fields is "Big Endian" (that is,
180 the lower byte offset holds the most significant byte, and higher offsets hold lesser significant bytes).

181 **6 MCTP over PCI Express VDM transport**

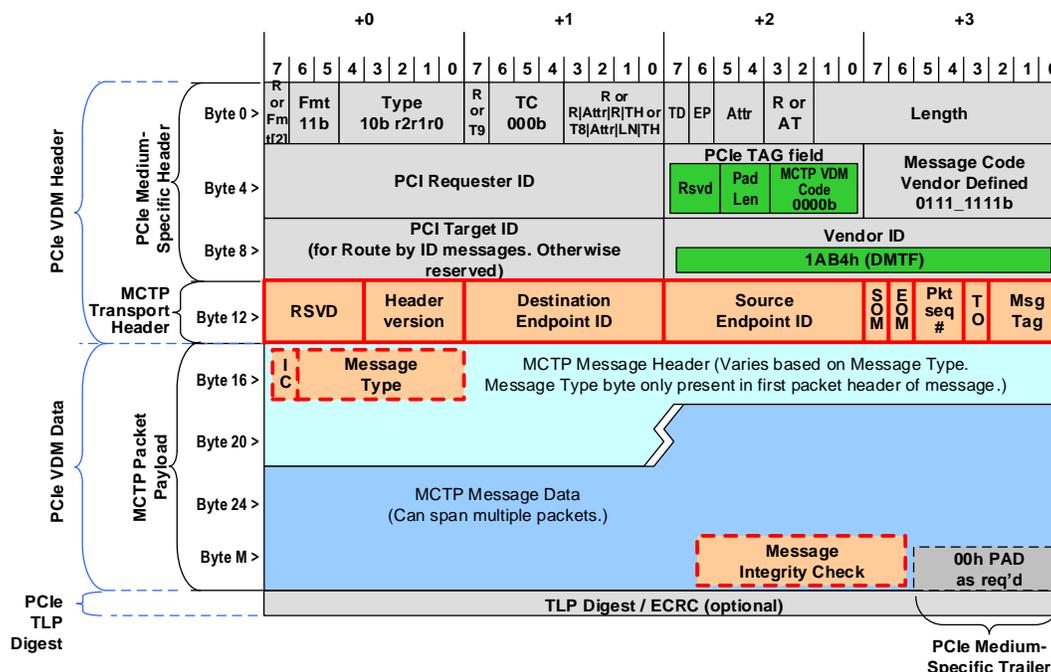
182 This document defines the medium-specific transport binding for transferring MCTP packets between
183 endpoints on PCI Express™ using PCIe Vendor Defined Messages (VDMs).

184 A MCTP over PCIe VDM compliant PCIe device shall support MCTP over PCIe VDM communications on
185 at least one PCIe Physical Function (PF) of the device. If a MCTP over PCIe VDM compliant PCI device
186 supports MCTP over PCIe VDM communications on more than one PCIe function, then MCTP over PCIe
187 VDM communication on each function shall be independent from MCTP over PCIe VDM communications
188 on other PCIe functions.

189 **6.1 Packet format**

190 The MCTP over PCI Express (PCIe) VDM transport binding transfers MCTP messages using PCIe Type
191 1 VDMs with data. MCTP messages use the MCTP VDM code value (0000b) that uniquely differentiates
192 MCTP messages from other DMTF VDMs.

193 Figure 1 shows the encapsulation of MCTP packet fields within a PCIe VDM.



194

195 **Figure 1 – MCTP over PCI Express Vendor Defined Message (VDM) packet format**

196 The fields labeled “PCIe Medium-Specific Header” and “PCIe Medium-Specific Trailer” are specific to
 197 carrying MCTP packets using PCIe VDMs. The fields labeled “MCTP Transport Header” and “MCTP
 198 Packet Payload” are common fields for all MCTP packets and messages and are specified in [MCTP](#). This
 199 document defines the location of those fields when they are carried in a PCIe VDM. The PCIe
 200 specification allows the last four bytes of the PCIe VDM header to be vendor defined. The MCTP over
 201 PCIe VDM transport binding specification uses these bytes for MCTP Transport header fields under the
 202 DMTF Vendor ID. This document also specifies the *medium-specific* use of the MCTP “Hdr Version” field.

203 Table 1 lists the PCIe medium-specific fields and field values that shall be used in MCTP over PCIe VDM
 204 communications. When not specified, field values shall be set according to PCIe specifications. Note that
 205 the presence of TLP prefixes in MCTP over PCIe VDM packets is implementation dependent and outside
 206 the scope of this specification.

