



# Virtual I/O Device (VIRTIO) Version 1.1

## Committee Specification 01

11 April 2019

### Specification URIs

#### This version:

<https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/tex/> (Authoritative)  
<https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/virtio-v1.1-cs01.pdf>  
<https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/virtio-v1.1-cs01.html>

#### Previous version:

N/A

#### Latest version:

<https://docs.oasis-open.org/virtio/virtio/v1.1/virtio-v1.1.pdf>  
<https://docs.oasis-open.org/virtio/virtio/v1.1/virtio-v1.1.html>

#### Technical Committee:

OASIS Virtual I/O Device (VIRTIO) TC

#### Chairs:

Michael S. Tsirkin ([mst@redhat.com](mailto:mst@redhat.com)), Red Hat  
Cornelia Huck ([cohuck@redhat.com](mailto:cohuck@redhat.com)), Red Hat

#### Editors:

Michael S. Tsirkin ([mst@redhat.com](mailto:mst@redhat.com)), Red Hat  
Cornelia Huck ([cohuck@redhat.com](mailto:cohuck@redhat.com)), Red Hat

#### Additional artifacts:

This prose specification is one component of a Work Product that also includes:

- Example Driver Listing:  
<https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/listings/>

#### Related work:

This specification replaces or supersedes:

- Virtual I/O Device (VIRTIO) Version 1.0. Edited by Rusty Russell, Michael S. Tsirkin, Cornelia Huck, and Pawel Moll. Latest version:  
<https://docs.oasis-open.org/virtio/virtio/v1.0/virtio-v1.0.html>
- Virtio PCI Card Specification Version 0.9.5:  
<http://ozlabs.org/~rusty/virtio-spec/virtio-0.9.5.pdf>

### Abstract:

This document describes the specifications of the “virtio” family of devices. These devices are found in virtual environments, yet by design they look like physical devices to the guest within the virtual machine - and this document treats them as such. This similarity allows the guest to use standard drivers and discovery mechanisms.

The purpose of virtio and this specification is that virtual environments and guests should have a straightforward, efficient, standard and extensible mechanism for virtual devices, rather than boutique per-environment or per-OS mechanisms.

### Status:

This document was last revised or approved by the Virtual I/O Device (VIRTIO) TC on the above date. The level of approval is also listed above. Check the “Latest version” location noted above for possible later revisions of this document. Any other numbered Versions and other technical work produced by the Technical Committee (TC) are listed at [https://www.oasis-open.org/committees/tc\\_home.php?wg\\_abbrev=virtio#technical](https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=virtio#technical).

Technical Committee members should send comments on this specification to the Technical Committee's email list. Others should send comments to the Technical Committee by using the “Send A Comment” button on the Technical Committee's web page at <https://www.oasis-open.org/committees/virtio/>.

This specification is provided under the [Non-Assertion](#) Mode of the [OASIS IPR Policy](#), the mode chosen when the Technical Committee was established. For information on whether any patents have been disclosed that may be essential to implementing this specification, and any offers of patent licensing terms, please refer to the Intellectual Property Rights section of the TC's web page (<https://github.com/oasis-tcs/virtio-admin/blob/master/IPR.md>).

Note that any machine-readable content ([Computer Language Definitions](#)) declared Normative for this Work Product is provided in separate plain text files. In the event of a discrepancy between any such plain text file and display content in the Work Product's prose narrative document(s), the content in the separate plain text file prevails.

### Citation format:

When referencing this specification the following citation format should be used:

#### [VIRTIO-v1.1]

*Virtual I/O Device (VIRTIO) Version 1.1*. Edited by Michael S. Tsirkin and Cornelia Huck. 11 April 2019. OASIS Committee Specification 01. <https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/virtio-v1.1-cs01.html>. Latest version: <https://docs.oasis-open.org/virtio/virtio/v1.1/virtio-v1.1.html>.

# Notices

Copyright © OASIS Open 2018. All Rights Reserved.

All capitalized terms in the following text have the meanings assigned to them in the OASIS Intellectual Property Rights Policy (the "OASIS IPR Policy"). The full [Policy](#) may be found at the OASIS website.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published, and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this section are included on all such copies and derivative works. However, this document itself may not be modified in any way, including by removing the copyright notice or references to OASIS, except as needed for the purpose of developing any document or deliverable produced by an OASIS Technical Committee (in which case the rules applicable to copyrights, as set forth in the OASIS IPR Policy, must be followed) or as required to translate it into languages other than English.

This specification is provided under the [Non-Assertion](#) Mode of the [OASIS IPR Policy](#), the mode chosen when the Technical Committee was established. For information on whether any patents have been disclosed that may be essential to implementing this specification, and any offers of patent licensing terms, please refer to the Intellectual Property Rights section of the TC's web page (<https://github.com/oasis-tcs/virtio-admin/blob/master/IPR.md>).

The limited permissions granted above are perpetual and will not be revoked by OASIS or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and OASIS DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY OWNERSHIP RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

OASIS requests that any OASIS Party or any other party that believes it has patent claims that would necessarily be infringed by implementations of this OASIS Committee Specification or OASIS Standard, to notify OASIS TC Administrator and provide an indication of its willingness to grant patent licenses to such patent claims in a manner consistent with the IPR Mode of the OASIS Technical Committee that produced this specification.

OASIS invites any party to contact the OASIS TC Administrator if it is aware of a claim of ownership of any patent claims that would necessarily be infringed by implementations of this specification by a patent holder that is not willing to provide a license to such patent claims in a manner consistent with the IPR Mode of the OASIS Technical Committee that produced this specification. OASIS may include such claims on its website, but disclaims any obligation to do so.

OASIS takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Information on OASIS' procedures with respect to rights in any document or deliverable produced by an OASIS Technical Committee can be found on the OASIS website. Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this OASIS Committee Specification or OASIS Standard, can be obtained from the OASIS TC Administrator. OASIS makes no representation that any information or list of intellectual property rights will at any time be complete, or that any claims in such list are, in fact, Essential Claims.

The name "OASIS" is a trademark of [OASIS](#), the owner and developer of this specification, and should be used only to refer to the organization and its official outputs. OASIS welcomes reference to, and implementation and use of, specifications, while reserving the right to enforce its marks against misleading uses. Please see <https://www.oasis-open.org/policies-guidelines/trademark> for above guidance.

---

# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>17</b>
1.1	Normative References	17
1.2	Non-Normative References	18
1.3	Terminology	18
1.3.1	Legacy Interface: Terminology	18
1.3.2	Transition from earlier specification drafts	19
1.4	Structure Specifications	19
1.5	Constant Specifications	20
<b>2</b>	<b>Basic Facilities of a Virtio Device</b>	<b>21</b>
2.1	Device Status Field	21
2.1.1	Driver Requirements: Device Status Field	21
2.1.2	Device Requirements: Device Status Field	22
2.2	Feature Bits	22
2.2.1	Driver Requirements: Feature Bits	22
2.2.2	Device Requirements: Feature Bits	22
2.2.3	Legacy Interface: A Note on Feature Bits	22
2.3	Notifications	23
2.4	Device Reset	23
2.4.1	Device Requirements: Device Reset	23
2.4.2	Driver Requirements: Device Reset	23
2.5	Device Configuration Space	23
2.5.1	Driver Requirements: Device Configuration Space	24
2.5.2	Device Requirements: Device Configuration Space	24
2.5.3	Legacy Interface: A Note on Device Configuration Space endian-ness	24
2.5.4	Legacy Interface: Device Configuration Space	24
2.6	Virtqueues	24
2.6.1	Virtqueue Reset	25
2.6.1.1	Virtqueue Reset	25
2.6.1.1.1	Device Requirements: Virtqueue Reset	25
2.6.1.1.2	Driver Requirements: Virtqueue Reset	25
2.6.1.2	Virtqueue Re-enable	25
2.6.1.2.1	Device Requirements: Virtqueue Re-enable	25
2.6.1.2.2	Driver Requirements: Virtqueue Re-enable	26
2.7	Split Virtqueues	26
2.7.1	Driver Requirements: Virtqueues	26
2.7.2	Legacy Interfaces: A Note on Virtqueue Layout	26
2.7.3	Legacy Interfaces: A Note on Virtqueue Endianness	27
2.7.4	Message Framing	27
2.7.4.1	Device Requirements: Message Framing	27
2.7.4.2	Driver Requirements: Message Framing	27
2.7.4.3	Legacy Interface: Message Framing	27
2.7.5	The Virtqueue Descriptor Table	28
2.7.5.1	Device Requirements: The Virtqueue Descriptor Table	28
2.7.5.2	Driver Requirements: The Virtqueue Descriptor Table	28
2.7.5.3	Indirect Descriptors	28
2.7.5.3.1	Driver Requirements: Indirect Descriptors	29

2.7.5.3.2	Device Requirements: Indirect Descriptors	29
2.7.6	The Virtqueue Available Ring	29
2.7.6.1	Driver Requirements: The Virtqueue Available Ring	29
2.7.7	Used Buffer Notification Suppression	30
2.7.7.1	Driver Requirements: Used Buffer Notification Suppression	30
2.7.7.2	Device Requirements: Used Buffer Notification Suppression	30
2.7.8	The Virtqueue Used Ring	30
2.7.8.1	Legacy Interface: The Virtqueue Used Ring	31
2.7.8.2	Device Requirements: The Virtqueue Used Ring	31
2.7.8.3	Driver Requirements: The Virtqueue Used Ring	31
2.7.9	In-order use of descriptors	31
2.7.10	Available Buffer Notification Suppression	32
2.7.10.1	Driver Requirements: Available Buffer Notification Suppression	32
2.7.10.2	Device Requirements: Available Buffer Notification Suppression	32
2.7.11	Helpers for Operating Virtqueues	32
2.7.12	Virtqueue Operation	33
2.7.13	Supplying Buffers to The Device	33
2.7.13.1	Placing Buffers Into The Descriptor Table	33
2.7.13.2	Updating The Available Ring	34
2.7.13.3	Updating <i>idx</i>	34
2.7.13.3.1	Driver Requirements: Updating <i>idx</i>	34
2.7.13.4	Notifying The Device	34
2.7.13.4.1	Driver Requirements: Notifying The Device	34
2.7.14	Receiving Used Buffers From The Device	34
2.8	Packed Virtqueues	35
2.8.1	Driver and Device Ring Wrap Counters	36
2.8.2	Polling of available and used descriptors	36
2.8.3	Write Flag	36
2.8.4	Element Address and Length	37
2.8.5	Scatter-Gather Support	37
2.8.6	Next Flag: Descriptor Chaining	37
2.8.7	Indirect Flag: Scatter-Gather Support	37
2.8.8	In-order use of descriptors	38
2.8.9	Multi-buffer requests	38
2.8.10	Driver and Device Event Suppression	38
2.8.10.1	Structure Size and Alignment	39
2.8.11	Driver Requirements: Virtqueues	39
2.8.12	Device Requirements: Virtqueues	39
2.8.13	The Virtqueue Descriptor Format	39
2.8.14	Event Suppression Structure Format	40
2.8.15	Device Requirements: The Virtqueue Descriptor Table	40
2.8.16	Driver Requirements: The Virtqueue Descriptor Table	40
2.8.17	Driver Requirements: Scatter-Gather Support	40
2.8.18	Device Requirements: Scatter-Gather Support	40
2.8.19	Driver Requirements: Indirect Descriptors	41
2.8.20	Virtqueue Operation	41
2.8.21	Supplying Buffers to The Device	41
2.8.21.1	Placing Available Buffers Into The Descriptor Ring	41
2.8.21.1.1	Driver Requirements: Updating flags	42
2.8.21.2	Sending Available Buffer Notifications	42
2.8.21.3	Implementation Example	42
2.8.21.3.1	Driver Requirements: Sending Available Buffer Notifications	43
2.8.22	Receiving Used Buffers From The Device	43
2.9	Driver Notifications	44
2.10	Shared Memory Regions	44
2.10.1	Addressing within regions	44
2.10.2	Device Requirements: Shared Memory Regions	45

2.11	Exporting Objects	45
<b>3</b>	<b>General Initialization And Device Operation</b>	<b>46</b>
3.1	Device Initialization	46
3.1.1	Driver Requirements: Device Initialization	46
3.1.2	Legacy Interface: Device Initialization	46
3.2	Device Operation	47
3.2.1	Notification of Device Configuration Changes	47
3.3	Device Cleanup	47
3.3.1	Driver Requirements: Device Cleanup	47
<b>4</b>	<b>Virtio Transport Options</b>	<b>48</b>
4.1	Virtio Over PCI Bus	48
4.1.1	Device Requirements: Virtio Over PCI Bus	48
4.1.2	PCI Device Discovery	48
4.1.2.1	Device Requirements: PCI Device Discovery	48
4.1.2.2	Driver Requirements: PCI Device Discovery	49
4.1.2.3	Legacy Interfaces: A Note on PCI Device Discovery	49
4.1.3	PCI Device Layout	49
4.1.3.1	Driver Requirements: PCI Device Layout	49
4.1.3.2	Device Requirements: PCI Device Layout	49
4.1.4	Virtio Structure PCI Capabilities	49
4.1.4.1	Driver Requirements: Virtio Structure PCI Capabilities	51
4.1.4.2	Device Requirements: Virtio Structure PCI Capabilities	51
4.1.4.3	Common configuration structure layout	51
4.1.4.3.1	Device Requirements: Common configuration structure layout	52
4.1.4.3.2	Driver Requirements: Common configuration structure layout	53
4.1.4.4	Notification structure layout	53
4.1.4.4.1	Device Requirements: Notification capability	54
4.1.4.5	ISR status capability	54
4.1.4.5.1	Device Requirements: ISR status capability	55
4.1.4.5.2	Driver Requirements: ISR status capability	55
4.1.4.6	Device-specific configuration	55
4.1.4.6.1	Device Requirements: Device-specific configuration	55
4.1.4.7	Shared memory capability	55
4.1.4.7.1	Device Requirements: Shared memory capability	55
4.1.4.8	Vendor data capability	55
4.1.4.8.1	Device Requirements: Vendor data capability	56
4.1.4.8.2	Driver Requirements: Vendor data capability	56
4.1.4.9	PCI configuration access capability	56
4.1.4.9.1	Device Requirements: PCI configuration access capability	56
4.1.4.9.2	Driver Requirements: PCI configuration access capability	57
4.1.4.10	Legacy Interfaces: A Note on PCI Device Layout	57
4.1.4.11	Non-transitional Device With Legacy Driver: A Note on PCI Device Layout	58
4.1.5	PCI-specific Initialization And Device Operation	58
4.1.5.1	Device Initialization	58
4.1.5.1.1	Virtio Device Configuration Layout Detection	58
4.1.5.1.2	MSI-X Vector Configuration	58
4.1.5.1.3	Virtqueue Configuration	59
4.1.5.2	Available Buffer Notifications	60
4.1.5.2.1	Driver Requirements: Available Buffer Notifications	60
4.1.5.3	Used Buffer Notifications	60
4.1.5.3.1	Device Requirements: Used Buffer Notifications	60
4.1.5.4	Notification of Device Configuration Changes	61
4.1.5.4.1	Device Requirements: Notification of Device Configuration Changes	61
4.1.5.4.2	Driver Requirements: Notification of Device Configuration Changes	61
4.1.5.5	Driver Handling Interrupts	61

4.2	Virtio Over MMIO	61
4.2.1	MMIO Device Discovery	61
4.2.2	MMIO Device Register Layout	62
4.2.2.1	Device Requirements: MMIO Device Register Layout	65
4.2.2.2	Driver Requirements: MMIO Device Register Layout	65
4.2.3	MMIO-specific Initialization And Device Operation	66
4.2.3.1	Device Initialization	66
4.2.3.1.1	Driver Requirements: Device Initialization	66
4.2.3.2	Virtqueue Configuration	66
4.2.3.3	Available Buffer Notifications	66
4.2.3.4	Notifications From The Device	66
4.2.3.4.1	Driver Requirements: Notifications From The Device	66
4.2.4	Legacy interface	67
4.3	Virtio Over Channel I/O	69
4.3.1	Basic Concepts	69
4.3.1.1	Channel Commands for Virtio	70
4.3.1.2	Notifications	70
4.3.1.3	Device Requirements: Basic Concepts	70
4.3.1.4	Driver Requirements: Basic Concepts	70
4.3.2	Device Initialization	70
4.3.2.1	Setting the Virtio Revision	70
4.3.2.1.1	Device Requirements: Setting the Virtio Revision	71
4.3.2.1.2	Driver Requirements: Setting the Virtio Revision	71
4.3.2.1.3	Legacy Interfaces: A Note on Setting the Virtio Revision	71
4.3.2.2	Configuring a Virtqueue	71
4.3.2.2.1	Device Requirements: Configuring a Virtqueue	72
4.3.2.2.2	Legacy Interface: A Note on Configuring a Virtqueue	72
4.3.2.3	Communicating Status Information	72
4.3.2.3.1	Driver Requirements: Communicating Status Information	72
4.3.2.3.2	Device Requirements: Communicating Status Information	73
4.3.2.4	Handling Device Features	73
4.3.2.5	Device Configuration	73
4.3.2.6	Setting Up Indicators	73
4.3.2.6.1	Setting Up Classic Queue Indicators	73
4.3.2.6.2	Setting Up Configuration Change Indicators	74
4.3.2.6.3	Setting Up Two-Stage Queue Indicators	74
4.3.2.6.4	Legacy Interfaces: A Note on Setting Up Indicators	74
4.3.3	Device Operation	74
4.3.3.1	Host->Guest Notification	74
4.3.3.1.1	Notification via Classic I/O Interrupts	74
4.3.3.1.2	Notification via Adapter I/O Interrupts	75
4.3.3.1.3	Legacy Interfaces: A Note on Host->Guest Notification	75
4.3.3.2	Guest->Host Notification	75
4.3.3.2.1	Device Requirements: Guest->Host Notification	76
4.3.3.2.2	Driver Requirements: Guest->Host Notification	76
4.3.3.3	Resetting Devices	76
4.3.3.3.1	Device Requirements: Resetting Devices	76
4.3.3.3.2	Driver Requirements: Resetting Devices	76
5	Device Types	77
5.1	Network Device	79
5.1.1	Device ID	79
5.1.2	Virtqueues	79
5.1.3	Feature bits	79
5.1.3.1	Feature bit requirements	80
5.1.3.2	Legacy Interface: Feature bits	81
5.1.4	Device configuration layout	81