207 **Table 1 – PCI Express medium-specific MCTP packet fields**

Field	Description
R or Fmt[2]	PCIe 1.1/2.0: PCIe reserved bit (1 bit). PCIe 2.1, 3.X, 4.X, 5.X: Fmt[2]. Set to 0b.
Fmt	Format (2 bits). Set to 11b to indicate 4 dword header with data.
Type	Type and Routing (5 bits). [4:3] Set to 10b to indicate a message [2:0] PCI message routing (r2r1r0) 000b : Route to Root Complex 010b : Route by ID 011b : Broadcast from Root Complex Other routing fields values are not supported for MCTP.

Field	Description
R or T9	PCIe 1.1/2.0/2.1/3.X: PCIe reserved bit (1 bit). Refer to the PCI Express™ specification (PCIe). Set to 0b. PCIe 4.X/5.X: T9 (1bit). Refer to the PCI Express™ specification (PCIe) Gen 4. Set to 0b.
TC	Traffic Class (3 bits). Set to 000b for MCTP over PCIe VDM.
R or R Attr R TH or T8 Attr LN TH	PCIe 1.1/2.0: PCIe reserved bits (4 bits). Set to 0000b PCIe 2.1/3.X: PCIe reserved bit (1 bit), Attr[2] (1 bit) – Set to 0b, reserved bit (1bit), and TH (1bit) – Set to 0b. PCIe 4.X/5.X: T8 bit (1 bit) – Set to 0b, Attr[2] (1 bit) – Set to 0b, LN (1bit) – Set to 0b, and TH (1bit) – Set to 0b
TD	TLP Digest (1 bit). 1b indicates the presence of the TLP Digest field at the end of the PCIe TLP (transaction layer packet). The TD bit should be set in accordance with the devices overall support for the TLP Digest capability, and whether that capability is enabled. See description of the TLP Digest / ECRC field, below, for additional information. Note that earlier versions of this specification erroneously required this bit to be set to 0b, which would have required devices to not support the TLP Digest capability.
EP	Error Poisoned (1 bit).
Attr[1:0]	Attributes (2 bits). Set to 00b or 01b for all MCTP over PCIe VDM.
R or AT	PCIe 1.1: PCIe reserved bits (2 bits). PCIe 2.0/2.1/3.X/4.X/5.X: Address Type (AT) field. Set to 00b.
Length	Length: Length of the PCIe VDM Data in dwords. Implementations shall support the baseline transmission unit defined in the MCTP Base Specification . For example, supporting a baseline transmission unit of 64 bytes requires supporting PCIe VDM data up to 16 dwords. An implementation may optionally support larger transfer unit sizes.
PCI Requester ID	Bus/device/function or bus/function number of the managed endpoint sending the message.
Pad Len	Pad Length (2-bits). 1-based count (0 to 3) of the number of 0x00 pad bytes that have been added to the end of the packet to make the packet dword aligned with respect to. PCIe . Because only packets with the EOM bit set to 1b are allowed to be less than the transfer unit size, packets that have the EOM bit set to 0b will already be dword aligned and will thus not require any pad bytes and will have a pad length of 00b.
MCTP VDM Code	Value that uniquely differentiates MCTP messages from other DMTF VDMs. Set to 0000b for this transport mapping as defined in this specification.
Message Code	(8 bits). Set to 0111_1111b to indicate a Type 1 VDM.
PCI Target ID	(16 bits). For Route By ID messages, this is the bus/device/function number or bus/function number that is the physical address of the target endpoint. This field is ignored for Broadcast and for Route to Root Complex messages.
Vendor ID	(16 bits). Set to 6836 (0x1AB4) for DMTF VDMs. The most significant byte is in byte 10, the least significant byte is byte 11.
RSVD	MCTP reserved (4 bits). Set these bits to 0 when generating a message. Ignore them on incoming messages.
Hdr Version	MCTP version (4 bits) 0001b : For MCTP devices that conform to the MCTP Base Specification and this version of the PCIe VDM transport binding. All other settings: Reserved to support future packet header field expansion or header version.

Field	Description
00h PAD	Pad bytes. 0 to 3 bytes of 00h as required to fill out the overall PCIe VDM data to be an integral number of dwords. Because only packets with the EOM bit set to 1b are allowed to be less than the transfer unit size, packets that have the EOM bit set to 0b will already be dword aligned, and will thus not require any pad bytes and will have a pad length of 00b.
TLP Digest / ECRC	(32 bits). TLP Digest / ECRC (End-to-end CRC). This field is defined for all PCIe TLPs (Transaction Layer Packets). Device support for this field is optional. However, per PCIe v2.1/3.X/4.X/5.X : "If a device Function is enabled to generate ECRC, it must calculate and apply ECRC for all TLPs originated by the Function. If the device supports generating this field, it must support it for all TLPs." Additionally, per PCIe v2.1/3.X/4.X/5.X , if the ultimate PCI Express Receiver of the TLP does not support ECRC checking, the receiver must ignore the TLP Digest.

208 6.2 Supported media

209 This physical transport binding has been designed to work with the following media as defined in
 210 [DSP0239](#) and listed in Table 2. Use of this binding with other types of physical media is not covered by
 211 this specification. Refer to DSP0239 for all supported physical media by MCTP transport bindings.

212 An implementation that is compliant with this specification shall at least support one of the PCIe media
 213 listed in Table 2. Note that the CXL is built on the [PCI Express](#) (PCIe) physical and electrical interface.

214

Table 2 – Supported media

Physical Media Identifier	Description
0x08	PCIe 1.1 compatible
0x09	PCIe 2.0 compatible
0x0A	PCIe 2.1 compatible
0x0B	PCIe 3.X compatible
0x0C	PCIe 4.X compatible
0x0D	PCIe 5.X compatible, CXL 1.X/2.X compatible

215 6.3 Physical address format for MCTP control messages

216 The address format shown in Table 3 is used for MCTP control commands that require a physical
 217 address parameter to be returned for a bus that uses this transport binding with one of the supported
 218 media types listed in 6.2. This includes commands such as the Resolve Endpoint ID, Routing Information
 219 Update, and Get Routing Table Entries commands.

220

Table 3 – Physical address format

Format Size	Address Type	Layout and Description	
2 bytes (BDF ID)	Bus Device Function (BDF)	byte 1	[7:0] – Bus number
		byte 2	[7:3] – Device number [2:0] – Function number
2 bytes (ARI ID)	Alternate Routing Identifier (ARI)	byte 1	[7:0] – Bus number
		byte 2	[7:0] – Function number

221 6.4 Message routing

222 Physical packet routing within a PCIe bus uses routing as defined by the PCIe specification. PCIe
223 physical routing/bridging is not the same as MCTP bridging. PCIe physical routing/bridging is generally
224 transparent to MCTP. There are no MCTP-defined functions for configuring or controlling the setup of a
225 PCIe bus. The following types of PCIe addressing are used with MCTP messages:

- 226 • **Route by ID**

227 All MCTP over PCIe VDM packets between endpoints that are not the bus owner shall use
228 Route by ID for message routing.

229 The MCTP bus owner shall use Route by ID for messages to individual MCTP endpoints.

230 MCTP endpoints are required to capture the PCIe requester ID and the MCTP source EID when
231 receiving an EID assignment MCTP control request message. This is because this request can
232 only be issued by the MCTP bus owner.

- 233 • **Route to root complex**

234 MCTP endpoints shall use this routing for the Discovery Notify request message to the MCTP
235 bus owner as part of the MCTP over PCIe VDM discovery process.

236 The MCTP endpoints shall use this routing for responding to the MCTP control request
237 messages that were sent using Broadcast from Root Complex.

238 Communication of MCTP PCIe VDM packets that are destined to MCTP bus owner using
239 routed to root complex is implementation specific and is outside the scope of this specification.

- 240 • **Broadcast from root complex**

241 The MCTP bus owner should use this routing for the Prepare for Endpoint Discovery and
242 Endpoint Discovery messages as part of the MCTP over PCIe VDM discovery process.

243

244 6.4.1 Routing peer transactions on a PCIe bus

245 Because the PCIe specification does not require peer-to-peer routing support in PCIe root complexes,
246 MCTP over PCIe VDM messages are not required to be routed to peer devices directly. When peer-to-
247 peer routing is not supported by a PCIe root complex, all MCTP over PCIe VDM messages between two
248 MCTP endpoints shall be routed to or through the MCTP bus owner as an MCTP bridge. If the PCIe root
249 complex, as the MCTP bus owner, supports peer-to-peer routing, it shall use direct physical addressing to
250 support routing between two MCTP endpoints on the PCIe bus.