5.1.4.1	Device Requirements: Device configuration layout	82
5.1.4.2	Driver Requirements: Device configuration layout	82
5.1.4.3	Legacy Interface: Device configuration layout	82
5.1.5	Device Initialization	83
5.1.6	Device Operation	83
5.1.6.1	Legacy Interface: Device Operation	84
5.1.6.2	Packet Transmission	84
5.1.6.2.1	Driver Requirements: Packet Transmission	85
5.1.6.2.2	Device Requirements: Packet Transmission	86
5.1.6.2.3	Packet Transmission Interrupt	86
5.1.6.3	Setting Up Receive Buffers	86
5.1.6.3.1	Driver Requirements: Setting Up Receive Buffers	86
5.1.6.3.2	Device Requirements: Setting Up Receive Buffers	86
5.1.6.4	Processing of Incoming Packets	87
5.1.6.4.1	Device Requirements: Processing of Incoming Packets	87
5.1.6.4.2	Driver Requirements: Processing of Incoming Packets	88
5.1.6.4.3	Hash calculation for incoming packets	88
5.1.6.4.4	Hash reporting for incoming packets	91
5.1.6.5	Control Virtqueue	91
5.1.6.5.1	Packet Receive Filtering	92
5.1.6.5.2	Setting MAC Address Filtering	93
5.1.6.5.3	VLAN Filtering	94
5.1.6.5.4	Gratuitous Packet Sending	94
5.1.6.5.5	Device operation in multiqueue mode	94
5.1.6.5.6	Automatic receive steering in multiqueue mode	95
5.1.6.5.7	Receive-side scaling (RSS)	96
5.1.6.5.8	Offloads State Configuration	97
5.1.6.6	Legacy Interface: Framing Requirements	98
5.2	Block Device	98
5.2.1	Device ID	98
5.2.2	Virtqueues	98
5.2.3	Feature bits	98
5.2.3.1	Legacy Interface: Feature bits	99
5.2.4	Device configuration layout	99
5.2.4.1	Legacy Interface: Device configuration layout	100
5.2.5	Device Initialization	100
5.2.5.1	Driver Requirements: Device Initialization	100
5.2.5.2	Device Requirements: Device Initialization	101
5.2.5.3	Legacy Interface: Device Initialization	101
5.2.6	Device Operation	101
5.2.6.1	Driver Requirements: Device Operation	102
5.2.6.2	Device Requirements: Device Operation	103
5.2.6.3	Legacy Interface: Device Operation	104
5.2.6.4	Legacy Interface: Framing Requirements	105
5.3	Console Device	105
5.3.1	Device ID	106
5.3.2	Virtqueues	106
5.3.3	Feature bits	106
5.3.4	Device configuration layout	106
5.3.4.1	Legacy Interface: Device configuration layout	106
5.3.5	Device Initialization	106
5.3.5.1	Device Requirements: Device Initialization	107
5.3.6	Device Operation	107
5.3.6.1	Driver Requirements: Device Operation	107
5.3.6.2	Multiport Device Operation	107
5.3.6.2.1	Device Requirements: Multiport Device Operation	108
5.3.6.2.2	Driver Requirements: Multiport Device Operation	108



5.3.6.3	Legacy Interface: Device Operation	108
5.3.6.4	Legacy Interface: Framing Requirements	108
5.4	Entropy Device	108
5.4.1	Device ID	109
5.4.2	Virtqueues	109
5.4.3	Feature bits	109
5.4.4	Device configuration layout	109
5.4.5	Device Initialization	109
5.4.6	Device Operation	109
5.4.6.1	Driver Requirements: Device Operation	109
5.4.6.2	Device Requirements: Device Operation	109
5.5	Traditional Memory Balloon Device	109
5.5.1	Device ID	109
5.5.2	Virtqueues	109
5.5.3	Feature bits	110
5.5.3.1	Driver Requirements: Feature bits	110
5.5.3.2	Device Requirements: Feature bits	110
5.5.4	Device configuration layout	110
5.5.5	Device Initialization	111
5.5.6	Device Operation	111
5.5.6.1	Driver Requirements: Device Operation	112
5.5.6.2	Device Requirements: Device Operation	112
5.5.6.2.1	Legacy Interface: Device Operation	112
5.5.6.3	Memory Statistics	113
5.5.6.3.1	Driver Requirements: Memory Statistics	113
5.5.6.3.2	Device Requirements: Memory Statistics	113
5.5.6.3.3	Legacy Interface: Memory Statistics	114
5.5.6.4	Memory Statistics Tags	114
5.5.6.5	Free Page Hinting	114
5.5.6.5.1	Driver Requirements: Free Page Hinting	115
5.5.6.5.2	Device Requirements: Free Page Hinting	115
5.5.6.5.3	Legacy Interface: Free Page Hinting	116
5.5.6.6	Page Poison	116
5.5.6.6.1	Driver Requirements: Page Poison	116
5.5.6.6.2	Device Requirements: Page Poison	116
5.5.6.7	Free Page Reporting	116
5.5.6.7.1	Driver Requirements: Free Page Reporting	117
5.5.6.7.2	Device Requirements: Free Page Reporting	117
5.6	SCSI Host Device	117
5.6.1	Device ID	117
5.6.2	Virtqueues	117
5.6.3	Feature bits	118
5.6.4	Device configuration layout	118
5.6.4.1	Driver Requirements: Device configuration layout	118
5.6.4.2	Device Requirements: Device configuration layout	118
5.6.4.3	Legacy Interface: Device configuration layout	119
5.6.5	Device Requirements: Device Initialization	119
5.6.6	Device Operation	119
5.6.6.0.1	Legacy Interface: Device Operation	119
5.6.6.1	Device Operation: Request Queues	119
5.6.6.1.1	Device Requirements: Device Operation: Request Queues	120
5.6.6.1.2	Driver Requirements: Device Operation: Request Queues	121
5.6.6.1.3	Legacy Interface: Device Operation: Request Queues	121
5.6.6.2	Device Operation: controlq	121
5.6.6.2.1	Legacy Interface: Device Operation: controlq	123
5.6.6.3	Device Operation: eventq	123
5.6.6.3.1	Driver Requirements: Device Operation: eventq	125

	5.6.6.3.2	Device Requirements: Device Operation: eventq	125
	5.6.6.3.3	Legacy Interface: Device Operation: eventq	125
	5.6.6.4	Legacy Interface: Framing Requirements	125
5.7		GPU Device	126
	5.7.1	Device ID	126
	5.7.2	Virtqueues	126
	5.7.3	Feature bits	126
	5.7.4	Device configuration layout	126
	5.7.4.1	Device configuration fields	126
	5.7.4.2	Events	127
	5.7.5	Device Requirements: Device Initialization	127
	5.7.6	Device Operation	127
	5.7.6.1	Device Operation: Create a framebuffer and configure scanout	127
	5.7.6.2	Device Operation: Update a framebuffer and scanout	127
	5.7.6.3	Device Operation: Using pageflip	127
	5.7.6.4	Device Operation: Multihead setup	128
	5.7.6.5	Device Requirements: Device Operation: Command lifecycle and fencing	128
	5.7.6.6	Device Operation: Configure mouse cursor	128
	5.7.6.7	Device Operation: Request header	128
	5.7.6.8	Device Operation: controlq	129
	5.7.6.9	Device Operation: controlq (3d)	134
	5.7.6.10	Device Operation: cursorq	135
	5.7.7	VGA Compatibility	136
5.8		Input Device	136
	5.8.1	Device ID	136
	5.8.2	Virtqueues	136
	5.8.3	Feature bits	136
	5.8.4	Device configuration layout	136
	5.8.5	Device Initialization	137
	5.8.5.1	Driver Requirements: Device Initialization	138
	5.8.5.2	Device Requirements: Device Initialization	138
	5.8.6	Device Operation	138
	5.8.6.1	Driver Requirements: Device Operation	138
	5.8.6.2	Device Requirements: Device Operation	138
5.9		Crypto Device	138
	5.9.1	Device ID	138
	5.9.2	Virtqueues	139
	5.9.3	Feature bits	139
	5.9.3.1	Feature bit requirements	139
	5.9.4	Supported crypto services	139
	5.9.4.1	CIPHER services	139
	5.9.4.2	HASH services	140
	5.9.4.3	MAC services	140
	5.9.4.4	AEAD services	141
	5.9.5	Device configuration layout	141
	5.9.5.1	Device Requirements: Device configuration layout	142
	5.9.5.2	Driver Requirements: Device configuration layout	142
	5.9.6	Device Initialization	142
	5.9.6.1	Driver Requirements: Device Initialization	142
	5.9.7	Device Operation	142
	5.9.7.1	Operation Status	143
	5.9.7.2	Control Virtqueue	143
	5.9.7.2.1	Session operation	144
	5.9.7.3	Data Virtqueue	148
	5.9.7.4	HASH Service Operation	150
	5.9.7.4.1	Driver Requirements: HASH Service Operation	150
	5.9.7.4.2	Device Requirements: HASH Service Operation	150

5.9.7.5	MAC Service Operation . . . . .	151
5.9.7.5.1	Driver Requirements: MAC Service Operation . . . . .	152
5.9.7.5.2	Device Requirements: MAC Service Operation . . . . .	152
5.9.7.6	Symmetric algorithms Operation . . . . .	152
5.9.7.6.1	Driver Requirements: Symmetric algorithms Operation . . . . .	156
5.9.7.6.2	Device Requirements: Symmetric algorithms Operation . . . . .	156
5.9.7.7	AEAD Service Operation . . . . .	156
5.9.7.7.1	Driver Requirements: AEAD Service Operation . . . . .	158
5.9.7.7.2	Device Requirements: AEAD Service Operation . . . . .	158
5.10	Socket Device . . . . .	158
5.10.1	Device ID . . . . .	158
5.10.2	Virtqueues . . . . .	158
5.10.3	Feature bits . . . . .	159
5.10.4	Device configuration layout . . . . .	159
5.10.5	Device Initialization . . . . .	159
5.10.6	Device Operation . . . . .	159
5.10.6.1	Virtqueue Flow Control . . . . .	160
5.10.6.1.1	Driver Requirements: Device Operation: Virtqueue Flow Control . . . . .	160
5.10.6.1.2	Device Requirements: Device Operation: Virtqueue Flow Control . . . . .	160
5.10.6.2	Addressing . . . . .	160
5.10.6.3	Buffer Space Management . . . . .	160
5.10.6.3.1	Driver Requirements: Device Operation: Buffer Space Management . . . . .	161
5.10.6.3.2	Device Requirements: Device Operation: Buffer Space Management . . . . .	161
5.10.6.4	Receive and Transmit . . . . .	161
5.10.6.4.1	Driver Requirements: Device Operation: Receive and Transmit . . . . .	161
5.10.6.4.2	Device Requirements: Device Operation: Receive and Transmit . . . . .	161
5.10.6.5	Stream Sockets . . . . .	161
5.10.6.6	Seqpacket Sockets . . . . .	162
5.10.6.6.1	Message and record boundaries . . . . .	162
5.10.6.7	Device Events . . . . .	162
5.10.6.7.1	Driver Requirements: Device Operation: Device Events . . . . .	162
5.11	File System Device . . . . .	163
5.11.1	Device ID . . . . .	163
5.11.2	Virtqueues . . . . .	163
5.11.3	Feature bits . . . . .	163
5.11.4	Device configuration layout . . . . .	163
5.11.4.1	Driver Requirements: Device configuration layout . . . . .	163
5.11.4.2	Device Requirements: Device configuration layout . . . . .	164
5.11.5	Device Initialization . . . . .	164
5.11.6	Device Operation . . . . .	164
5.11.6.1	Device Operation: Request Queues . . . . .	164
5.11.6.2	Device Operation: High Priority Queue . . . . .	165
5.11.6.2.1	Device Requirements: Device Operation: High Priority Queue . . . . .	165
5.11.6.2.2	Driver Requirements: Device Operation: High Priority Queue . . . . .	165
5.11.6.3	Device Operation: Notification Queue . . . . .	165
5.11.6.3.1	Driver Requirements: Device Operation: Notification Queue . . . . .	166
5.11.6.4	Device Operation: DAX Window . . . . .	166
5.11.6.4.1	Device Requirements: Device Operation: DAX Window . . . . .	166
5.11.6.4.2	Driver Requirements: Device Operation: DAX Window . . . . .	166
5.11.6.5	Security Considerations . . . . .	167
5.11.6.6	Live migration considerations . . . . .	167
5.12	RPMB Device . . . . .	167
5.12.1	Device ID . . . . .	167
5.12.2	Virtqueues . . . . .	167
5.12.3	Feature bits . . . . .	168
5.12.4	Device configuration layout . . . . .	168
5.12.5	Device Requirements: Device Initialization . . . . .	168

5.12.6	Device Operation	168
5.12.6.1	Device Operation: Request Queue	169
5.12.6.1.1	Device Requirements: Device Operation: Program Key	169
5.12.6.1.2	Device Requirements: Device Operation: Get Write Counter	170
5.12.6.1.3	Device Requirements: Device Operation: Data Write	170
5.12.6.1.4	Device Requirements: Device Operation: Data Read	170
5.12.6.1.5	Device Requirements: Device Operation: Result Read	171
5.12.6.2	Driver Requirements: Device Operation	171
5.12.6.3	Device Requirements: Device Operation	171
5.13	IOMMU device	171
5.13.1	Device ID	172
5.13.2	Virtqueues	172
5.13.3	Feature bits	172
5.13.3.1	Driver Requirements: Feature bits	172
5.13.3.2	Device Requirements: Feature bits	172
5.13.4	Device configuration layout	172
5.13.4.1	Driver Requirements: Device configuration layout	173
5.13.4.2	Device Requirements: Device configuration layout	173
5.13.5	Device initialization	173
5.13.5.1	Driver Requirements: Device Initialization	173
5.13.6	Device operations	173
5.13.6.1	Driver Requirements: Device operations	174
5.13.6.2	Device Requirements: Device operations	175
5.13.6.3	ATTACH request	175
5.13.6.3.1	Driver Requirements: ATTACH request	175
5.13.6.3.2	Device Requirements: ATTACH request	176
5.13.6.4	DETACH request	176
5.13.6.4.1	Driver Requirements: DETACH request	176
5.13.6.4.2	Device Requirements: DETACH request	176
5.13.6.5	MAP request	177
5.13.6.5.1	Driver Requirements: MAP request	177
5.13.6.5.2	Device Requirements: MAP request	178
5.13.6.6	UNMAP request	178
5.13.6.6.1	Driver Requirements: UNMAP request	179
5.13.6.6.2	Device Requirements: UNMAP request	179
5.13.6.7	PROBE request	179
5.13.6.7.1	Driver Requirements: PROBE request	180
5.13.6.7.2	Device Requirements: PROBE request	180
5.13.6.8	PROBE properties	180
5.13.6.8.1	Property RESV_MEM	180
5.13.6.9	Fault reporting	181
5.13.6.9.1	Driver Requirements: Fault reporting	182
5.13.6.9.2	Device Requirements: Fault reporting	182
5.14	Sound Device	182
5.14.1	Device ID	182
5.14.2	Virtqueues	183
5.14.3	Feature Bits	183
5.14.4	Device Configuration Layout	183
5.14.5	Device Initialization	183
5.14.5.1	Driver Requirements: Device Initialization	183
5.14.6	Device Operation	183
5.14.6.1	Item Information Request	185
5.14.6.2	Driver Requirements: Item Information Request	185
5.14.6.3	Relationships with the High Definition Audio Specification	185
5.14.6.4	Jack Control Messages	186
5.14.6.4.1	VIRTIO_SND_R_JACK_INFO	186
5.14.6.4.2	VIRTIO_SND_R_JACK_REMAP	186

5.14.6.5	Jack Notifications	187
5.14.6.6	PCM Control Messages	187
5.14.6.6.1	PCM Command Lifecycle	187
5.14.6.6.2	VIRTIO_SND_R_PCM_INFO	188
5.14.6.6.3	VIRTIO_SND_R_PCM_SET_PARAMS	189
5.14.6.6.4	VIRTIO_SND_R_PCM_PREPARE	190
5.14.6.6.5	VIRTIO_SND_R_PCM_RELEASE	190
5.14.6.6.6	VIRTIO_SND_R_PCM_START	190
5.14.6.6.7	VIRTIO_SND_R_PCM_STOP	190
5.14.6.7	PCM Notifications	191
5.14.6.8	PCM I/O Messages	191
5.14.6.8.1	Output Stream	191
5.14.6.8.2	Input Stream	192
5.14.6.9	Channel Map Control Messages	192
5.14.6.9.1	VIRTIO_SND_R_CHMAP_INFO	192
5.15	Memory Device	193
5.15.1	Device ID	194
5.15.2	Virtqueues	194
5.15.3	Feature bits	194
5.15.4	Device configuration layout	194
5.15.4.1	Driver Requirements: Device configuration layout	194
5.15.4.2	Device Requirements: Device configuration layout	195
5.15.5	Device Initialization	195
5.15.5.1	Driver Requirements: Device Initialization	195
5.15.5.2	Device Requirements: Device Initialization	195
5.15.6	Device Operation	195
5.15.6.1	Driver Requirements: Device Operation	196
5.15.6.2	Device Requirements: Device Operation	197
5.15.6.3	PLUG request	197
5.15.6.3.1	Driver Requirements: PLUG request	197
5.15.6.3.2	Device Requirements: PLUG request	197
5.15.6.4	UNPLUG request	198
5.15.6.4.1	Driver Requirements: UNPLUG request	198
5.15.6.4.2	Device Requirements: UNPLUG request	198
5.15.6.5	UNPLUG ALL request	198
5.15.6.5.1	Driver Requirements: UNPLUG request	198
5.15.6.5.2	Device Requirements: UNPLUG request	198
5.15.6.6	STATE request	199
5.15.6.6.1	Driver Requirements: STATE request	199
5.15.6.6.2	Device Requirements: STATE request	199
5.16	I2C Adapter Device	199
5.16.1	Device ID	200
5.16.2	Virtqueues	200
5.16.3	Feature bits	200
5.16.4	Device configuration layout	200
5.16.5	Device Initialization	200
5.16.6	Device Operation	200
5.16.6.1	Device Operation: Request Queue	200
5.16.6.2	Device Operation: Operation Status	201
5.16.6.3	Driver Requirements: Device Operation	201
5.16.6.4	Device Requirements: Device Operation	202
5.17	SCMI Device	202
5.17.1	Device ID	202
5.17.2	Virtqueues	202
5.17.3	Feature bits	202
5.17.3.1	Device Requirements: Feature bits	203
5.17.4	Device configuration layout	203

5.17.5	Device Initialization	203
5.17.6	Device Operation	203
5.17.6.1	cmdq Operation	203
5.17.6.1.1	Device Requirements: cmdq Operation	204
5.17.6.1.2	Driver Requirements: cmdq Operation	204
5.17.6.2	Setting Up eventq Buffers	204
5.17.6.2.1	Driver Requirements: Setting Up eventq Buffers	204
5.17.6.3	eventq Operation	204
5.17.6.3.1	Device Requirements: eventq Operation	205
5.17.6.4	Shared Memory Operation	205
5.17.6.4.1	Device Requirements: Shared Memory Operation	205
5.18	GPIO Device	205
5.18.1	Device ID	205
5.18.2	Virtqueues	205
5.18.3	Feature bits	206
5.18.4	Device configuration layout	206
5.18.5	Device Initialization	206
5.18.6	Device Operation: requestq	206
5.18.6.1	requestq Operation: Get Line Names	207
5.18.6.2	requestq Operation: Get Direction	208
5.18.6.3	requestq Operation: Set Direction	208
5.18.6.4	requestq Operation: Get Value	208
5.18.6.5	requestq Operation: Set Value	208
5.18.6.6	requestq Operation: Set IRQ Type	209
5.18.6.7	requestq Operation: Message Flow	209
5.18.6.8	Driver Requirements: requestq Operation	210
5.18.6.9	Device Requirements: requestq Operation	210
5.18.7	Device Operation: eventq	210
5.18.7.1	eventq Operation: Message Flow	211
5.18.7.2	Driver Requirements: eventq Operation	212
5.18.7.3	Device Requirements: eventq Operation	212
5.19	PMEM Device	212
5.19.1	Device ID	212
5.19.2	Virtqueues	212
5.19.3	Feature bits	212
5.19.4	Device configuration layout	212
5.19.5	Device Initialization	213
5.19.5.1	Device Requirements: Device Initialization	213
5.19.5.2	Driver Requirements: Device Initialization	213
5.19.6	Driver Operations	213
5.19.7	Device Operations	213
5.19.7.1	Device Requirements: Device Operation: Virtqueue flush	213
5.19.7.2	Device Operations	213
5.19.7.3	Device Requirements: Device Operation: Virtqueue return	214
5.19.8	Possible security implications	214
5.19.9	Countermeasures	214
5.19.9.1	With SHARED mapping	214
5.19.9.2	With PRIVATE mapping	214
5.19.9.3	Workload specific mapping	214
5.19.9.4	Prevent cache eviction	214
6	Reserved Feature Bits	215
6.1	Driver Requirements: Reserved Feature Bits	216
6.2	Device Requirements: Reserved Feature Bits	216
6.3	Legacy Interface: Reserved Feature Bits	216
7	Conformance	218

7.1	Conformance Targets	218
7.2	Clause 1: Driver Conformance	218
7.2.1	Clause 2: PCI Driver Conformance	219
7.2.2	Clause 3: MMIO Driver Conformance	219
7.2.3	Clause 4: Channel I/O Driver Conformance	219
7.2.4	Clause 5: Network Driver Conformance	219
7.2.5	Clause 6: Block Driver Conformance	220
7.2.6	Clause 7: Console Driver Conformance	220
7.2.7	Clause 8: Entropy Driver Conformance	220
7.2.8	Clause 9: Traditional Memory Balloon Driver Conformance	220
7.2.9	Clause 10: SCSI Host Driver Conformance	220
7.2.10	Clause 11: Input Driver Conformance	220
7.2.11	Clause 12: Crypto Driver Conformance	221
7.2.12	Clause 13: Socket Driver Conformance	221
7.2.13	Clause 14: File System Driver Conformance	221
7.2.14	Clause 15: RPMB Driver Conformance	221
7.2.15	Clause 16: IOMMU Driver Conformance	221
7.2.16	Clause 17: Sound Driver Conformance	222
7.2.17	Clause 18: Memory Driver Conformance	222
7.2.18	Clause 19: I2C Adapter Driver Conformance	222
7.2.19	Clause 20: SCMI Driver Conformance	222
7.2.20	Clause 21: GPIO Driver Conformance	222
7.2.21	Clause 22: PMEM Driver Conformance	222
7.3	Clause 23: Device Conformance	223
7.3.1	Clause 24: PCI Device Conformance	223
7.3.2	Clause 25: MMIO Device Conformance	223
7.3.3	Clause 26: Channel I/O Device Conformance	224
7.3.4	Clause 27: Network Device Conformance	224
7.3.5	Clause 28: Block Device Conformance	224
7.3.6	Clause 29: Console Device Conformance	224
7.3.7	Clause 30: Entropy Device Conformance	224
7.3.8	Clause 31: Traditional Memory Balloon Device Conformance	225
7.3.9	Clause 32: SCSI Host Device Conformance	225
7.3.10	Clause 33: Input Device Conformance	225
7.3.11	Clause 34: Crypto Device Conformance	225
7.3.12	Clause 35: Socket Device Conformance	225
7.3.13	Clause 36: File System Device Conformance	225
7.3.14	Clause 37: RPMB Device Conformance	226
7.3.15	Clause 38: IOMMU Device Conformance	226
7.3.16	Clause 39: Sound Device Conformance	226
7.3.17	Clause 40: Memory Device Conformance	226
7.3.18	Clause 41: I2C Adapter Device Conformance	227
7.3.19	Clause 42: SCMI Device Conformance	227
7.3.20	Clause 43: GPIO Device Conformance	227
7.3.21	Clause 44: PMEM Device Conformance	227
7.4	Clause 45: Legacy Interface: Transitional Device and Transitional Driver Conformance	227
A	<b>virtio_queue.h</b>	<b>229</b>
B	<b>Creating New Device Types</b>	<b>231</b>
B.1	How Many Virtqueues?	231
B.2	What Device Configuration Space Layout?	231
B.3	What Device Number?	231
B.4	How many MSI-X vectors? (for PCI)	231
B.5	Device Improvements	231
C	<b>Acknowledgements</b>	<b>232</b>





---

# 1 Introduction

This document describes the specifications of the “virtio” family of devices. These devices are found in virtual environments, yet by design they look like physical devices to the guest within the virtual machine - and this document treats them as such. This similarity allows the guest to use standard drivers and discovery mechanisms.

The purpose of virtio and this specification is that virtual environments and guests should have a straightforward, efficient, standard and extensible mechanism for virtual devices, rather than boutique per-environment or per-OS mechanisms.

**Straightforward:** Virtio devices use normal bus mechanisms of interrupts and DMA which should be familiar to any device driver author. There is no exotic page-flipping or COW mechanism: it’s just a normal device.<sup>1</sup>

**Efficient:** Virtio devices consist of rings of descriptors for both input and output, which are neatly laid out to avoid cache effects from both driver and device writing to the same cache lines.

**Standard:** Virtio makes no assumptions about the environment in which it operates, beyond supporting the bus to which device is attached. In this specification, virtio devices are implemented over MMIO, Channel I/O and PCI bus transports <sup>2</sup>, earlier drafts have been implemented on other buses not included here.

**Extensible:** Virtio devices contain feature bits which are acknowledged by the guest operating system during device setup. This allows forwards and backwards compatibility: the device offers all the features it knows about, and the driver acknowledges those it understands and wishes to use.

## 1.1 Normative References

- |                          |   |
|--------------------------|---|
| <b>[RFC2119]</b>         | Bradner S., “Key words for use in RFCs to Indicate Requirement Levels”, BCP 14, RFC 2119, March 1997.<br><a href="http://www.ietf.org/rfc/rfc2119.txt">http://www.ietf.org/rfc/rfc2119.txt</a>  |
| <b>[RFC4122]</b>         | Leach, P., Mealling, M., and R. Salz, “A Universally Unique IDentifier (UUID) URN Namespace”, RFC 4122, DOI 10.17487/RFC4122, July 2005.<br><a href="http://www.ietf.org/rfc/rfc4122.txt">http://www.ietf.org/rfc/rfc4122.txt</a>                                   |
| <b>[S390 PoP]</b>        | z/Architecture Principles of Operation, IBM Publication SA22-7832,<br><a href="http://publibfi.boulder.ibm.com/epubs/pdf/dz9zr009.pdf">http://publibfi.boulder.ibm.com/epubs/pdf/dz9zr009.pdf</a> , and any future revisions  |
| <b>[S390 Common I/O]</b> | ESA/390 Common I/O-Device and Self-Description, IBM Publication SA22-7204,<br><a href="http://publibfp.dhe.ibm.com/cgi-bin/bookmgr/BOOKS/dz9ar501/CCONTENTS">http://publibfp.dhe.ibm.com/cgi-bin/bookmgr/BOOKS/dz9ar501/CCONTENTS</a> ,<br>and any future revisions |
| <b>[PCI]</b>             | Conventional PCI Specifications,<br><a href="http://www.pcisig.com/specifications/conventional/">http://www.pcisig.com/specifications/conventional/</a> , PCI-SIG   |
| <b>[PCIe]</b>            | PCI Express Specifications<br><a href="http://www.pcisig.com/specifications/pciexpress/">http://www.pcisig.com/specifications/pciexpress/</a> , PCI-SIG   |

---

<sup>1</sup>This lack of page-sharing implies that the implementation of the device (e.g. the hypervisor or host) needs full access to the guest memory. Communication with untrusted parties (i.e. inter-guest communication) requires copying.

<sup>2</sup>The Linux implementation further separates the virtio transport code from the specific virtio drivers: these drivers are shared between different transports.

[IEEE 802]	IEEE Standard for Local and Metropolitan Area Networks: Overview and Architecture, <a href="http://www.ieee802.org/">http://www.ieee802.org/</a> , IEEE
[SAM]	SCSI Architectural Model, <a href="http://www.t10.org/cgi-bin/ac.pl?t=f&amp;f=sam4r05.pdf">http://www.t10.org/cgi-bin/ac.pl?t=f&amp;f=sam4r05.pdf</a>
[SCSI MMC]	SCSI Multimedia Commands, <a href="http://www.t10.org/cgi-bin/ac.pl?t=f&amp;f=mmc6r00.pdf">http://www.t10.org/cgi-bin/ac.pl?t=f&amp;f=mmc6r00.pdf</a>
[FUSE]	Linux FUSE interface, <a href="https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/include/uapi/linux/fuse.h">https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/include/uapi/linux/fuse.h</a>
[eMMC]	eMMC Electrical Standard (5.1), JESD84-B51, <a href="http://www.jedec.org/sites/default/files/docs/JESD84-B51.pdf">http://www.jedec.org/sites/default/files/docs/JESD84-B51.pdf</a>
[HDA]	High Definition Audio Specification, <a href="https://www.intel.com/content/dam/www/public/us/en/documents/product-specifications/high-definition-audio-specification.pdf">https://www.intel.com/content/dam/www/public/us/en/documents/product-specifications/high-definition-audio-specification.pdf</a>
[I2C]	I2C-bus specification and user manual, <a href="https://www.nxp.com/docs/en/user-guide/UM10204.pdf">https://www.nxp.com/docs/en/user-guide/UM10204.pdf</a>
[SCMI]	Arm System Control and Management Interface, DEN0056, <a href="https://developer.arm.com/docs/den0056/c">https://developer.arm.com/docs/den0056/c</a> , version C and any future revisions

## 1.2 Non-Normative References

[Virtio PCI Draft]	Virtio PCI Draft Specification <a href="http://ozlabs.org/~rusty/virtio-spec/virtio-0.9.5.pdf">http://ozlabs.org/~rusty/virtio-spec/virtio-0.9.5.pdf</a>
--------------------	---

## 1.3 Terminology

The key words “MUST”, “MUST NOT”, “REQUIRED”, “SHALL”, “SHALL NOT”, “SHOULD”, “SHOULD NOT”, “RECOMMENDED”, “MAY”, and “OPTIONAL” in this document are to be interpreted as described in [RFC2119].

### 1.3.1 Legacy Interface: Terminology

Specification drafts preceding version 1.0 of this specification (e.g. see [Virtio PCI Draft]) defined a similar, but different interface between the driver and the device. Since these are widely deployed, this specification accommodates OPTIONAL features to simplify transition from these earlier draft interfaces.

Specifically devices and drivers MAY support:

**Legacy Interface** is an interface specified by an earlier draft of this specification (before 1.0)

**Legacy Device** is a device implemented before this specification was released, and implementing a legacy interface on the host side

**Legacy Driver** is a driver implemented before this specification was released, and implementing a legacy interface on the guest side

Legacy devices and legacy drivers are not compliant with this specification.

To simplify transition from these earlier draft interfaces, a device MAY implement:

**Transitional Device** a device supporting both drivers conforming to this specification, and allowing legacy drivers.

Similarly, a driver MAY implement:

**Transitional Driver** a driver supporting both devices conforming to this specification, and legacy devices.

**Note:** Legacy interfaces are not required; ie. don't implement them unless you have a need for backwards compatibility!

Devices or drivers with no legacy compatibility are referred to as non-transitional devices and drivers, respectively.

### 1.3.2 Transition from earlier specification drafts

For devices and drivers already implementing the legacy interface, some changes will have to be made to support this specification.

In this case, it might be beneficial for the reader to focus on sections tagged "Legacy Interface" in the section title. These highlight the changes made since the earlier drafts.

## 1.4 Structure Specifications

Many device and driver in-memory structure layouts are documented using the C struct syntax. All structures are assumed to be without additional padding. To stress this, cases where common C compilers are known to insert extra padding within structures are tagged using the GNU C `__attribute__((packed))` syntax.

For the integer data types used in the structure definitions, the following conventions are used:

**u8, u16, u32, u64** An unsigned integer of the specified length in bits.

**le16, le32, le64** An unsigned integer of the specified length in bits, in little-endian byte order.

**be16, be32, be64** An unsigned integer of the specified length in bits, in big-endian byte order.

Some of the fields to be defined in this specification don't start or don't end on a byte boundary. Such fields are called bit-fields. A set of bit-fields is always a sub-division of an integer typed field.

Bit-fields within integer fields are always listed in order, from the least significant to the most significant bit. The bit-fields are considered unsigned integers of the specified width with the next in significance relationship of the bits preserved.

For example:

```
struct S {
    be16 {
        A : 15;
        B : 1;
    } x;
    be16 y;
};
```

documents the value A stored in the low 15 bit of x and the value B stored in the high bit of x, the 16-bit integer x in turn stored using the big-endian byte order at the beginning of the structure S, and being followed immediately by an unsigned integer y stored in big-endian byte order at an offset of 2 bytes (16 bits) from the beginning of the structure.

Note that this notation somewhat resembles the C bitfield syntax but should not be naively converted to a bitfield notation for portable code: it matches the way bitfields are packed by C compilers on little-endian architectures but not the way bitfields are packed by C compilers on big-endian architectures.

Assuming that CPU\_TO\_BE16 converts a 16-bit integer from a native CPU to the big-endian byte order, the following is the equivalent portable C code to generate a value to be stored into x:

```
CPU_TO_BE16(B << 15 | A)
```

## 1.5 Constant Specifications

In many cases, numeric values used in the interface between the device and the driver are documented using the C `#define` and `/* */` comment syntax. Multiple related values are grouped together with a common name as a prefix, using `_` as a separator. Using `_XXX` as a suffix refers to all values in a group. For example:

```
/* Field Fld value A description */
#define VIRTIO_FLD_A      (1 << 0)
/* Field Fld value B description */
#define VIRTIO_FLD_B      (1 << 1)
```

documents two numeric values for a field *Fld*, with *Fld* having value 1 referring to *A* and *Fld* having value 2 referring to *B*. Note that `<<` refers to the shift-left operation.

Further, in this case `VIRTIO_FLD_A` and `VIRTIO_FLD_B` refer to values 1 and 2 of *Fld* respectively. Further, `VIRTIO_FLD_XXX` refers to either `VIRTIO_FLD_A` or `VIRTIO_FLD_B`.

---

## 2 Basic Facilities of a Virtio Device

A virtio device is discovered and identified by a bus-specific method (see the bus specific sections: [4.1 Virtio Over PCI Bus](#), [4.2 Virtio Over MMIO](#) and [4.3 Virtio Over Channel I/O](#)). Each device consists of the following parts:

- Device status field
- Feature bits
- Notifications
- Device Configuration space
- One or more virtqueues

### 2.1 Device Status Field

During device initialization by a driver, the driver follows the sequence of steps specified in [3.1](#).

The *device status* field provides a simple low-level indication of the completed steps of this sequence. It's most useful to imagine it hooked up to traffic lights on the console indicating the status of each device. The following bits are defined (listed below in the order in which they would be typically set):

**ACKNOWLEDGE (1)** Indicates that the guest OS has found the device and recognized it as a valid virtio device.

**DRIVER (2)** Indicates that the guest OS knows how to drive the device.

**Note:** There could be a significant (or infinite) delay before setting this bit. For example, under Linux, drivers can be loadable modules.

**FAILED (128)** Indicates that something went wrong in the guest, and it has given up on the device. This could be an internal error, or the driver didn't like the device for some reason, or even a fatal error during device operation.

**FEATURES\_OK (8)** Indicates that the driver has acknowledged all the features it understands, and feature negotiation is complete.

**DRIVER\_OK (4)** Indicates that the driver is set up and ready to drive the device.

**DEVICE\_NEEDS\_RESET (64)** Indicates that the device has experienced an error from which it can't recover.

The *device status* field starts out as 0, and is reinitialized to 0 by the device during reset.

#### 2.1.1 Driver Requirements: Device Status Field

The driver **MUST** update *device status*, setting bits to indicate the completed steps of the driver initialization sequence specified in [3.1](#). The driver **MUST NOT** clear a *device status* bit. If the driver sets the FAILED bit, the driver **MUST** later reset the device before attempting to re-initialize.

The driver **SHOULD NOT** rely on completion of operations of a device if **DEVICE\_NEEDS\_RESET** is set.

**Note:** For example, the driver can't assume requests in flight will be completed if **DEVICE\_NEEDS\_RESET** is set, nor can it assume that they have not been completed. A good implementation will try to recover by issuing a reset.

## 2.1.2 Device Requirements: Device Status Field

The device **MUST NOT** consume buffers or send any used buffer notifications to the driver before `DRIVER_OK`.

The device **SHOULD** set `DEVICE_NEEDS_RESET` when it enters an error state that a reset is needed. If `DRIVER_OK` is set, after it sets `DEVICE_NEEDS_RESET`, the device **MUST** send a device configuration change notification to the driver.

## 2.2 Feature Bits

Each virtio device offers all the features it understands. During device initialization, the driver reads this and tells the device the subset that it accepts. The only way to renegotiate is to reset the device.

This allows for forwards and backwards compatibility: if the device is enhanced with a new feature bit, older drivers will not write that feature bit back to the device. Similarly, if a driver is enhanced with a feature that the device doesn't support, it sees the new feature is not offered.

Feature bits are allocated as follows:

**0 to 23** Feature bits for the specific device type

**24 to 40** Feature bits reserved for extensions to the queue and feature negotiation mechanisms

**41 and above** Feature bits reserved for future extensions.

**Note:** For example, feature bit 0 for a network device (i.e. Device ID 1) indicates that the device supports checksumming of packets.

In particular, new fields in the device configuration space are indicated by offering a new feature bit.