251 6.4.2 Routing messages between PCIe and other buses

252 All MCTP messages that span between PCIe and other buses shall be sent through the MCTP bus
253 owner. The MCTP bus owner has the destination EID routing tables necessary to route messages
254 between the two bus segments.

255 If an endpoint is aware of multiple routes to a destination over multiple bus types, a higher level
256 algorithm/protocol above MCTP shall be used to determine which bus/route to use. Typically this decision
257 can be based on things like power state and MCTP discovery state.

258 6.5 Bus owner address

259 The MCTP PCIe VDM bus owner functionality shall be accessible through "Route-to-Root Complex"
260 addressing.

261 6.6 Bus address assignment for PCIe

262 PCIe bus addresses are assigned per the mechanisms specified in [PCIe](#).

263 6.7 Host dependencies

264 MCTP over PCIe VDM, when used in a typical “PC” computer system, has a dependency on the host
265 CPU, host software, power management states, link states, and reset. Some of these dependencies are
266 described as follows:

267 • **Reset**

268 Assertion of “Fundamental Reset” on the bus causes both the host functionality as well as the
269 MCTP PCIe VDM communication on an MCTP PCIe endpoint to be reset. From the assertion
270 “Fundamental Reset” until the PCIe fabric has been configured and enumerated, no “MCTP
271 over PCI Express” messages can be sent.

272 Similarly, if MCTP PCIe VDM communication is supported on a function, a function level reset
273 (FLR) could reset MCTP PCIe VDM endpoint as well as MCTP PCIe VDM communication on
274 that function.

275 • **Configuration and enumeration**

276 Following the de-assertion “Fundamental Reset”, the software running on the host CPU
277 configures and enumerates the PCIe fabric. Failure of the host CPU or boot software to properly
278 configure and enumerate the PCIe fabric prevents it from being used for MCTP over PCIe VDM
279 messaging.

280 • **Power management states**

281 The host (as defined in the context of the [PCI Express™ specification](#)) controls PCIe bus power
282 management. The host may power down PCIe devices and links, or place them in sleep states,
283 independent of management controllers, which may cause MCTP PCIe VDM communication to
284 be unavailable. Depending on the device usage in the system, a PCIe device may retain or lose
285 states such as EID, “discovered” state, and routing information (if the device is a bridge). A
286 PCIe device that loses MCTP PCIe VDM communication state needs to be reinitialized and/or
287 rediscovered after it returns to a power state that supports MCTP over PCIe VDM
288 communication.

289 • **Link states**

290 The PCIe link states affect MCTP over PCIe VDM communications. MCTP over PCIe VDM
291 communication can be performed only when the PCIe link is in a state that allows VDM
292 communications. The mechanisms for PCIe link state transitions are outside the scope of this
293 specification.

294 • **PCIe Root Complex**

295 PCIe Root Complex (RC) is responsible for communicating route-to-root complex MCTP over
296 PCIe VDM discovery messages to the MCTP bus owner.

297 6.8 Discovery Notify message use for PCIe

298 An MCTP control Discovery Notify message shall be sent from a PCIe endpoint to the MCTP bus owner
299 whenever the physical address for the device changes (that is, the endpoint receives a Type 0
300 configuration write request and the bus number is different than the currently stored bus number). This
301 occurs on the first Type 0 configuration write following a PCIe bus reset during initial enumeration, or
302 during re-enumeration where the bus number has changed (for example, because of a hot plug event,
303 bus reset, and so on).

304 Endpoints use the Discovery Notify command to inform the MCTP bus owner that it needs to update the
305 endpoint's ID. The Discovery Notify command shall be sent with the PCIe message routing set to 000b
306 (Route-to-Root Complex), the Destination Endpoint ID for the Discovery Notify message shall be set to
307 the Null Destination EID. The Source Endpoint ID field shall be set to the Null Source EID if the device
308 has not yet been assigned an EID; otherwise, it shall contain the assigned EID value.

309 **6.9 MCTP over PCIe endpoint discovery**

310 This clause describes the steps used to support discovering MCTP endpoints on PCIe.

311 **6.9.1 Discovered flag**

312 Each endpoint (except the bus owner) on the PCIe bus maintains an internal flag called the *Discovered*
313 flag.

314 The flag is set to the *discovered* state when the Set Endpoint ID command is received.

315 The Prepare for Endpoint Discovery message causes each recipient endpoint on the PCIe bus to set their
316 respective Discovered flag to the *undiscovered* state. For the Prepare for Endpoint Discovery request
317 message, the routing in the physical transport header should be set to 011b (Broadcast from Root
318 Complex).

319 An endpoint also sets the flag to the *undiscovered* state at the following times:

- 320 • Whenever the PCI bus/device/function or bus/function number associated with the endpoint is
321 initially assigned or is changed to a different value.
- 322 • Whenever an endpoint first appears on the bus and requires an EID assignment. A device shall
323 have been enumerated on PCI and have a bus/device/function or bus/function number before it
324 can do this.
- 325 • During operation if an endpoint enters a state that causes it to lose its EID assignment.
- 326 • For endpoints that have already received an EID assignment but are in any temporary state
327 where the endpoint was unable to respond to MCTP control requests for more than T_{RECLAIM}
328 seconds.

329 Only endpoints that have their Discovered flag set to *undiscovered* shall respond to the Endpoint
330 Discovery message. Endpoints that have the flag set to *discovered* shall not respond to the Endpoint
331 Discovery message.

332 For PCIe endpoints, an Endpoint Discovery broadcast request message can be sent by the MCTP bus
333 owner to discover all MCTP-capable devices. MCTP-capable endpoints respond with an Endpoint
334 Discovery response message.

335 **6.9.2 PCIe endpoint announcement**

336 One or more endpoints may announce their presence and their need for an EID assignment by
337 autonomously sending a Discovery Notify message to the bus owner. This would typically trigger the
338 MCTP bus owner to perform the PCIe endpoint discovery/enumeration processes described in the
339 following subclauses.

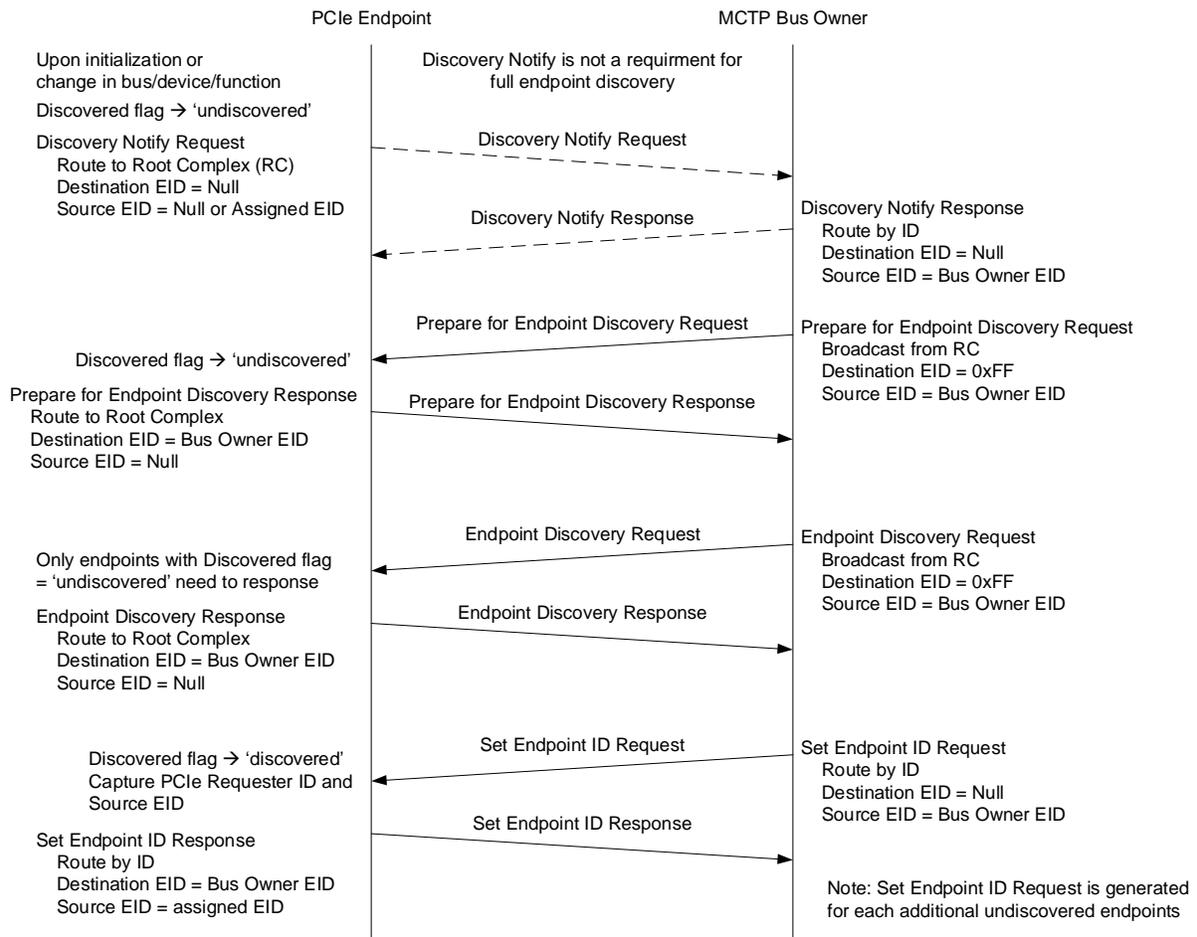
340 **6.9.3 Full endpoint Discovery/Enumeration**