### 2.2.1 Driver Requirements: Feature Bits

The driver **MUST NOT** accept a feature which the device did not offer, and **MUST NOT** accept a feature which requires another feature which was not accepted.

The driver **SHOULD** go into backwards compatibility mode if the device does not offer a feature it understands, otherwise **MUST** set the `FAILED device status` bit and cease initialization.

### 2.2.2 Device Requirements: Feature Bits

The device **MUST NOT** offer a feature which requires another feature which was not offered. The device **SHOULD** accept any valid subset of features the driver accepts, otherwise it **MUST** fail to set the `FEATURES_OK device status` bit when the driver writes it.

If a device has successfully negotiated a set of features at least once (by accepting the `FEATURES_OK device status` bit during device initialization), then it **SHOULD NOT** fail re-negotiation of the same set of features after a device or system reset. Failure to do so would interfere with resuming from suspend and error recovery.

### 2.2.3 Legacy Interface: A Note on Feature Bits

Transitional Drivers **MUST** detect Legacy Devices by detecting that the feature bit `VIRTIO_F_VERSION_1` is not offered. Transitional devices **MUST** detect Legacy drivers by detecting that `VIRTIO_F_VERSION_1` has not been acknowledged by the driver.

In this case device is used through the legacy interface.

Legacy interface support is **OPTIONAL**. Thus, both transitional and non-transitional devices and drivers are compliant with this specification.

Requirements pertaining to transitional devices and drivers is contained in sections named 'Legacy Interface' like this one.



When device is used through the legacy interface, transitional devices and transitional drivers **MUST** operate according to the requirements documented within these legacy interface sections. Specification text within these sections generally does not apply to non-transitional devices.

## 2.3 Notifications

The notion of sending a notification (driver to device or device to driver) plays an important role in this specification. The *modus operandi* of the notifications is transport specific.

There are three types of notifications:

- configuration change notification
- available buffer notification
- used buffer notification.

Configuration change notifications and used buffer notifications are sent by the device, the recipient is the driver. A configuration change notification indicates that the device configuration space has changed; a used buffer notification indicates that a buffer may have been made used on the virtqueue designated by the notification.

Available buffer notifications are sent by the driver, the recipient is the device. This type of notification indicates that a buffer may have been made available on the virtqueue designated by the notification.

The semantics, the transport-specific implementations, and other important aspects of the different notifications are specified in detail in the following chapters.

Most transports implement notifications sent by the device to the driver using interrupts. Therefore, in previous versions of this specification, these notifications were often called interrupts. Some names defined in this specification still retain this interrupt terminology. Occasionally, the term event is used to refer to a notification or a receipt of a notification.

## 2.4 Device Reset

The driver may want to initiate a device reset at various times; notably, it is required to do so during device initialization and device cleanup.

The mechanism used by the driver to initiate the reset is transport specific.

### 2.4.1 Device Requirements: Device Reset

A device **MUST** reinitialize *device status* to 0 after receiving a reset.

A device **MUST NOT** send notifications or interact with the queues after indicating completion of the reset by reinitializing *device status* to 0, until the driver re-initializes the device.

### 2.4.2 Driver Requirements: Device Reset

The driver **SHOULD** consider a driver-initiated reset complete when it reads *device status* as 0.

## 2.5 Device Configuration Space

Device configuration space is generally used for rarely-changing or initialization-time parameters. Where configuration fields are optional, their existence is indicated by feature bits: Future versions of this specification will likely extend the device configuration space by adding extra fields at the tail.

**Note:** The device configuration space uses the little-endian format for multi-byte fields.

Each transport also provides a generation count for the device configuration space, which will change whenever there is a possibility that two accesses to the device configuration space can see different versions of that space.

### 2.5.1 Driver Requirements: Device Configuration Space

Drivers **MUST NOT** assume reads from fields greater than 32 bits wide are atomic, nor are reads from multiple fields: drivers **SHOULD** read device configuration space fields like so:

```
u32 before, after;
do {
    before = get_config_generation(device);
    // read config entry/entries.
    after = get_config_generation(device);
} while (after != before);
```

For optional configuration space fields, the driver **MUST** check that the corresponding feature is offered before accessing that part of the configuration space.

**Note:** See section 3.1 for details on feature negotiation.

Drivers **MUST NOT** limit structure size and device configuration space size. Instead, drivers **SHOULD** only check that device configuration space is *large enough* to contain the fields necessary for device operation.

**Note:** For example, if the specification states that device configuration space 'includes a single 8-bit field' drivers should understand this to mean that the device configuration space might also include an arbitrary amount of tail padding, and accept any device configuration space size equal to or greater than the specified 8-bit size.

### 2.5.2 Device Requirements: Device Configuration Space

The device **MUST** allow reading of any device-specific configuration field before `FEATURES_OK` is set by the driver. This includes fields which are conditional on feature bits, as long as those feature bits are offered by the device.

### 2.5.3 Legacy Interface: A Note on Device Configuration Space endian-ness

Note that for legacy interfaces, device configuration space is generally the guest's native endian, rather than PCI's little-endian. The correct endian-ness is documented for each device.

### 2.5.4 Legacy Interface: Device Configuration Space

Legacy devices did not have a configuration generation field, thus are susceptible to race conditions if configuration is updated. This affects the block *capacity* (see 5.2.4) and network *mac* (see 5.1.4) fields; when using the legacy interface, drivers **SHOULD** read these fields multiple times until two reads generate a consistent result.

## 2.6 Virtqueues

The mechanism for bulk data transport on virtio devices is pretentiously called a virtqueue. Each device can have zero or more virtqueues<sup>1</sup>.

Driver makes requests available to device by adding an available buffer to the queue, i.e., adding a buffer describing the request to a virtqueue, and optionally triggering a driver event, i.e., sending an available buffer notification to the device.

Device executes the requests and - when complete - adds a used buffer to the queue, i.e., lets the driver know by marking the buffer as used. Device can then trigger a device event, i.e., send a used buffer notification to the driver.

<sup>1</sup>For example, the simplest network device has one virtqueue for transmit and one for receive.

Device reports the number of bytes it has written to memory for each buffer it uses. This is referred to as “used length”.

Device is not generally required to use buffers in the same order in which they have been made available by the driver.

Some devices always use descriptors in the same order in which they have been made available. These devices can offer the VIRTIO\_F\_IN\_ORDER feature. If negotiated, this knowledge might allow optimizations or simplify driver and/or device code.

Each virtqueue can consist of up to 3 parts:

- Descriptor Area - used for describing buffers
- Driver Area - extra data supplied by driver to the device
- Device Area - extra data supplied by device to driver

**Note:** Note that previous versions of this spec used different names for these parts (following 2.7):

- Descriptor Table - for the Descriptor Area
- Available Ring - for the Driver Area
- Used Ring - for the Device Area

Two formats are supported: Split Virtqueues (see 2.7 Split Virtqueues) and Packed Virtqueues (see 2.8 Packed Virtqueues).

Every driver and device supports either the Packed or the Split Virtqueue format, or both.

## 2.6.1 Virtqueue Reset

When VIRTIO\_F\_RING\_RESET is negotiated, the driver can reset a virtqueue individually. The way to reset the virtqueue is transport specific.

Virtqueue reset is divided into two parts. The driver first resets a queue and can afterwards optionally re-enable it.

### 2.6.1.1 Virtqueue Reset

#### 2.6.1.1.1 Device Requirements: Virtqueue Reset

After a queue has been reset by the driver, the device MUST NOT execute any requests from that virtqueue, or notify the driver for it.

The device MUST reset any state of a virtqueue to the default state, including the available state and the used state.

#### 2.6.1.1.2 Driver Requirements: Virtqueue Reset

After the driver tells the device to reset a queue, the driver MUST verify that the queue has actually been reset.

After the queue has been successfully reset, the driver MAY release any resource associated with that virtqueue.

### 2.6.1.2 Virtqueue Re-enable

This process is the same as the initialization process of a single queue during the initialization of the entire device.

#### 2.6.1.2.1 Device Requirements: Virtqueue Re-enable

The device MUST observe any queue configuration that may have been changed by the driver, like the maximum queue size.

### 2.6.1.2.2 Driver Requirements: Virtqueue Re-enable

When re-enabling a queue, the driver MUST configure the queue resources as during initial virtqueue discovery, but optionally with different parameters.

## 2.7 Split Virtqueues

The split virtqueue format was the only format supported by the version 1.0 (and earlier) of this standard.

The split virtqueue format separates the virtqueue into several parts, where each part is write-able by either the driver or the device, but not both. Multiple parts and/or locations within a part need to be updated when making a buffer available and when marking it as used.

Each queue has a 16-bit queue size parameter, which sets the number of entries and implies the total size of the queue.

Each virtqueue consists of three parts:

- Descriptor Table - occupies the Descriptor Area
- Available Ring - occupies the Driver Area
- Used Ring - occupies the Device Area

where each part is physically-contiguous in guest memory, and has different alignment requirements.

The memory alignment and size requirements, in bytes, of each part of the virtqueue are summarized in the following table:

Virtqueue Part	Alignment	Size
Descriptor Table	16	16*(Queue Size)
Available Ring	2	6 + 2*(Queue Size)
Used Ring	4	6 + 8*(Queue Size)

The Alignment column gives the minimum alignment for each part of the virtqueue.

The Size column gives the total number of bytes for each part of the virtqueue.

Queue Size corresponds to the maximum number of buffers in the virtqueue<sup>2</sup>. Queue Size value is always a power of 2. The maximum Queue Size value is 32768. This value is specified in a bus-specific way.

When the driver wants to send a buffer to the device, it fills in a slot in the descriptor table (or chains several together), and writes the descriptor index into the available ring. It then notifies the device. When the device has finished a buffer, it writes the descriptor index into the used ring, and sends a used buffer notification.

### 2.7.1 Driver Requirements: Virtqueues

The driver MUST ensure that the physical address of the first byte of each virtqueue part is a multiple of the specified alignment value in the above table.

### 2.7.2 Legacy Interfaces: A Note on Virtqueue Layout

For Legacy Interfaces, several additional restrictions are placed on the virtqueue layout:

Each virtqueue occupies two or more physically-contiguous pages (usually defined as 4096 bytes, but depending on the transport; henceforth referred to as Queue Align) and consists of three parts:

Descriptor Table	Available Ring (...padding...)	Used Ring
------------------	--------------------------------	-----------

The bus-specific Queue Size field controls the total number of bytes for the virtqueue. When using the legacy interface, the transitional driver MUST retrieve the Queue Size field from the device and MUST allocate the

<sup>2</sup>For example, if Queue Size is 4 then at most 4 buffers can be queued at any given time.

total number of bytes for the virtqueue according to the following formula (Queue Align given in qalign and Queue Size given in qsz):

```
#define ALIGN(x) (((x) + qalign) & ~qalign)
static inline unsigned virtq_size(unsigned int qsz)
{
    return ALIGN(sizeof(struct virtq_desc)*qsz + sizeof(u16)*(3 + qsz))
        + ALIGN(sizeof(u16)*3 + sizeof(struct virtq_used_elem)*qsz);
}
```

This wastes some space with padding. When using the legacy interface, both transitional devices and drivers **MUST** use the following virtqueue layout structure to locate elements of the virtqueue:

```
struct virtq {
    // The actual descriptors (16 bytes each)
    struct virtq_desc desc[ Queue Size ];

    // A ring of available descriptor heads with free-running index.
    struct virtq_avail avail;

    // Padding to the next Queue Align boundary.
    u8 pad[ Padding ];

    // A ring of used descriptor heads with free-running index.
    struct virtq_used used;
};
```

### 2.7.3 Legacy Interfaces: A Note on Virtqueue Endianness

Note that when using the legacy interface, transitional devices and drivers **MUST** use the native endian of the guest as the endian of fields and in the virtqueue. This is opposed to little-endian for non-legacy interface as specified by this standard. It is assumed that the host is already aware of the guest endian.

### 2.7.4 Message Framing

The framing of messages with descriptors is independent of the contents of the buffers. For example, a network transmit buffer consists of a 12 byte header followed by the network packet. This could be most simply placed in the descriptor table as a 12 byte output descriptor followed by a 1514 byte output descriptor, but it could also consist of a single 1526 byte output descriptor in the case where the header and packet are adjacent, or even three or more descriptors (possibly with loss of efficiency in that case).

Note that, some device implementations have large-but-reasonable restrictions on total descriptor size (such as based on IOV\_MAX in the host OS). This has not been a problem in practice: little sympathy will be given to drivers which create unreasonably-sized descriptors such as by dividing a network packet into 1500 single-byte descriptors!

#### 2.7.4.1 Device Requirements: Message Framing

The device **MUST NOT** make assumptions about the particular arrangement of descriptors. The device **MAY** have a reasonable limit of descriptors it will allow in a chain.

#### 2.7.4.2 Driver Requirements: Message Framing

The driver **MUST** place any device-writable descriptor elements after any device-readable descriptor elements.

The driver **SHOULD NOT** use an excessive number of descriptors to describe a buffer.

#### 2.7.4.3 Legacy Interface: Message Framing

Regrettably, initial driver implementations used simple layouts, and devices came to rely on it, despite this specification wording. In addition, the specification for virtio\_blk SCSI commands required intuiting field lengths from frame boundaries (see [5.2.6.3 Legacy Interface: Device Operation](#))

Thus when using the legacy interface, the `VIRTIO_F_ANY_LAYOUT` feature indicates to both the device and the driver that no assumptions were made about framing. Requirements for transitional drivers when this is not negotiated are included in each device section.

### 2.7.5 The Virtqueue Descriptor Table

The descriptor table refers to the buffers the driver is using for the device. *addr* is a physical address, and the buffers can be chained via *next*. Each descriptor describes a buffer which is read-only for the device (“device-readable”) or write-only for the device (“device-writable”), but a chain of descriptors can contain both device-readable and device-writable buffers.

The actual contents of the memory offered to the device depends on the device type. Most common is to begin the data with a header (containing little-endian fields) for the device to read, and postfix it with a status trailer for the device to write.

```
struct virtq_desc {
    /* Address (guest-physical). */
    le64 addr;
    /* Length. */
    le32 len;

    /* This marks a buffer as continuing via the next field. */
#define VIRTQ_DESC_F_NEXT 1
    /* This marks a buffer as device write-only (otherwise device read-only). */
#define VIRTQ_DESC_F_WRITE 2
    /* This means the buffer contains a list of buffer descriptors. */
#define VIRTQ_DESC_F_INDIRECT 4
    /* The flags as indicated above. */
    le16 flags;
    /* Next field if flags & NEXT */
    le16 next;
};
```

The number of descriptors in the table is defined by the queue size for this virtqueue: this is the maximum possible descriptor chain length.

If `VIRTIO_F_IN_ORDER` has been negotiated, driver uses descriptors in ring order: starting from offset 0 in the table, and wrapping around at the end of the table.

**Note:** The legacy [\[Virtio PCI Draft\]](#) referred to this structure as `vring_desc`, and the constants as `VRING_DESC_F_NEXT`, etc, but the layout and values were identical.

#### 2.7.5.1 Device Requirements: The Virtqueue Descriptor Table

A device **MUST NOT** write to a device-readable buffer, and a device **SHOULD NOT** read a device-writable buffer (it **MAY** do so for debugging or diagnostic purposes). A device **MUST NOT** write to any descriptor table entry.

#### 2.7.5.2 Driver Requirements: The Virtqueue Descriptor Table

Drivers **MUST NOT** add a descriptor chain longer than  $2^{32}$  bytes in total; this implies that loops in the descriptor chain are forbidden!

If `VIRTIO_F_IN_ORDER` has been negotiated, and when making a descriptor with `VRING_DESC_F_NEXT` set in *flags* at offset *x* in the table available to the device, driver **MUST** set *next* to 0 for the last descriptor in the table (where  $x = \text{queue\_size} - 1$ ) and to  $x + 1$  for the rest of the descriptors.

#### 2.7.5.3 Indirect Descriptors

Some devices benefit by concurrently dispatching a large number of large requests. The `VIRTIO_F_INDIRECT_DESC` feature allows this (see [A virtio\\_queue.h](#)). To increase ring capacity the driver can store a table of indirect descriptors anywhere in memory, and insert a descriptor in main virtqueue (with *flags* & `VIRTQ_DESC_F_INDIRECT` on) that refers to memory buffer containing this indirect descriptor table; *addr* and *len* refer to the indirect table address and length in bytes, respectively.

The indirect table layout structure looks like this (*len* is the length of the descriptor that refers to this table, which is a variable, so this code won't compile):

```
struct indirect_descriptor_table {
    /* The actual descriptors (16 bytes each) */
    struct virtq_desc desc[len / 16];
};
```

The first indirect descriptor is located at start of the indirect descriptor table (index 0), additional indirect descriptors are chained by *next*. An indirect descriptor without a valid *next* (with *flags&VIRTQ\_DESC\_F\_NEXT* off) signals the end of the descriptor. A single indirect descriptor table can include both device-readable and device-writable descriptors.

If *VIRTIO\_F\_IN\_ORDER* has been negotiated, indirect descriptors use sequential indices, in-order: index 0 followed by index 1 followed by index 2, etc.

#### 2.7.5.3.1 Driver Requirements: Indirect Descriptors

The driver **MUST NOT** set the *VIRTQ\_DESC\_F\_INDIRECT* flag unless the *VIRTIO\_F\_INDIRECT\_DESC* feature was negotiated. The driver **MUST NOT** set the *VIRTQ\_DESC\_F\_INDIRECT* flag within an indirect descriptor (ie. only one table per descriptor).

A driver **MUST NOT** create a descriptor chain longer than the Queue Size of the device.

A driver **MUST NOT** set both *VIRTQ\_DESC\_F\_INDIRECT* and *VIRTQ\_DESC\_F\_NEXT* in *flags*.

If *VIRTIO\_F\_IN\_ORDER* has been negotiated, indirect descriptors **MUST** appear sequentially, with *next* taking the value of 1 for the 1st descriptor, 2 for the 2nd one, etc.

#### 2.7.5.3.2 Device Requirements: Indirect Descriptors

The device **MUST** ignore the write-only flag (*flags&VIRTQ\_DESC\_F\_WRITE*) in the descriptor that refers to an indirect table.

The device **MUST** handle the case of zero or more normal chained descriptors followed by a single descriptor with *flags&VIRTQ\_DESC\_F\_INDIRECT*.

**Note:** While unusual (most implementations either create a chain solely using non-indirect descriptors, or use a single indirect element), such a layout is valid.

### 2.7.6 The Virtqueue Available Ring

The available ring has the following layout structure:

```
struct virtq_avail {
#define VIRTQ_AVAIL_F_NO_INTERRUPT    1
    le16 flags;
    le16 idx;
    le16 ring[ /* Queue Size */ ];
    le16 used_event; /* Only if VIRTIO_F_EVENT_IDX */
};
```

The driver uses the available ring to offer buffers to the device: each ring entry refers to the head of a descriptor chain. It is only written by the driver and read by the device.

*idx* field indicates where the driver would put the next descriptor entry in the ring (modulo the queue size). This starts at 0, and increases.

**Note:** The legacy [\[Virtio PCI Draft\]](#) referred to this structure as *vring\_avail*, and the constant as *VRING\_AVAIL\_F\_NO\_INTERRUPT*, but the layout and value were identical.

#### 2.7.6.1 Driver Requirements: The Virtqueue Available Ring

A driver **MUST NOT** decrement the available *idx* on a virtqueue (ie. there is no way to “unexpose” buffers).



## 2.7.7 Used Buffer Notification Suppression

If the `VIRTIO_F_EVENT_IDX` feature bit is not negotiated, the `flags` field in the available ring offers a crude mechanism for the driver to inform the device that it doesn't want notifications when buffers are used. Otherwise `used_event` is a more performant alternative where the driver specifies how far the device can progress before a notification is required.

Neither of these notification suppression methods are reliable, as they are not synchronized with the device, but they serve as useful optimizations.

### 2.7.7.1 Driver Requirements: Used Buffer Notification Suppression

If the `VIRTIO_F_EVENT_IDX` feature bit is not negotiated:

- The driver **MUST** set `flags` to 0 or 1.
- The driver **MAY** set `flags` to 1 to advise the device that notifications are not needed.

Otherwise, if the `VIRTIO_F_EVENT_IDX` feature bit is negotiated:

- The driver **MUST** set `flags` to 0.
- The driver **MAY** use `used_event` to advise the device that notifications are unnecessary until the device writes an entry with an index specified by `used_event` into the used ring (equivalently, until `idx` in the used ring will reach the value `used_event + 1`).

The driver **MUST** handle spurious notifications from the device.

### 2.7.7.2 Device Requirements: Used Buffer Notification Suppression

If the `VIRTIO_F_EVENT_IDX` feature bit is not negotiated:

- The device **MUST** ignore the `used_event` value.
- After the device writes a descriptor index into the used ring:
  - If `flags` is 1, the device **SHOULD NOT** send a notification.
  - If `flags` is 0, the device **MUST** send a notification.

Otherwise, if the `VIRTIO_F_EVENT_IDX` feature bit is negotiated:

- The device **MUST** ignore the lower bit of `flags`.
- After the device writes a descriptor index into the used ring:
  - If the `idx` field in the used ring (which determined where that descriptor index was placed) was equal to `used_event`, the device **MUST** send a notification.
  - Otherwise the device **SHOULD NOT** send a notification.

**Note:** For example, if `used_event` is 0, then a device using

`VIRTIO_F_EVENT_IDX` would send a used buffer notification to the driver after the first buffer is used (and again after the 65536th buffer, etc).

## 2.7.8 The Virtqueue Used Ring

The used ring has the following layout structure:

```
struct virtq_used {
#define VIRTQ_USED_F_NO_NOTIFY 1
    le16 flags;
    le16 idx;
    struct virtq_used_elem ring[ /* Queue Size */];
    le16 avail_event; /* Only if VIRTIO_F_EVENT_IDX */
};

/* le32 is used here for ids for padding reasons. */
struct virtq_used_elem {
```

```

/* Index of start of used descriptor chain. */
le32 id;
/*
 * The number of bytes written into the device writable portion of
 * the buffer described by the descriptor chain.
 */
le32 len;
};

```

The used ring is where the device returns buffers once it is done with them: it is only written to by the device, and read by the driver.

Each entry in the ring is a pair: *id* indicates the head entry of the descriptor chain describing the buffer (this matches an entry placed in the available ring by the guest earlier), and *len* the total of bytes written into the buffer.

**Note:** *len* is particularly useful for drivers using untrusted buffers: if a driver does not know exactly how much has been written by the device, the driver would have to zero the buffer in advance to ensure no data leakage occurs.

For example, a network driver may hand a received buffer directly to an unprivileged userspace application. If the network device has not overwritten the bytes which were in that buffer, this could leak the contents of freed memory from other processes to the application.

*idx* field indicates where the device would put the next descriptor entry in the ring (modulo the queue size). This starts at 0, and increases.

**Note:** The legacy [\[Virtio PCI Draft\]](#) referred to these structures as `vring_used` and `vring_used_elem`, and the constant as `VRING_USED_F_NO_NOTIFY`, but the layout and value were identical.

### 2.7.8.1 Legacy Interface: The Virtqueue Used Ring

Historically, many drivers ignored the *len* value, as a result, many devices set *len* incorrectly. Thus, when using the legacy interface, it is generally a good idea to ignore the *len* value in used ring entries if possible. Specific known issues are listed per device type.

### 2.7.8.2 Device Requirements: The Virtqueue Used Ring

The device **MUST** set *len* prior to updating the used *idx*.

The device **MUST** write at least *len* bytes to descriptor, beginning at the first device-writable buffer, prior to updating the used *idx*.

The device **MAY** write more than *len* bytes to descriptor.

**Note:** There are potential error cases where a device might not know what parts of the buffers have been written. This is why *len* is permitted to be an underestimate: that's preferable to the driver believing that uninitialized memory has been overwritten when it has not.

### 2.7.8.3 Driver Requirements: The Virtqueue Used Ring

The driver **MUST NOT** make assumptions about data in device-writable buffers beyond the first *len* bytes, and **SHOULD** ignore this data.

## 2.7.9 In-order use of descriptors

Some devices always use descriptors in the same order in which they have been made available. These devices can offer the `VIRTIO_F_IN_ORDER` feature. If negotiated, this knowledge allows devices to notify the use of a batch of buffers to the driver by only writing out a single used ring entry with the *id* corresponding to the head entry of the descriptor chain describing the last buffer in the batch.

The device then skips forward in the ring according to the size of the batch. Accordingly, it increments the used *idx* by the size of the batch.

The driver needs to look up the used *id* and calculate the batch size to be able to advance to where the next used ring entry will be written by the device.

This will result in the used ring entry at an offset matching the first available ring entry in the batch, the used ring entry for the next batch at an offset matching the first available ring entry in the next batch, etc.

The skipped buffers (for which no used ring entry was written) are assumed to have been used (read or written) by the device completely.

## 2.7.10 Available Buffer Notification Suppression

The device can suppress available buffer notifications in a manner analogous to the way drivers can suppress used buffer notifications as detailed in section 2.7.7. The device manipulates *flags* or *avail\_event* in the used ring the same way the driver manipulates *flags* or *used\_event* in the available ring.

### 2.7.10.1 Driver Requirements: Available Buffer Notification Suppression

The driver MUST initialize *flags* in the used ring to 0 when allocating the used ring.

If the VIRTIO\_F\_EVENT\_IDX feature bit is not negotiated:

- The driver MUST ignore the *avail\_event* value.
- After the driver writes a descriptor index into the available ring:
  - If *flags* is 1, the driver SHOULD NOT send a notification.
  - If *flags* is 0, the driver MUST send a notification.

Otherwise, if the VIRTIO\_F\_EVENT\_IDX feature bit is negotiated:

- The driver MUST ignore the lower bit of *flags*.
- After the driver writes a descriptor index into the available ring:
  - If the *idx* field in the available ring (which determined where that descriptor index was placed) was equal to *avail\_event*, the driver MUST send a notification.
  - Otherwise the driver SHOULD NOT send a notification.

### 2.7.10.2 Device Requirements: Available Buffer Notification Suppression

If the VIRTIO\_F\_EVENT\_IDX feature bit is not negotiated:

- The device MUST set *flags* to 0 or 1.
- The device MAY set *flags* to 1 to advise the driver that notifications are not needed.

Otherwise, if the VIRTIO\_F\_EVENT\_IDX feature bit is negotiated:

- The device MUST set *flags* to 0.
- The device MAY use *avail\_event* to advise the driver that notifications are unnecessary until the driver writes entry with an index specified by *avail\_event* into the available ring (equivalently, until *idx* in the available ring will reach the value *avail\_event* + 1).

The device MUST handle spurious notifications from the driver.

## 2.7.11 Helpers for Operating Virtqueues

The Linux Kernel Source code contains the definitions above and helper routines in a more usable form, in `include/uapi/linux/virtio_ring.h`. This was explicitly licensed by IBM and Red Hat under the (3-clause) BSD license so that it can be freely used by all other projects, and is reproduced (with slight variation) in [A virtio\\_queue.h](#).

## 2.7.12 Virtqueue Operation

There are two parts to virtqueue operation: supplying new available buffers to the device, and processing used buffers from the device.

**Note:** As an example, the simplest virtio network device has two virtqueues: the transmit virtqueue and the receive virtqueue. The driver adds outgoing (device-readable) packets to the transmit virtqueue, and then frees them after they are used. Similarly, incoming (device-writable) buffers are added to the receive virtqueue, and processed after they are used.

What follows is the requirements of each of these two parts when using the split virtqueue format in more detail.

## 2.7.13 Supplying Buffers to The Device

The driver offers buffers to one of the device's virtqueues as follows:

1. The driver places the buffer into free descriptor(s) in the descriptor table, chaining as necessary (see [2.7.5 The Virtqueue Descriptor Table](#)).
2. The driver places the index of the head of the descriptor chain into the next ring entry of the available ring.
3. Steps 1 and 2 MAY be performed repeatedly if batching is possible.
4. The driver performs a suitable memory barrier to ensure the device sees the updated descriptor table and available ring before the next step.
5. The available *idx* is increased by the number of descriptor chain heads added to the available ring.
6. The driver performs a suitable memory barrier to ensure that it updates the *idx* field before checking for notification suppression.
7. The driver sends an available buffer notification to the device if such notifications are not suppressed.

Note that the above code does not take precautions against the available ring buffer wrapping around: this is not possible since the ring buffer is the same size as the descriptor table, so step (1) will prevent such a condition.

In addition, the maximum queue size is 32768 (the highest power of 2 which fits in 16 bits), so the 16-bit *idx* value can always distinguish between a full and empty buffer.

What follows is the requirements of each stage in more detail.

### 2.7.13.1 Placing Buffers Into The Descriptor Table

A buffer consists of zero or more device-readable physically-contiguous elements followed by zero or more physically-contiguous device-writable elements (each has at least one element). This algorithm maps it into the descriptor table to form a descriptor chain:

for each buffer element, *b*:

1. Get the next free descriptor table entry, *d*
2. Set *d.addr* to the physical address of the start of *b*
3. Set *d.len* to the length of *b*.
4. If *b* is device-writable, set *d.flags* to `VIRTQ_DESC_F_WRITE`, otherwise 0.
5. If there is a buffer element after this:
  - (a) Set *d.next* to the index of the next free descriptor element.
  - (b) Set the `VIRTQ_DESC_F_NEXT` bit in *d.flags*.

In practice, *d.next* is usually used to chain free descriptors, and a separate count kept to check there are enough free descriptors before beginning the mappings.

### 2.7.13.2 Updating The Available Ring

The descriptor chain head is the first *d* in the algorithm above, ie. the index of the descriptor table entry referring to the first part of the buffer. A naive driver implementation MAY do the following (with the appropriate conversion to-and-from little-endian assumed):

```
avail->ring[avail->idx % qsz] = head;
```

However, in general the driver MAY add many descriptor chains before it updates *idx* (at which point they become visible to the device), so it is common to keep a counter of how many the driver has added:

```
avail->ring[(avail->idx + added++) % qsz] = head;
```

### 2.7.13.3 Updating *idx*

*idx* always increments, and wraps naturally at 65536:

```
avail->idx += added;
```

Once available *idx* is updated by the driver, this exposes the descriptor and its contents. The device MAY access the descriptor chains the driver created and the memory they refer to immediately.

#### 2.7.13.3.1 Driver Requirements: Updating *idx*

The driver MUST perform a suitable memory barrier before the *idx* update, to ensure the device sees the most up-to-date copy.

### 2.7.13.4 Notifying The Device

The actual method of device notification is bus-specific, but generally it can be expensive. So the device MAY suppress such notifications if it doesn't need them, as detailed in section 2.7.10.

The driver has to be careful to expose the new *idx* value before checking if notifications are suppressed.

#### 2.7.13.4.1 Driver Requirements: Notifying The Device

The driver MUST perform a suitable memory barrier before reading *flags* or *avail\_event*, to avoid missing a notification.

## 2.7.14 Receiving Used Buffers From The Device

Once the device has used buffers referred to by a descriptor (read from or written to them, or parts of both, depending on the nature of the virtqueue and the device), it sends a used buffer notification to the driver as detailed in section 2.7.7.

**Note:** For optimal performance, a driver MAY disable used buffer notifications while processing the used ring, but beware the problem of missing notifications between emptying the ring and reenabling notifications. This is usually handled by re-checking for more used buffers after notifications are re-enabled:

```
virtq_disable_used_buffer_notifications(vq);

for (;;) {
    if (vq->last_seen_used != le16_to_cpu(virtq->used.idx)) {
        virtq_enable_used_buffer_notifications(vq);
        mb();

        if (vq->last_seen_used != le16_to_cpu(virtq->used.idx))
            break;

        virtq_disable_used_buffer_notifications(vq);
    }
}
```

```

    struct virtq_used_elem *e = virtq.used->ring[vq->last_seen_used%vsz];
    process_buffer(e);
    vq->last_seen_used++;
}

```

## 2.8 Packed Virtqueues

Packed virtqueues is an alternative compact virtqueue layout using read-write memory, that is memory that is both read and written by both host and guest.

Use of packed virtqueues is negotiated by the VIRTIO\_F\_RING\_PACKED feature bit.

Packed virtqueues support up to  $2^{15}$  entries each.

With current transports, virtqueues are located in guest memory allocated by the driver. Each packed virtqueue consists of three parts:

- Descriptor Ring - occupies the Descriptor Area
- Driver Event Suppression - occupies the Driver Area
- Device Event Suppression - occupies the Device Area

Where the Descriptor Ring in turn consists of descriptors, and where each descriptor can contain the following parts:

- Buffer ID
- Element Address
- Element Length
- Flags

A buffer consists of zero or more device-readable physically-contiguous elements followed by zero or more physically-contiguous device-writable elements (each buffer has at least one element).

When the driver wants to send such a buffer to the device, it writes at least one available descriptor describing elements of the buffer into the Descriptor Ring. The descriptor(s) are associated with a buffer by means of a Buffer ID stored within the descriptor.

The driver then notifies the device. When the device has finished processing the buffer, it writes a used device descriptor including the Buffer ID into the Descriptor Ring (overwriting a driver descriptor previously made available), and sends a used event notification.

The Descriptor Ring is used in a circular manner: the driver writes descriptors into the ring in order. After reaching the end of the ring, the next descriptor is placed at the head of the ring. Once the ring is full of driver descriptors, the driver stops sending new requests and waits for the device to start processing descriptors and to write out some used descriptors before making new driver descriptors available.

Similarly, the device reads descriptors from the ring in order and detects that a driver descriptor has been made available. As processing of descriptors is completed, used descriptors are written by the device back into the ring.

Note: after reading driver descriptors and starting their processing in order, the device might complete their processing out of order. Used device descriptors are written in the order in which their processing is complete.

The Device Event Suppression data structure is write-only by the device. It includes information for reducing the number of device events, i.e., sending fewer available buffer notifications to the device.

The Driver Event Suppression data structure is read-only by the device. It includes information for reducing the number of driver events, i.e., sending fewer used buffer notifications to the driver.

## 2.8.1 Driver and Device Ring Wrap Counters

Each of the driver and the device are expected to maintain, internally, a single-bit ring wrap counter initialized to 1.

The counter maintained by the driver is called the Driver Ring Wrap Counter. The driver changes the value of this counter each time it makes available the last descriptor in the ring (after making the last descriptor available).

The counter maintained by the device is called the Device Ring Wrap Counter. The device changes the value of this counter each time it uses the last descriptor in the ring (after marking the last descriptor used).

It is easy to see that the Driver Ring Wrap Counter in the driver matches the Device Ring Wrap Counter in the device when both are processing the same descriptor, or when all available descriptors have been used.

To mark a descriptor as available and used, both the driver and the device use the following two flags:

```
#define VIRTQ_DESC_F_AVAIL    (1 << 7)
#define VIRTQ_DESC_F_USED    (1 << 15)
```

To mark a descriptor as available, the driver sets the VIRTQ\_DESC\_F\_AVAIL bit in Flags to match the internal Driver Ring Wrap Counter. It also sets the VIRTQ\_DESC\_F\_USED bit to match the *inverse* value (i.e. to not match the internal Driver Ring Wrap Counter).

To mark a descriptor as used, the device sets the VIRTQ\_DESC\_F\_USED bit in Flags to match the internal Device Ring Wrap Counter. It also sets the VIRTQ\_DESC\_F\_AVAIL bit to match the *same* value.

Thus VIRTQ\_DESC\_F\_AVAIL and VIRTQ\_DESC\_F\_USED bits are different for an available descriptor and equal for a used descriptor.

Note that this observation is mostly useful for sanity-checking as these are necessary but not sufficient conditions - for example, all descriptors are zero-initialized. To detect used and available descriptors it is possible for drivers and devices to keep track of the last observed value of VIRTQ\_DESC\_F\_USED/VIRTQ\_DESC\_F\_AVAIL. Other techniques to detect VIRTQ\_DESC\_F\_AVAIL/VIRTQ\_DESC\_F\_USED bit changes might also be possible.

## 2.8.2 Polling of available and used descriptors

Writes of device and driver descriptors can generally be reordered, but each side (driver and device) are only required to poll (or test) a single location in memory: the next device descriptor after the one they processed previously, in circular order.

Sometimes the device needs to only write out a single used descriptor after processing a batch of multiple available descriptors. As described in more detail below, this can happen when using descriptor chaining or with in-order use of descriptors. In this case, the device writes out a used descriptor with the buffer id of the last descriptor in the group. After processing the used descriptor, both device and driver then skip forward in the ring the number of the remaining descriptors in the group until processing (reading for the driver and writing for the device) the next used descriptor.

## 2.8.3 Write Flag

In an available descriptor, the VIRTQ\_DESC\_F\_WRITE bit within Flags is used to mark a descriptor as corresponding to a write-only or read-only element of a buffer.

```
/* This marks a descriptor as device write-only (otherwise device read-only). */
#define VIRTQ_DESC_F_WRITE    2
```

In a used descriptor, this bit is used to specify whether any data has been written by the device into any parts of the buffer.



## 2.8.4 Element Address and Length

In an available descriptor, Element Address corresponds to the physical address of the buffer element. The length of the element assumed to be physically contiguous is stored in Element Length.

In a used descriptor, Element Address is unused. Element Length specifies the length of the buffer that has been initialized (written to) by the device.