341 The following process is typically used when the MCTP bus owner wishes to discover and enumerate all
342 MCTP endpoints on the PCIe bus.

- 343 1) The MCTP bus owner issues a broadcast Prepare for Endpoint Discovery message. This
344 message causes each discoverable endpoint on the bus to set its PCIe endpoint Discovered

- 345 flag to undiscovered. Depending on the number of endpoints and the buffer space available in
346 the MCTP bus owner, the MCTP bus owner may not receive all of the response messages. The
347 discovery process does not require the MCTP bus owner to receive all the response messages
348 to the Prepare for Endpoint Discovery request. Because the MCTP bus owner cannot determine
349 that all endpoints have received the Prepare for Endpoint Discovery request, it is recommended
350 that Prepare for Endpoint Discovery request is retried MN1 times to help ensure that all
351 endpoints have received the request. The MCTP bus owner is not required to wait for MT2 time
352 interval between the retries.
- 353 2) The MCTP bus owner should wait for MT2 time interval to help ensure that all endpoints that
354 received the Prepare for Endpoint Discovery request have processed the request.
- 355 3) The MCTP bus owner issues a broadcast Endpoint Discovery request message. All MCTP-
356 capable devices that have their Discovered flag set to undiscovered will respond with an
357 Endpoint Discovery response message.
- 358 4) Depending on the number of endpoints and the buffer space available in the MCTP bus owner,
359 the MCTP bus owner receives some or all of these response messages. For each response
360 message received from an undiscovered MCTP-capable device PCIe bus/device/function or
361 bus/function number, the MCTP bus owner issues a Set Endpoint ID command to the physical
362 address for the endpoint. This causes the endpoint to set its Discovered flag to *discovered*.
363 From this point, the endpoint shall not respond to the Endpoint Discovery command until
364 another Prepare for Endpoint Discovery command is received or some other condition causes
365 the Discovered flag to be set back to *undiscovered*.
- 366 5) If the MCTP bus owner received any responses to the Endpoint Discovery request issued in
367 Step 3, then it shall repeat steps 3 and 4 until it no longer gets any responses to the Endpoint
368 Discovery request. In this case, then the MCTP bus owner is allowed to send the next Endpoint
369 Discovery request without waiting for MT2 time interval. If no responses were received by the
370 MCTP bus owner to the Endpoint Discovery request within the MT2 time interval, then the
371 discovery process is completed.
- 372 After the initial endpoint enumeration, it is recommended that the MCTP bus owner maintains a list of the
373 unique IDs for the endpoints it has discovered, and reassigns the same IDs to those endpoints if a
374 bus/device/function or bus/function number changes during system operation.
- 375 Figure 2 provides an example flow of operations for full endpoint discovery.

Full PCIe MCTP Endpoint Discovery



376

377

Figure 2 – Flow of operations for full MCTP Discovery over PCIe

378 6.9.4 Partial endpoint Discovery/Enumeration

379 This process is used when the MCTP bus owner wishes to discover endpoints that may have been added
 380 to the bus after a full enumeration has been done. This situation can occur if a device has its
 381 bus/device/function or bus/function number change after the full enumeration has been done, or when a
 382 hot-plug device is added to the system, or if a device that is already present in the system — but was in a
 383 disabled or powered-down state — comes on-line.