Element Length is reserved for used descriptors without the `VIRTQ_DESC_F_WRITE` flag, and is ignored by drivers.

## 2.8.5 Scatter-Gather Support

Some drivers need an ability to supply a list of multiple buffer elements (also known as a scatter/gather list) with a request. Two features support this: descriptor chaining and indirect descriptors.

If neither feature is in use by the driver, each buffer is physically-contiguous, either read-only or write-only and is described completely by a single descriptor.

While unusual (most implementations either create all lists solely using non-indirect descriptors, or always use a single indirect element), if both features have been negotiated, mixing indirect and non-indirect descriptors in a ring is valid, as long as each list only contains descriptors of a given type.

Scatter/gather lists only apply to available descriptors. A single used descriptor corresponds to the whole list.

The device limits the number of descriptors in a list through a transport-specific and/or device-specific value. If not limited, the maximum number of descriptors in a list is the virt queue size.

## 2.8.6 Next Flag: Descriptor Chaining

The packed ring format allows the driver to supply a scatter/gather list to the device by using multiple descriptors, and setting the `VIRTQ_DESC_F_NEXT` bit in Flags for all but the last available descriptor.

```
/* This marks a buffer as continuing. */
#define VIRTQ_DESC_F_NEXT 1
```

Buffer ID is included in the last descriptor in the list.

The driver always makes the first descriptor in the list available after the rest of the list has been written out into the ring. This guarantees that the device will never observe a partial scatter/gather list in the ring.

Note: all flags, including `VIRTQ_DESC_F_AVAIL`, `VIRTQ_DESC_F_USED`, `VIRTQ_DESC_F_WRITE` must be set/cleared correctly in all descriptors in the list, not just the first one.

The device only writes out a single used descriptor for the whole list. It then skips forward according to the number of descriptors in the list. The driver needs to keep track of the size of the list corresponding to each buffer ID, to be able to skip to where the next used descriptor is written by the device.

For example, if descriptors are used in the same order in which they are made available, this will result in the used descriptor overwriting the first available descriptor in the list, the used descriptor for the next list overwriting the first available descriptor in the next list, etc.

`VIRTQ_DESC_F_NEXT` is reserved in used descriptors, and should be ignored by drivers.

## 2.8.7 Indirect Flag: Scatter-Gather Support

Some devices benefit by concurrently dispatching a large number of large requests. The `VIRTIO_F_INDIRECT_DESC` feature allows this. To increase ring capacity the driver can store a (read-only by the device) table of indirect descriptors anywhere in memory, and insert a descriptor in the main virtqueue (with *Flags* bit `VIRTQ_DESC_F_INDIRECT` on) that refers to a buffer element containing this indirect descriptor table; *addr* and *len* refer to the indirect table address and length in bytes, respectively.

```
/* This means the element contains a table of descriptors. */
#define VIRTQ_DESC_F_INDIRECT 4
```

The indirect table layout structure looks like this (*len* is the Buffer Length of the descriptor that refers to this table, which is a variable):

```
struct pvirtq_indirect_descriptor_table {
    /* The actual descriptor structures (struct pvirtq_desc each) */
    struct pvirtq_desc desc[len / sizeof(struct pvirtq_desc)];
};
```

The first descriptor is located at the start of the indirect descriptor table, additional indirect descriptors come immediately afterwards. The `VIRTQ_DESC_F_WRITE` *flags* bit is the only valid flag for descriptors in the indirect table. Others are reserved and are ignored by the device. Buffer ID is also reserved and is ignored by the device.

In descriptors with `VIRTQ_DESC_F_INDIRECT` set `VIRTQ_DESC_F_WRITE` is reserved and is ignored by the device.

## 2.8.8 In-order use of descriptors

Some devices always use descriptors in the same order in which they have been made available. These devices can offer the `VIRTIO_F_IN_ORDER` feature. If negotiated, this knowledge allows devices to notify the use of a batch of buffers to the driver by only writing out a single used descriptor with the Buffer ID corresponding to the last descriptor in the batch.

The device then skips forward in the ring according to the size of the batch. The driver needs to look up the used Buffer ID and calculate the batch size to be able to advance to where the next used descriptor will be written by the device.

This will result in the used descriptor overwriting the first available descriptor in the batch, the used descriptor for the next batch overwriting the first available descriptor in the next batch, etc.

The skipped buffers (for which no used descriptor was written) are assumed to have been used (read or written) by the device completely.

## 2.8.9 Multi-buffer requests

Some devices combine multiple buffers as part of processing of a single request. These devices always mark the descriptor corresponding to the first buffer in the request used after the rest of the descriptors (corresponding to rest of the buffers) in the request - which follow the first descriptor in ring order - has been marked used and written out into the ring. This guarantees that the driver will never observe a partial request in the ring.

## 2.8.10 Driver and Device Event Suppression

In many systems used and available buffer notifications involve significant overhead. To mitigate this overhead, each virtqueue includes two identical structures used for controlling notifications between the device and the driver.

The Driver Event Suppression structure is read-only by the device and controls the used buffer notifications sent by the device to the driver.

The Device Event Suppression structure is read-only by the driver and controls the available buffer notifications sent by the driver to the device.

Each of these Event Suppression structures includes the following fields:

**Descriptor Ring Change Event Flags** Takes values:

```
/* Enable events */
#define RING_EVENT_FLAGS_ENABLE 0x0
/* Disable events */
#define RING_EVENT_FLAGS_DISABLE 0x1
/*
 * Enable events for a specific descriptor
 * (as specified by Descriptor Ring Change Event Offset/Wrap Counter).
```

```

* Only valid if VIRTIO_F_EVENT_IDX has been negotiated.
*/
#define RING_EVENT_FLAGS_DESC 0x2
/* The value 0x3 is reserved */

```

**Descriptor Ring Change Event Offset** If Event Flags set to descriptor specific event: offset within the ring (in units of descriptor size). Event will only trigger when this descriptor is made available/used respectively.

**Descriptor Ring Change Event Wrap Counter** If Event Flags set to descriptor specific event: offset within the ring (in units of descriptor size). Event will only trigger when Ring Wrap Counter matches this value and a descriptor is made available/used respectively.

After writing out some descriptors, both the device and the driver are expected to consult the relevant structure to find out whether a used respectively an available buffer notification should be sent.

### 2.8.10.1 Structure Size and Alignment

Each part of the virtqueue is physically-contiguous in guest memory, and has different alignment requirements.

The memory alignment and size requirements, in bytes, of each part of the virtqueue are summarized in the following table:

Virtqueue Part	Alignment	Size
Descriptor Ring	16	16*(Queue Size)
Device Event Suppression	4	4
Driver Event Suppression	4	4

The Alignment column gives the minimum alignment for each part of the virtqueue.

The Size column gives the total number of bytes for each part of the virtqueue.

Queue Size corresponds to the maximum number of descriptors in the virtqueue<sup>3</sup>. The Queue Size value does not have to be a power of 2.

## 2.8.11 Driver Requirements: Virtqueues

The driver MUST ensure that the physical address of the first byte of each virtqueue part is a multiple of the specified alignment value in the above table.

## 2.8.12 Device Requirements: Virtqueues

The device MUST start processing driver descriptors in the order in which they appear in the ring. The device MUST start writing device descriptors into the ring in the order in which they complete. The device MAY reorder descriptor writes once they are started.

## 2.8.13 The Virtqueue Descriptor Format

The available descriptor refers to the buffers the driver is sending to the device. *addr* is a physical address, and the descriptor is identified with a buffer using the *id* field.

```

struct pvirtq_desc {
    /* Buffer Address. */
    le64 addr;
    /* Buffer Length. */
    le32 len;
    /* Buffer ID. */
    le16 id;
    /* The flags depending on descriptor type. */
    le16 flags;
}

```

<sup>3</sup>For example, if Queue Size is 4 then at most 4 buffers can be queued at any given time.

```
};
```

The descriptor ring is zero-initialized.

### 2.8.14 Event Suppression Structure Format

The following structure is used to reduce the number of notifications sent between driver and device.

```
struct pvirtq_event_suppress {
    le16 {
        desc_event_off : 15; /* Descriptor Ring Change Event Offset */
        desc_event_wrap : 1; /* Descriptor Ring Change Event Wrap Counter */
    } desc; /* If desc_event_flags set to RING_EVENT_FLAGS_DESC */
    le16 {
        desc_event_flags : 2, /* Descriptor Ring Change Event Flags */
        reserved : 14; /* Reserved, set to 0 */
    } flags;
};
```

### 2.8.15 Device Requirements: The Virtqueue Descriptor Table

A device MUST NOT write to a device-readable buffer, and a device SHOULD NOT read a device-writable buffer. A device MUST NOT use a descriptor unless it observes the VIRTQ\_DESC\_F\_AVAIL bit in its *flags* being changed (e.g. as compared to the initial zero value). A device MUST NOT change a descriptor after changing it's the VIRTQ\_DESC\_F\_USED bit in its *flags*.

### 2.8.16 Driver Requirements: The Virtqueue Descriptor Table

A driver MUST NOT change a descriptor unless it observes the VIRTQ\_DESC\_F\_USED bit in its *flags* being changed. A driver MUST NOT change a descriptor after changing the VIRTQ\_DESC\_F\_AVAIL bit in its *flags*. When notifying the device, driver MUST set *next\_off* and *next\_wrap* to match the next descriptor not yet made available to the device. A driver MAY send multiple available buffer notifications without making any new descriptors available to the device.

### 2.8.17 Driver Requirements: Scatter-Gather Support

A driver MUST NOT create a descriptor list longer than allowed by the device.

A driver MUST NOT create a descriptor list longer than the Queue Size.

This implies that loops in the descriptor list are forbidden!

The driver MUST place any device-writable descriptor elements after any device-readable descriptor elements.

A driver MUST NOT depend on the device to use more descriptors to be able to write out all descriptors in a list. A driver MUST make sure there's enough space in the ring for the whole list before making the first descriptor in the list available to the device.

A driver MUST NOT make the first descriptor in the list available before all subsequent descriptors comprising the list are made available.

### 2.8.18 Device Requirements: Scatter-Gather Support

The device MUST use descriptors in a list chained by the VIRTQ\_DESC\_F\_NEXT flag in the same order that they were made available by the driver.

The device MAY limit the number of buffers it will allow in a list.

## 2.8.19 Driver Requirements: Indirect Descriptors

The driver MUST NOT set the `DESC_F_INDIRECT` flag unless the `VIRTIO_F_INDIRECT_DESC` feature was negotiated. The driver MUST NOT set any flags except `DESC_F_WRITE` within an indirect descriptor.

A driver MUST NOT create a descriptor chain longer than allowed by the device.

A driver MUST NOT write direct descriptors with `DESC_F_INDIRECT` set in a scatter-gather list linked by `VIRTQ_DESC_F_NEXT`. *flags*.

## 2.8.20 Virtqueue Operation

There are two parts to virtqueue operation: supplying new available buffers to the device, and processing used buffers from the device.

What follows is the requirements of each of these two parts when using the packed virtqueue format in more detail.

## 2.8.21 Supplying Buffers to The Device

The driver offers buffers to one of the device's virtqueues as follows:

1. The driver places the buffer into free descriptor(s) in the Descriptor Ring.
2. The driver performs a suitable memory barrier to ensure that it updates the descriptor(s) before checking for notification suppression.
3. If notifications are not suppressed, the driver notifies the device of the new available buffers.

What follows are the requirements of each stage in more detail.

### 2.8.21.1 Placing Available Buffers Into The Descriptor Ring

For each buffer element, *b*:

1. Get the next descriptor table entry, *d*
2. Get the next free buffer id value
3. Set *d.addr* to the physical address of the start of *b*
4. Set *d.len* to the length of *b*.
5. Set *d.id* to the buffer id
6. Calculate the flags as follows:
  - (a) If *b* is device-writable, set the `VIRTQ_DESC_F_WRITE` bit to 1, otherwise 0
  - (b) Set the `VIRTQ_DESC_F_AVAIL` bit to the current value of the Driver Ring Wrap Counter
  - (c) Set the `VIRTQ_DESC_F_USED` bit to inverse value
7. Perform a memory barrier to ensure that the descriptor has been initialized
8. Set *d.flags* to the calculated flags value
9. If *d* is the last descriptor in the ring, toggle the Driver Ring Wrap Counter
10. Otherwise, increment *d* to point at the next descriptor

This makes a single descriptor buffer available. However, in general the driver MAY make use of a batch of descriptors as part of a single request. In that case, it defers updating the descriptor flags for the first descriptor (and the previous memory barrier) until after the rest of the descriptors have been initialized.

Once the descriptor *flags* field is updated by the driver, this exposes the descriptor and its contents. The device MAY access the descriptor and any following descriptors the driver created and the memory they refer to immediately.

### 2.8.21.1.1 Driver Requirements: Updating flags

The driver MUST perform a suitable memory barrier before the *flags* update, to ensure the device sees the most up-to-date copy.

### 2.8.21.2 Sending Available Buffer Notifications

The actual method of device notification is bus-specific, but generally it can be expensive. So the device MAY suppress such notifications if it doesn't need them, using the Event Suppression structure comprising the Device Area as detailed in section 2.8.14.

The driver has to be careful to expose the new *flags* value before checking if notifications are suppressed.

### 2.8.21.3 Implementation Example

Below is a driver code example. It does not attempt to reduce the number of available buffer notifications, neither does it support the VIRTIO\_F\_EVENT\_IDX feature.

```
/* Note: vq->avail_wrap_count is initialized to 1 */
/* Note: vq->sgs is an array same size as the ring */

id = alloc_id(vq);

first = vq->next_avail;
sgs = 0;
for (each buffer element b) {
    sgs++;

    vq->ids[vq->next_avail] = -1;
    vq->desc[vq->next_avail].address = get_addr(b);
    vq->desc[vq->next_avail].len = get_len(b);

    avail = vq->avail_wrap_count ? VIRTQ_DESC_F_AVAIL : 0;
    used = !vq->avail_wrap_count ? VIRTQ_DESC_F_USED : 0;
    f = get_flags(b) | avail | used;
    if (b is not the last buffer element) {
        f |= VIRTQ_DESC_F_NEXT;
    }

    /* Don't mark the 1st descriptor available until all of them are ready. */
    if (vq->next_avail == first) {
        flags = f;
    } else {
        vq->desc[vq->next_avail].flags = f;
    }

    last = vq->next_avail;

    vq->next_avail++;

    if (vq->next_avail >= vq->size) {
        vq->next_avail = 0;
        vq->avail_wrap_count ^= 1;
    }
}
vq->sgs[id] = sgs;
/* ID included in the last descriptor in the list */
vq->desc[last].id = id;
write_memory_barrier();
vq->desc[first].flags = flags;

memory_barrier();

if (vq->device_event.flags != RING_EVENT_FLAGS_DISABLE) {
    notify_device(vq);
}
```

### 2.8.21.3.1 Driver Requirements: Sending Available Buffer Notifications

The driver MUST perform a suitable memory barrier before reading the Event Suppression structure occupying the Device Area. Failing to do so could result in mandatory available buffer notifications not being sent.

## 2.8.22 Receiving Used Buffers From The Device

Once the device has used buffers referred to by a descriptor (read from or written to them, or parts of both, depending on the nature of the virtqueue and the device), it sends a used buffer notification to the driver as detailed in section 2.8.14.

**Note:** For optimal performance, a driver MAY disable used buffer notifications while processing the used buffers, but beware the problem of missing notifications between emptying the ring and reenabling used buffer notifications. This is usually handled by re-checking for more used buffers after notifications are re-enabled:

```
/* Note: vq->used_wrap_count is initialized to 1 */
vq->driver_event.flags = RING_EVENT_FLAGS_DISABLE;

for (;;) {
    struct pvirtq_desc *d = vq->desc[vq->next_used];

    /*
     * Check that
     * 1. Descriptor has been made available. This check is necessary
     *    if the driver is making new descriptors available in parallel
     *    with this processing of used descriptors (e.g. from another thread).
     *    Note: there are many other ways to check this, e.g.
     *    track the number of outstanding available descriptors or buffers
     *    and check that it's not 0.
     * 2. Descriptor has been used by the device.
     */
    flags = d->flags;
    bool avail = flags & VIRTQ_DESC_F_AVAIL;
    bool used = flags & VIRTQ_DESC_F_USED;
    if (avail != vq->used_wrap_count || used != vq->used_wrap_count) {
        vq->driver_event.flags = RING_EVENT_FLAGS_ENABLE;
        memory_barrier();

        /*
         * Re-test in case the driver made more descriptors available in
         * parallel with the used descriptor processing (e.g. from another
         * thread) and/or the device used more descriptors before the driver
         * enabled events.
         */
        flags = d->flags;
        bool avail = flags & VIRTQ_DESC_F_AVAIL;
        bool used = flags & VIRTQ_DESC_F_USED;
        if (avail != vq->used_wrap_count || used != vq->used_wrap_count) {
            break;
        }

        vq->driver_event.flags = RING_EVENT_FLAGS_DISABLE;
    }

    read_memory_barrier();

    /* skip descriptors until the next buffer */
    id = d->id;
    assert(id < vq->size);
    sgs = vq->sgs[id];
    vq->next_used += sgs;
    if (vq->next_used >= vq->size) {
        vq->next_used -= vq->size;
        vq->used_wrap_count ^= 1;
    }
}
```



```
    free_id(vq, id);  
    process_buffer(d);  
}
```

## 2.9 Driver Notifications

The driver is sometimes required to send an available buffer notification to the device.

When `VIRTIO_F_NOTIFICATION_DATA` has not been negotiated, this notification involves sending the virtqueue number to the device (method depending on the transport).

However, some devices benefit from the ability to find out the amount of available data in the queue without accessing the virtqueue in memory: for efficiency or as a debugging aid.

To help with these optimizations, when `VIRTIO_F_NOTIFICATION_DATA` has been negotiated, driver notifications to the device include the following information:

**vqn** VQ number to be notified.

**next\_off** Offset within the ring where the next available ring entry will be written. When `VIRTIO_F_RING_PACKED` has not been negotiated this refers to the 15 least significant bits of the available index. When `VIRTIO_F_RING_PACKED` has been negotiated this refers to the offset (in units of descriptor entries) within the descriptor ring where the next available descriptor will be written.

**next\_wrap** Wrap Counter. With `VIRTIO_F_RING_PACKED` this is the wrap counter referring to the next available descriptor. Without `VIRTIO_F_RING_PACKED` this is the most significant bit (bit 15) of the available index.

Note that the driver can send multiple notifications even without making any more buffers available. When `VIRTIO_F_NOTIFICATION_DATA` has been negotiated, these notifications would then have identical *next\_off* and *next\_wrap* values.

## 2.10 Shared Memory Regions

Shared memory regions are an additional facility available to devices that need a region of memory that's continuously shared between the device and the driver, rather than passed between them in the way virtqueue elements are.

Example uses include shared caches and version pools for versioned data structures.

The memory region is allocated by the device and presented to the driver. Where the device is implemented in software on a host, this arrangement allows the memory region to be allocated by a library on the host, which the device may not have full control over.

A device may have multiple shared memory regions associated with it. Each region has a *shmid* to identify it, the meaning of which is device-specific.

Enumeration and location of shared memory regions is performed in a transport-specific way.

Memory consistency rules vary depending on the region and the device and they will be specified as required by each device.

### 2.10.1 Addressing within regions

References into shared memory regions are represented as offsets from the beginning of the region instead of absolute memory addresses. Offsets are used both for references between structures stored within shared memory and for requests placed in virtqueues that refer to shared memory. The *shmid* may be explicit or may be inferred from the context of the reference.

### 2.10.2 Device Requirements: Shared Memory Regions

Shared memory regions MUST NOT expose shared memory regions which are used to control the operation of the device, nor to stream data.

## 2.11 Exporting Objects

When an object created by one virtio device needs to be shared with a separate virtio device, the first device can export the object by generating a UUID which can then be passed to the second device to identify the object.

What constitutes an object, how to export objects, and how to import objects are defined by the individual device types. It is RECOMMENDED that devices generate version 4 UUIDs as specified by [\[RFC4122\]](#).

---

## 3 General Initialization And Device Operation

We start with an overview of device initialization, then expand on the details of the device and how each step is preformed. This section is best read along with the bus-specific section which describes how to communicate with the specific device.

### 3.1 Device Initialization

#### 3.1.1 Driver Requirements: Device Initialization

The driver **MUST** follow this sequence to initialize a device:

1. Reset the device.
2. Set the ACKNOWLEDGE status bit: the guest OS has noticed the device.
3. Set the DRIVER status bit: the guest OS knows how to drive the device.
4. Read device feature bits, and write the subset of feature bits understood by the OS and driver to the device. During this step the driver **MAY** read (but **MUST NOT** write) the device-specific configuration fields to check that it can support the device before accepting it.
5. Set the FEATURES\_OK status bit. The driver **MUST NOT** accept new feature bits after this step.
6. Re-read *device status* to ensure the FEATURES\_OK bit is still set: otherwise, the device does not support our subset of features and the device is unusable.
7. Perform device-specific setup, including discovery of virtqueues for the device, optional per-bus setup, reading and possibly writing the device's virtio configuration space, and population of virtqueues.
8. Set the DRIVER\_OK status bit. At this point the device is “live”.

If any of these steps go irrecoverably wrong, the driver **SHOULD** set the FAILED status bit to indicate that it has given up on the device (it can reset the device later to restart if desired). The driver **MUST NOT** continue initialization in that case.

The driver **MUST NOT** send any buffer available notifications to the device before setting DRIVER\_OK.

#### 3.1.2 Legacy Interface: Device Initialization

Legacy devices did not support the FEATURES\_OK status bit, and thus did not have a graceful way for the device to indicate unsupported feature combinations. They also did not provide a clear mechanism to end feature negotiation, which meant that devices finalized features on first-use, and no features could be introduced which radically changed the initial operation of the device.

Legacy driver implementations often used the device before setting the DRIVER\_OK bit, and sometimes even before writing the feature bits to the device.

The result was the steps 5 and 6 were omitted, and steps 4, 7 and 8 were conflated.

Therefore, when using the legacy interface:

- The transitional driver **MUST** execute the initialization sequence as described in 3.1 but omitting the steps 5 and 6.
- The transitional device **MUST** support the driver writing device configuration fields before the step 4.
- The transitional device **MUST** support the driver using the device before the step 8.

## 3.2 Device Operation

When operating the device, each field in the device configuration space can be changed by either the driver or the device.

Whenever such a configuration change is triggered by the device, driver is notified. This makes it possible for drivers to cache device configuration, avoiding expensive configuration reads unless notified.

### 3.2.1 Notification of Device Configuration Changes

For devices where the device-specific configuration information can be changed, a configuration change notification is sent when a device-specific configuration change occurs.

In addition, this notification is triggered by the device setting `DEVICE_NEEDS_RESET` (see [2.1.2](#)).

## 3.3 Device Cleanup

Once the driver has set the `DRIVER_OK` status bit, all the configured virtqueue of the device are considered live. None of the virtqueues of a device are live once the device has been reset.

### 3.3.1 Driver Requirements: Device Cleanup

A driver **MUST NOT** alter virtqueue entries for exposed buffers, i.e., buffers which have been made available to the device (and not been used by the device) of a live virtqueue.

Thus a driver **MUST** ensure a virtqueue isn't live (by device reset) before removing exposed buffers.

---

## 4 Virtio Transport Options

Virtio can use various different buses, thus the standard is split into virtio general and bus-specific sections.

### 4.1 Virtio Over PCI Bus

Virtio devices are commonly implemented as PCI devices.

A Virtio device can be implemented as any kind of PCI device: a Conventional PCI device or a PCI Express device. To assure designs meet the latest level requirements, see the PCI-SIG home page at <http://www.pcisig.com> for any approved changes.

#### 4.1.1 Device Requirements: Virtio Over PCI Bus

A Virtio device using Virtio Over PCI Bus MUST expose to guest an interface that meets the specification requirements of the appropriate PCI specification: [PCI] and [PCIe] respectively.

#### 4.1.2 PCI Device Discovery

Any PCI device with PCI Vendor ID 0x1AF4, and PCI Device ID 0x1000 through 0x107F inclusive is a virtio device. The actual value within this range indicates which virtio device is supported by the device. The PCI Device ID is calculated by adding 0x1040 to the Virtio Device ID, as indicated in section 5. Additionally, devices MAY utilize a Transitional PCI Device ID range, 0x1000 to 0x103F depending on the device type.

##### 4.1.2.1 Device Requirements: PCI Device Discovery

Devices MUST have the PCI Vendor ID 0x1AF4. Devices MUST either have the PCI Device ID calculated by adding 0x1040 to the Virtio Device ID, as indicated in section 5 or have the Transitional PCI Device ID depending on the device type, as follows:

Transitional PCI Device ID	Virtio Device
0x1000	network card
0x1001	block device
0x1002	memory ballooning (traditional)
0x1003	console
0x1004	SCSI host
0x1005	entropy source
0x1009	9P transport

For example, the network card device with the Virtio Device ID 1 has the PCI Device ID 0x1041 or the Transitional PCI Device ID 0x1000.

The PCI Subsystem Vendor ID and the PCI Subsystem Device ID MAY reflect the PCI Vendor and Device ID of the environment (for informational purposes by the driver).

Non-transitional devices SHOULD have a PCI Device ID in the range 0x1040 to 0x107f. Non-transitional devices SHOULD have a PCI Revision ID of 1 or higher. Non-transitional devices SHOULD have a PCI Subsystem Device ID of 0x40 or higher.

This is to reduce the chance of a legacy driver attempting to drive the device.

#### 4.1.2.2 Driver Requirements: PCI Device Discovery

Drivers MUST match devices with the PCI Vendor ID 0x1AF4 and the PCI Device ID in the range 0x1040 to 0x107f, calculated by adding 0x1040 to the Virtio Device ID, as indicated in section 5. Drivers for device types listed in section 4.1.2 MUST match devices with the PCI Vendor ID 0x1AF4 and the Transitional PCI Device ID indicated in section 4.1.2.

Drivers MUST match any PCI Revision ID value. Drivers MAY match any PCI Subsystem Vendor ID and any PCI Subsystem Device ID value.

#### 4.1.2.3 Legacy Interfaces: A Note on PCI Device Discovery

Transitional devices MUST have a PCI Revision ID of 0. Transitional devices MUST have the PCI Subsystem Device ID matching the Virtio Device ID, as indicated in section 5. Transitional devices MUST have the Transitional PCI Device ID in the range 0x1000 to 0x103f.

This is to match legacy drivers.

### 4.1.3 PCI Device Layout

The device is configured via I/O and/or memory regions (though see 4.1.4.9 for access via the PCI configuration space), as specified by Virtio Structure PCI Capabilities.

Fields of different sizes are present in the device configuration regions. All 64-bit, 32-bit and 16-bit fields are little-endian. 64-bit fields are to be treated as two 32-bit fields, with low 32 bit part followed by the high 32 bit part.

#### 4.1.3.1 Driver Requirements: PCI Device Layout

For device configuration access, the driver MUST use 8-bit wide accesses for 8-bit wide fields, 16-bit wide and aligned accesses for 16-bit wide fields and 32-bit wide and aligned accesses for 32-bit and 64-bit wide fields. For 64-bit fields, the driver MAY access each of the high and low 32-bit parts of the field independently.

#### 4.1.3.2 Device Requirements: PCI Device Layout

For 64-bit device configuration fields, the device MUST allow driver independent access to high and low 32-bit parts of the field.

### 4.1.4 Virtio Structure PCI Capabilities

The virtio device configuration layout includes several structures:

- Common configuration
- Notifications
- ISR Status
- Device-specific configuration (optional)
- PCI configuration access

Each structure can be mapped by a Base Address register (BAR) belonging to the function, or accessed via the special VIRTIO\_PCI\_CAP\_PCI\_CFG field in the PCI configuration space.

The location of each structure is specified using a vendor-specific PCI capability located on the capability list in PCI configuration space of the device. This virtio structure capability uses little-endian format; all fields are read-only for the driver unless stated otherwise:

```
struct virtio_pci_cap {
    u8 cap_vndr; /* Generic PCI field: PCI_CAP_ID_VNDR */
    u8 cap_next; /* Generic PCI field: next ptr. */
    u8 cap_len; /* Generic PCI field: capability length */
    u8 cfg_type; /* Identifies the structure. */
    u8 bar; /* Where to find it. */
}
```

```

    u8 id;          /* Multiple capabilities of the same type */
    u8 padding[2];  /* Pad to full dword. */
    le32 offset;    /* Offset within bar. */
    le32 length;    /* Length of the structure, in bytes. */
};

```

This structure can be followed by extra data, depending on *cfg\_type*, as documented below.

The fields are interpreted as follows:

**cap\_vndr** 0x09; Identifies a vendor-specific capability.

**cap\_next** Link to next capability in the capability list in the PCI configuration space.

**cap\_len** Length of this capability structure, including the whole of struct virtio\_pci\_cap, and extra data if any. This length MAY include padding, or fields unused by the driver.

**cfg\_type** identifies the structure, according to the following table:

```

/* Common configuration */
#define VIRTIO_PCI_CAP_COMMON_CFG      1
/* Notifications */
#define VIRTIO_PCI_CAP_NOTIFY_CFG      2
/* ISR Status */
#define VIRTIO_PCI_CAP_ISR_CFG         3
/* Device specific configuration */
#define VIRTIO_PCI_CAP_DEVICE_CFG      4
/* PCI configuration access */
#define VIRTIO_PCI_CAP_PCI_CFG         5
/* Shared memory region */
#define VIRTIO_PCI_CAP_SHARED_MEMORY_CFG 8
/* Vendor-specific data */
#define VIRTIO_PCI_CAP_VENDOR_CFG      9

```

Any other value is reserved for future use.

Each structure is detailed individually below.

The device MAY offer more than one structure of any type - this makes it possible for the device to expose multiple interfaces to drivers. The order of the capabilities in the capability list specifies the order of preference suggested by the device. A device may specify that this ordering mechanism be overridden by the use of the *id* field.

**Note:** For example, on some hypervisors, notifications using IO accesses are faster than memory accesses. In this case, the device would expose two capabilities with *cfg\_type* set to VIRTIO\_PCI\_CAP\_NOTIFY\_CFG: the first one addressing an I/O BAR, the second one addressing a memory BAR. In this example, the driver would use the I/O BAR if I/O resources are available, and fall back on memory BAR when I/O resources are unavailable.

**bar** values 0x0 to 0x5 specify a Base Address register (BAR) belonging to the function located beginning at 10h in PCI Configuration Space and used to map the structure into Memory or I/O Space. The BAR is permitted to be either 32-bit or 64-bit, it can map Memory Space or I/O Space.

Any other value is reserved for future use.

**id** Used by some device types to uniquely identify multiple capabilities of a certain type. If the device type does not specify the meaning of this field, its contents are undefined.

**offset** indicates where the structure begins relative to the base address associated with the BAR. The alignment requirements of *offset* are indicated in each structure-specific section below.

**length** indicates the length of the structure.

*length* MAY include padding, or fields unused by the driver, or future extensions.

**Note:** For example, a future device might present a large structure size of several MBytes. As current devices never utilize structures larger than 4KBytes in size, driver MAY limit the mapped structure size to e.g. 4KBytes (thus ignoring parts of structure after the first 4KBytes) to allow forward compatibility with such devices without loss of functionality and without wasting resources.



A variant of this type, struct `virtio_pci_cap64`, is defined for those capabilities that require offsets or lengths larger than 4GiB:

```
struct virtio_pci_cap64 {
    struct virtio_pci_cap cap;
    u32 offset_hi;
    u32 length_hi;
};
```

Given that the `cap.length` and `cap.offset` fields are only 32 bit, the additional `offset_hi` and `length_hi` fields provide the most significant 32 bits of a total 64 bit offset and length within the BAR specified by `cap.bar`.

#### 4.1.4.1 Driver Requirements: Virtio Structure PCI Capabilities

The driver MUST ignore any vendor-specific capability structure which has a reserved `cfg_type` value.

The driver SHOULD use the first instance of each virtio structure type they can support.

The driver MUST accept a `cap_len` value which is larger than specified here.

The driver MUST ignore any vendor-specific capability structure which has a reserved `bar` value.

The drivers SHOULD only map part of configuration structure large enough for device operation. The drivers MUST handle an unexpectedly large `length`, but MAY check that `length` is large enough for device operation.

The driver MUST NOT write into any field of the capability structure, with the exception of those with `cap_type` `VIRTIO_PCI_CAP_PCI_CFG` as detailed in 4.1.4.9.2.

#### 4.1.4.2 Device Requirements: Virtio Structure PCI Capabilities

The device MUST include any extra data (from the beginning of the `cap_vndr` field through end of the extra data fields if any) in `cap_len`. The device MAY append extra data or padding to any structure beyond that.

If the device presents multiple structures of the same type, it SHOULD order them from optimal (first) to least-optimal (last).

#### 4.1.4.3 Common configuration structure layout

The common configuration structure is found at the `bar` and `offset` within the `VIRTIO_PCI_CAP_COMMON_CFG` capability; its layout is below.

```
struct virtio_pci_common_cfg {
    /* About the whole device. */
    le32 device_feature_select; /* read-write */
    le32 device_feature; /* read-only for driver */
    le32 driver_feature_select; /* read-write */
    le32 driver_feature; /* read-write */
    le16 config_msix_vector; /* read-write */
    le16 num_queues; /* read-only for driver */
    u8 device_status; /* read-write */
    u8 config_generation; /* read-only for driver */

    /* About a specific virtqueue. */
    le16 queue_select; /* read-write */
    le16 queue_size; /* read-write */
    le16 queue_msix_vector; /* read-write */
    le16 queue_enable; /* read-write */
    le16 queue_notify_off; /* read-only for driver */
    le64 queue_desc; /* read-write */
    le64 queue_driver; /* read-write */
    le64 queue_device; /* read-write */
    le16 queue_notify_data; /* read-only for driver */
    le16 queue_reset; /* read-write */
};
```

**device\_feature\_select** The driver uses this to select which feature bits `device_feature` shows. Value 0x0 selects Feature Bits 0 to 31, 0x1 selects Feature Bits 32 to 63, etc.

**device\_feature** The device uses this to report which feature bits it is offering to the driver: the driver writes to *device\_feature\_select* to select which feature bits are presented.

**driver\_feature\_select** The driver uses this to select which feature bits *driver\_feature* shows. Value 0x0 selects Feature Bits 0 to 31, 0x1 selects Feature Bits 32 to 63, etc.

**driver\_feature** The driver writes this to accept feature bits offered by the device. Driver Feature Bits selected by *driver\_feature\_select*.

**config\_msix\_vector** The driver sets the Configuration Vector for MSI-X.

**num\_queues** The device specifies the maximum number of virtqueues supported here.

**device\_status** The driver writes the device status here (see 2.1). Writing 0 into this field resets the device.

**config\_generation** Configuration atomicity value. The device changes this every time the configuration noticeably changes.

**queue\_select** Queue Select. The driver selects which virtqueue the following fields refer to.

**queue\_size** Queue Size. On reset, specifies the maximum queue size supported by the device. This can be modified by the driver to reduce memory requirements. A 0 means the queue is unavailable.

**queue\_msix\_vector** The driver uses this to specify the queue vector for MSI-X.

**queue\_enable** The driver uses this to selectively prevent the device from executing requests from this virtqueue. 1 - enabled; 0 - disabled.

**queue\_notify\_off** The driver reads this to calculate the offset from start of Notification structure at which this virtqueue is located.

**Note:** this is *not an offset in bytes*. See 4.1.4.4 below.

**queue\_desc** The driver writes the physical address of Descriptor Area here. See section 2.6.

**queue\_driver** The driver writes the physical address of Driver Area here. See section 2.6.

**queue\_device** The driver writes the physical address of Device Area here. See section 2.6.

**queue\_notify\_data** This field exists only if VIRTIO\_F\_NOTIF\_CONFIG\_DATA has been negotiated. The driver will use this value to put it in the 'virtqueue number' field in the available buffer notification structure. See section 4.1.5.2.

**Note:** This field provides the device with flexibility to determine how virtqueues will be referred to in available buffer notifications. In a trivial case the device can set *queue\_notify\_data*=vqn. Some devices may benefit from providing another value, for example an internal virtqueue identifier, or an internal offset related to the virtqueue number.

**queue\_reset** The driver uses this to selectively reset the queue. This field exists only if VIRTIO\_F\_RING\_RESET has been negotiated. (see 2.6.1).

#### 4.1.4.3.1 Device Requirements: Common configuration structure layout

*offset* MUST be 4-byte aligned.

The device MUST present at least one common configuration capability.

The device MUST present the feature bits it is offering in *device\_feature*, starting at bit *device\_feature\_select* \* 32 for any *device\_feature\_select* written by the driver.

**Note:** This means that it will present 0 for any *device\_feature\_select* other than 0 or 1, since no feature defined here exceeds 63.

The device MUST present any valid feature bits the driver has written in *driver\_feature*, starting at bit *driver\_feature\_select* \* 32 for any *driver\_feature\_select* written by the driver. Valid feature bits are those which are subset of the corresponding *device\_feature* bits. The device MAY present invalid bits written by the driver.

**Note:** This means that a device can ignore writes for feature bits it never offers, and simply present 0 on reads. Or it can just mirror what the driver wrote (but it will still have to check them when the driver sets `FEATURES_OK`).