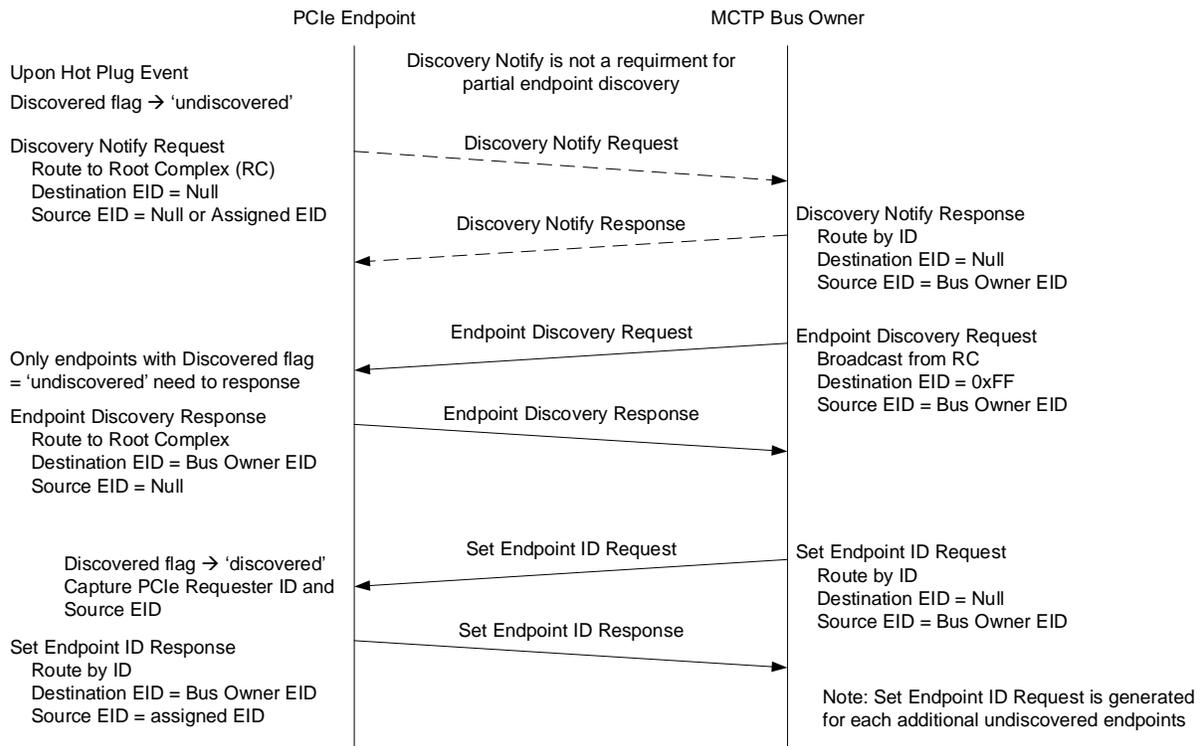
384 The partial discovery process is the same as the full discovery process except that the MCTP bus owner
 385 skips the step of broadcasting a Prepare for Endpoint Discovery command in order to avoid clearing the
 386 Discovered flags of already discovered endpoints.

387 The partial discovery process may be initiated when a device that is added or enabled for MCTP sends a
 388 Discovery Notify message to the MCTP bus owner. The MCTP bus owner may also elect to periodically
 389 issue a broadcast Endpoint Discovery message to test for whether any undiscovered endpoints have
 390 been missed. The Discovery Notify message provides the MCTP bus owner with the bus/device/function
 391 or bus/function number of the MCTP PCIe endpoint. The MCTP bus owner can then send a directed
 392 Endpoint Discovery message to the endpoint to confirm that the device has not been discovered. The
 393 MCTP bus owner then issues a Set Endpoint ID command to the physical address for the endpoint which
 394 causes the endpoint to set its Discovered flag to *discovered*.

395 It is recommended that the MCTP bus owner maintains a list of the unique MCTP EIDs for the endpoints
 396 it has discovered and reassigns the same MCTP EIDs to those endpoints if a bus/device/function or
 397 bus/function number changes during system operation.

398 Figure 3 provides an example flow of operations for partial endpoint discovery.

Partial PCIe MCTP Endpoint Discovery



399

400

Figure 3 – Flow of operations for Partial Endpoint Discovery

401 **6.9.5 Endpoint re-enumeration**

402 If the bus implementation includes hot-plug devices, the bus owner shall perform a full or partial endpoint
 403 discovery any time the MCTP bus owner goes into a temporary state where the MCTP bus owner can
 404 miss receiving a Discovery Notify message (for example, if the bus owner device is reset or receives a
 405 firmware update). Whether a full or partial endpoint discovery is required is dependent on how much
 406 information the MCTP bus owner retains from prior enumerations.

407 **6.10 MCTP messages timing requirements**

408 Table 4 lists MCTP-specific timing requirements for MCTP Control messages and operation on the PCIe
 409 VDM medium. All MCTP Control Messages over PCIe VDM shall comply to the timing specification listed
 410 in Table 4.

411 **Table 4 – Timing specifications for MCTP Control messages on PCIe VDM**

Timing Specification	Symbol	Min	Max	Description
Endpoint ID reclaim	T _{RECLAIM}	–	5 sec	Maximum interval that an endpoint is allowed to be non-responsive to MCTP control messages before its EID may be reclaimed by the bus owner. A bus owner shall wait at least for this interval before an EID of the non-responsive endpoint is reclaimed.
Number of request retries	MN1	2	See Description column	Total of three tries, minimum: the original try plus two retries. The maximum number of retries for a given request is limited by the requirement that all retries shall occur within MT4, max of the initial request.
Request-to-response time	MT1	–	120 ms	This interval is measured at the responder from the end of the reception of an MCTP control request to the beginning of the transmission of the corresponding MCTP control response. This requirement is tested under the condition where the responder can successfully transmit the response on the first try.
Time-out waiting for a response	MT2	MT1 max ^[1] + 6 ms	MT4, min ^[1]	This interval at the requester sets the minimum amount of time that a requester should wait before retrying a MCTP control request. This interval is measured at the requester from the end of the successful transmission of the MCTP control request to the beginning of the reception of the corresponding MCTP control response. NOTE: This specification does not preclude an implementation from adjusting the minimum time-out waiting for a response to a smaller number than MT2 based on the measured response times from responders. The mechanism for doing so is outside the scope of this specification.
Instance ID expiration interval	MT4	5 sec ^[2]	6 sec	Interval after which the instance ID for a given response will expire and become reusable if a response has not been received for the request. This is also the maximum time that a responder tracks an instance ID for a given request from a given requester.

Timing Specification	Symbol	Min	Max	Description
<p>NOTE 1: Unless otherwise specified, this timing applies to the mandatory and optional MCTP commands.</p> <p>NOTE 2: If a requester is reset, it may produce the same sequence number for a request as one that was previously issued. To guard against this, it is recommended that sequence number expiration be implemented. Any request from a given requester that is received more than MT4 seconds after a previous, matching request should be treated as a new request, not a retry.</p>				

ANNEX A (informative)

Notations and conventions

416 Notations

417 Examples of notations used in this document are as follows: list into text needed

- 418 • 2:N In field descriptions, this will typically be used to represent a range of byte offsets
419 starting from byte two and continuing to and including byte N. The lowest offset is on
420 the left, the highest is on the right.
- 421 • (6) Parentheses around a single number can be used in message field descriptions to
422 indicate a byte field that may be present or absent.
- 423 • (3:6) Parentheses around a field consisting of a range of bytes indicates the entire range
424 may be present or absent. The lowest offset is on the left, the highest is on the right.
- 425 • [PCIe](#) Underlined, blue text is typically used to indicate a reference to a document or
426 specification called out in Clause 2, "Normative References" or to items hyperlinked
427 within the document.
- 428 • rsvd Abbreviation for Reserved. Case insensitive.
- 429 • [4] Square brackets around a number are typically used to indicate a bit offset. Bit offsets
430 are given as 0-based values (that is, the least significant bit [LSb] offset = 0).
- 431 • [7:5] A range of bit offsets. The most significant bit is on the left, the least significant bit is
432 on the right.
- 433 • 1b The lower case "b" following a number consisting of 0s and 1s is used to indicate the
434 number is being given in binary format.
- 435 • 0x12A A leading "0x" is used to indicate a number given in hexadecimal format.

436
437
438
439
440

ANNEX B (informative)

Change log

Version	Date	Description
1.0.0	2009-07-28	
1.0.1	2009-10-30	Created erratum to clarify Length field definition of PCIe VDM header for MCTP PCIe VDM transport binding, modify introduction section, and clean up references section.
1.0.2	2014-12-07	Clarifications to TD bit usage. Added TLP Digest/ECRC to packet figure and to field descriptions table.
1.1.0	2018-10-24	Added support for PCIe Gen 3, PCIe Gen 4, and ARI. Fixed Figure 1 to cover PCIe 1.0/2.0/2.1/3.X/4.0. Clarified MCTP over PCIe VDM compliant management device requirements. Clarified Endpoint ID reclaim definition. Clarified MCTP bus owner requirements in the specification. Eliminated PCIe bus owner term and replaced it with PCIe root complex where applicable.
1.2.0	2021-03-02	Added support for PCIe Gen 5.X, CXL 1.X, and CXL 2.X.

441
442