**Note:** A driver shouldn't write invalid bits anyway, as per [3.1.1](#), but this attempts to handle it.

The device **MUST** present a changed `config_generation` after the driver has read a device-specific configuration value which has changed since any part of the device-specific configuration was last read.

**Note:** As `config_generation` is an 8-bit value, simply incrementing it on every configuration change could violate this requirement due to wrap. Better would be to set an internal flag when it has changed, and if that flag is set when the driver reads from the device-specific configuration, increment `config_generation` and clear the flag.

The device **MUST** reset when 0 is written to `device_status`, and present a 0 in `device_status` once that is done.

The device **MUST** present a 0 in `queue_enable` on reset.

If `VIRTIO_F_RING_RESET` has been negotiated, the device **MUST** present a 0 in `queue_enable` after the driver has reset the virtqueue via `queue_reset`.

If `VIRTIO_F_RING_RESET` has been negotiated, the device **MUST** present a 0 in `queue_reset` on reset.

If `VIRTIO_F_RING_RESET` has been negotiated, the device **MUST** present a 0 in `queue_reset` after the virtqueue is enabled with `queue_enable`.

The device **MUST** reset the queue when 1 is written to `queue_reset`, and present a 1 in `queue_reset` after the queue has been reset, until the driver re-enables the queue via `queue_enable` or the device is reset. The device **MUST** present consistent default values after queue reset. (see [2.6.1](#)).

The device **MUST** present a 0 in `queue_size` if the virtqueue corresponding to the current `queue_select` is unavailable.

If `VIRTIO_F_RING_PACKED` has not been negotiated, the device **MUST** present either a value of 0 or a power of 2 in `queue_size`.

#### 4.1.4.3.2 Driver Requirements: Common configuration structure layout

The driver **MUST NOT** write to `device_feature`, `num_queues`, `config_generation`, `queue_notify_off` or `queue_notify_data`.

If `VIRTIO_F_RING_PACKED` has been negotiated, the driver **MUST NOT** write the value 0 to `queue_size`. If `VIRTIO_F_RING_PACKED` has not been negotiated, the driver **MUST NOT** write a value which is not a power of 2 to `queue_size`.

The driver **MUST** configure the other virtqueue fields before enabling the virtqueue with `queue_enable`.

After writing 0 to `device_status`, the driver **MUST** wait for a read of `device_status` to return 0 before reinitializing the device.

The driver **MUST NOT** write a 0 to `queue_enable`.

If `VIRTIO_F_RING_RESET` has been negotiated, after the driver writes 1 to `queue_reset` to reset the queue, it **MUST** verify that the queue has been reset by reading back `queue_reset` and ensuring that it is 1. The driver **MAY** re-enable the queue by writing a 1 to `queue_enable` after ensuring that the other virtqueue fields have been set up correctly. The driver **MAY** set driver-writeable queue configuration values to different values than those that were used before the queue reset. (see [2.6.1](#)).

#### 4.1.4.4 Notification structure layout

The notification location is found using the `VIRTIO_PCI_CAP_NOTIFY_CFG` capability. This capability is immediately followed by an additional field, like so:

```
struct virtio_pci_notify_cap {
    struct virtio_pci_cap cap;
    le32 notify_off_multiplier; /* Multiplier for queue_notify_off. */
};
```

*notify\_off\_multiplier* is combined with the *queue\_notify\_off* to derive the Queue Notify address within a BAR for a virtqueue:

$$\text{cap.offset} + \text{queue\_notify\_off} * \text{notify\_off\_multiplier}$$

The *cap.offset* and *notify\_off\_multiplier* are taken from the notification capability structure above, and the *queue\_notify\_off* is taken from the common configuration structure.

**Note:** For example, if *notify\_off\_multiplier* is 0, the device uses the same Queue Notify address for all queues.

#### 4.1.4.4.1 Device Requirements: Notification capability

The device **MUST** present at least one notification capability.

For devices not offering VIRTIO\_F\_NOTIFICATION\_DATA:

The *cap.offset* **MUST** be 2-byte aligned.

The device **MUST** either present *notify\_off\_multiplier* as an even power of 2, or present *notify\_off\_multiplier* as 0.

The value *cap.length* presented by the device **MUST** be at least 2 and **MUST** be large enough to support queue notification offsets for all supported queues in all possible configurations.

For all queues, the value *cap.length* presented by the device **MUST** satisfy:

$$\text{cap.length} \geq \text{queue\_notify\_off} * \text{notify\_off\_multiplier} + 2$$

For devices offering VIRTIO\_F\_NOTIFICATION\_DATA:

The device **MUST** either present *notify\_off\_multiplier* as a number that is a power of 2 that is also a multiple 4, or present *notify\_off\_multiplier* as 0.

The *cap.offset* **MUST** be 4-byte aligned.

The value *cap.length* presented by the device **MUST** be at least 4 and **MUST** be large enough to support queue notification offsets for all supported queues in all possible configurations.

For all queues, the value *cap.length* presented by the device **MUST** satisfy:

$$\text{cap.length} \geq \text{queue\_notify\_off} * \text{notify\_off\_multiplier} + 4$$

#### 4.1.4.5 ISR status capability

The VIRTIO\_PCI\_CAP\_ISR\_CFG capability refers to at least a single byte, which contains the 8-bit ISR status field to be used for INT#x interrupt handling.

The *offset* for the *ISR status* has no alignment requirements.

The ISR bits allow the driver to distinguish between device-specific configuration change interrupts and normal virtqueue interrupts:

Bits	0	1	2 to 31
Purpose	Queue Interrupt	Device Configuration Interrupt	Reserved

To avoid an extra access, simply reading this register resets it to 0 and causes the device to de-assert the interrupt.

In this way, driver read of ISR status causes the device to de-assert an interrupt.

See sections 4.1.5.3 and 4.1.5.4 for how this is used.

#### 4.1.4.5.1 Device Requirements: ISR status capability

The device MUST present at least one VIRTIO\_PCI\_CAP\_ISR\_CFG capability.

The device MUST set the Device Configuration Interrupt bit in *ISR status* before sending a device configuration change notification to the driver.

If MSI-X capability is disabled, the device MUST set the Queue Interrupt bit in *ISR status* before sending a virtqueue notification to the driver.

If MSI-X capability is disabled, the device MUST set the Interrupt Status bit in the PCI Status register in the PCI Configuration Header of the device to the logical OR of all bits in *ISR status* of the device. The device then asserts/deasserts INT#x interrupts unless masked according to standard PCI rules [PCI].

The device MUST reset *ISR status* to 0 on driver read.

#### 4.1.4.5.2 Driver Requirements: ISR status capability

If MSI-X capability is enabled, the driver SHOULD NOT access *ISR status* upon detecting a Queue Interrupt.

#### 4.1.4.6 Device-specific configuration

The device MUST present at least one VIRTIO\_PCI\_CAP\_DEVICE\_CFG capability for any device type which has a device-specific configuration.

##### 4.1.4.6.1 Device Requirements: Device-specific configuration

The *offset* for the device-specific configuration MUST be 4-byte aligned.

#### 4.1.4.7 Shared memory capability

Shared memory regions 2.10 are enumerated on the PCI transport as a sequence of VIRTIO\_PCI\_CAP\_SHARED\_MEMORY\_CFG capabilities, one per region.

The capability is defined by a struct `virtio_pci_cap64` and utilises the *cap.id* to allow multiple shared memory regions per device. The identifier in *cap.id* does not denote a certain order of preference; it is only used to uniquely identify a region.

##### 4.1.4.7.1 Device Requirements: Shared memory capability

The region defined by the combination of the *cap.offset*, *offset\_hi*, and *cap.length*, *length\_hi* fields MUST be contained within the BAR specified by *cap.bar*.

The *cap.id* MUST be unique for any one device instance.

#### 4.1.4.8 Vendor data capability

The optional Vendor data capability allows the device to present vendor-specific data to the driver, without conflicts, for debugging and/or reporting purposes, and without conflicting with standard functionality.

This capability augments but does not replace the standard subsystem ID and subsystem vendor ID fields (offsets 0x2C and 0x2E in the PCI configuration space header) as specified by [PCI].

Vendor data capability is enumerated on the PCI transport as a VIRTIO\_PCI\_CAP\_VENDOR\_CFG capability.

The capability has the following structure:

```
struct virtio_pci_vndr_data {
    u8 cap_vndr;    /* Generic PCI field: PCI_CAP_ID_VNDR */
    u8 cap_next;    /* Generic PCI field: next ptr. */
    u8 cap_len;     /* Generic PCI field: capability length */
};
```

```

    u8 cfg_type; /* Identifies the structure. */
    u16 vendor_id; /* Identifies the vendor-specific format. */
    /* For Vendor Definition */
    /* Pads structure to a multiple of 4 bytes */
    /* Reads must not have side effects */
};

```

Where *vendor\_id* identifies the PCI-SIG assigned Vendor ID as specified by [PCI].

Note that the capability size is required to be a multiple of 4.

To make it safe for a generic driver to access the capability, reads from this capability MUST NOT have any side effects.

#### 4.1.4.8.1 Device Requirements: Vendor data capability

Devices CAN present *vendor\_id* that does not match either the PCI Vendor ID or the PCI Subsystem Vendor ID.

Devices CAN present multiple Vendor data capabilities with either different or identical *vendor\_id* values.

The value *vendor\_id* MUST NOT equal 0x1AF4.

The size of the Vendor data capability MUST be a multiple of 4 bytes.

Reads of the Vendor data capability by the driver MUST NOT have any side effects.

#### 4.1.4.8.2 Driver Requirements: Vendor data capability

The driver SHOULD NOT use the Vendor data capability except for debugging and reporting purposes.

The driver MUST qualify the *vendor\_id* before interpreting or writing into the Vendor data capability.

#### 4.1.4.9 PCI configuration access capability

The VIRTIO\_PCI\_CAP\_PCI\_CFG capability creates an alternative (and likely suboptimal) access method to the common configuration, notification, ISR and device-specific configuration regions.

The capability is immediately followed by an additional field like so:

```

struct virtio_pci_cfg_cap {
    struct virtio_pci_cap cap;
    u8 pci_cfg_data[4]; /* Data for BAR access. */
};

```

The fields *cap.bar*, *cap.length*, *cap.offset* and *pci\_cfg\_data* are read-write (RW) for the driver.

To access a device region, the driver writes into the capability structure (ie. within the PCI configuration space) as follows:

- The driver sets the BAR to access by writing to *cap.bar*.
- The driver sets the size of the access by writing 1, 2 or 4 to *cap.length*.
- The driver sets the offset within the BAR by writing to *cap.offset*.

At that point, *pci\_cfg\_data* will provide a window of size *cap.length* into the given *cap.bar* at offset *cap.offset*.

#### 4.1.4.9.1 Device Requirements: PCI configuration access capability

The device MUST present at least one VIRTIO\_PCI\_CAP\_PCI\_CFG capability.

Upon detecting driver write access to *pci\_cfg\_data*, the device MUST execute a write access at offset *cap.offset* at BAR selected by *cap.bar* using the first *cap.length* bytes from *pci\_cfg\_data*.



Upon detecting driver read access to *pci\_cfg\_data*, the device MUST execute a read access of length *cap.length* at offset *cap.offset* at BAR selected by *cap.bar* and store the first *cap.length* bytes in *pci\_cfg\_data*.

#### 4.1.4.9.2 Driver Requirements: PCI configuration access capability

The driver MUST NOT write a *cap.offset* which is not a multiple of *cap.length* (ie. all accesses MUST be aligned).

The driver MUST NOT read or write *pci\_cfg\_data* unless *cap.bar*, *cap.length* and *cap.offset* address *cap.length* bytes within a BAR range specified by some other Virtio Structure PCI Capability of type other than *VIRTIO\_PCI\_CAP\_PCI\_CFG*.

#### 4.1.4.10 Legacy Interfaces: A Note on PCI Device Layout

Transitional devices MUST present part of configuration registers in a legacy configuration structure in BAR0 in the first I/O region of the PCI device, as documented below. When using the legacy interface, transitional drivers MUST use the legacy configuration structure in BAR0 in the first I/O region of the PCI device, as documented below.

When using the legacy interface the driver MAY access the device-specific configuration region using any width accesses, and a transitional device MUST present driver with the same results as when accessed using the “natural” access method (i.e. 32-bit accesses for 32-bit fields, etc).

Note that this is possible because while the virtio common configuration structure is PCI (i.e. little) endian, when using the legacy interface the device-specific configuration region is encoded in the native endian of the guest (where such distinction is applicable).

When used through the legacy interface, the virtio common configuration structure looks as follows:

Bits	32	32	32	16	16	16	8	8
Read / Write	R	R+W	R+W	R	R+W	R+W	R+W	R
Purpose	Device Features bits 0:31	Driver Features bits 0:31	Queue Address	<i>queue_size</i>	<i>queue_select</i>	Queue Notify	Device Status	ISR Status

If MSI-X is enabled for the device, two additional fields immediately follow this header:

Bits	16	16
Read/Write	R+W	R+W
Purpose (MSI-X)	<i>config_msix_vector</i>	<i>queue_msix_vector</i>

Note: When MSI-X capability is enabled, device-specific configuration starts at byte offset 24 in virtio common configuration structure. When MSI-X capability is not enabled, device-specific configuration starts at byte offset 20 in virtio header. ie. once you enable MSI-X on the device, the other fields move. If you turn it off again, they move back!

Any device-specific configuration space immediately follows these general headers:

Bits	Device Specific	...
Read / Write	Device Specific	
Purpose	Device Specific	

When accessing the device-specific configuration space using the legacy interface, transitional drivers MUST access the device-specific configuration space at an offset immediately following the general headers.

When using the legacy interface, transitional devices MUST present the device-specific configuration space if any at an offset immediately following the general headers.



Note that only Feature Bits 0 to 31 are accessible through the Legacy Interface. When used through the Legacy Interface, Transitional Devices MUST assume that Feature Bits 32 to 63 are not acknowledged by Driver.

As legacy devices had no *config\_generation* field, see [2.5.4 Legacy Interface: Device Configuration Space](#) for workarounds.

#### 4.1.4.11 Non-transitional Device With Legacy Driver: A Note on PCI Device Layout

All known legacy drivers check either the PCI Revision or the Device and Vendor IDs, and thus won't attempt to drive a non-transitional device.

A buggy legacy driver might mistakenly attempt to drive a non-transitional device. If support for such drivers is required (as opposed to fixing the bug), the following would be the recommended way to detect and handle them.

**Note:** Such buggy drivers are not currently known to be used in production.

##### 4.1.4.11.0.1 Device Requirements: Non-transitional Device With Legacy Driver

Non-transitional devices, on a platform where a legacy driver for a legacy device with the same ID (including PCI Revision, Device and Vendor IDs) is known to have previously existed, SHOULD take the following steps to cause the legacy driver to fail gracefully when it attempts to drive them:

1. Present an I/O BAR in BAR0, and
2. Respond to a single-byte zero write to offset 18 (corresponding to Device Status register in the legacy layout) of BAR0 by presenting zeroes on every BAR and ignoring writes.

### 4.1.5 PCI-specific Initialization And Device Operation

#### 4.1.5.1 Device Initialization

This documents PCI-specific steps executed during Device Initialization.

##### 4.1.5.1.1 Virtio Device Configuration Layout Detection

As a prerequisite to device initialization, the driver scans the PCI capability list, detecting virtio configuration layout using Virtio Structure PCI capabilities as detailed in [4.1.4](#)

##### 4.1.5.1.1.1 Legacy Interface: A Note on Device Layout Detection

Legacy drivers skipped the Device Layout Detection step, assuming legacy device configuration space in BAR0 in I/O space unconditionally.

Legacy devices did not have the Virtio PCI Capability in their capability list.

Therefore:

Transitional devices MUST expose the Legacy Interface in I/O space in BAR0.

Transitional drivers MUST look for the Virtio PCI Capabilities on the capability list. If these are not present, driver MUST assume a legacy device, and use it through the legacy interface.

Non-transitional drivers MUST look for the Virtio PCI Capabilities on the capability list. If these are not present, driver MUST assume a legacy device, and fail gracefully.

##### 4.1.5.1.2 MSI-X Vector Configuration

When MSI-X capability is present and enabled in the device (through standard PCI configuration space) *config\_msix\_vector* and *queue\_msix\_vector* are used to map configuration change and queue interrupts to MSI-X vectors. In this case, the ISR Status is unused.

Writing a valid MSI-X Table entry number, 0 to 0x7FF, to *config\_msix\_vector/queue\_msix\_vector* maps interrupts triggered by the configuration change/selected queue events respectively to the corresponding MSI-X vector. To disable interrupts for an event type, the driver unmaps this event by writing a special NO\_VECTOR value:

```
/* Vector value used to disable MSI for queue */  
#define VIRTIO_MSI_NO_VECTOR          0xffff
```

Note that mapping an event to vector might require device to allocate internal device resources, and thus could fail.

#### 4.1.5.1.2.1 Device Requirements: MSI-X Vector Configuration

A device that has an MSI-X capability SHOULD support at least 2 and at most 0x800 MSI-X vectors. Device MUST report the number of vectors supported in *Table Size* in the MSI-X Capability as specified in [PCI]. The device SHOULD restrict the reported MSI-X Table Size field to a value that might benefit system performance.

**Note:** For example, a device which does not expect to send interrupts at a high rate might only specify 2 MSI-X vectors.

Device MUST support mapping any event type to any valid vector 0 to MSI-X *Table Size*. Device MUST support unmapping any event type.

The device MUST return vector mapped to a given event, (NO\_VECTOR if unmapped) on read of *config\_msix\_vector/queue\_msix\_vector*. The device MUST have all queue and configuration change events are unmapped upon reset.

Devices SHOULD NOT cause mapping an event to vector to fail unless it is impossible for the device to satisfy the mapping request. Devices MUST report mapping failures by returning the NO\_VECTOR value when the relevant *config\_msix\_vector/queue\_msix\_vector* field is read.

#### 4.1.5.1.2.2 Driver Requirements: MSI-X Vector Configuration

Driver MUST support device with any MSI-X Table Size 0 to 0x7FF. Driver MAY fall back on using INT#x interrupts for a device which only supports one MSI-X vector (MSI-X Table Size = 0).

Driver MAY interpret the Table Size as a hint from the device for the suggested number of MSI-X vectors to use.

Driver MUST NOT attempt to map an event to a vector outside the MSI-X Table supported by the device, as reported by *Table Size* in the MSI-X Capability.

After mapping an event to vector, the driver MUST verify success by reading the Vector field value: on success, the previously written value is returned, and on failure, NO\_VECTOR is returned. If a mapping failure is detected, the driver MAY retry mapping with fewer vectors, disable MSI-X or report device failure.

#### 4.1.5.1.3 Virtqueue Configuration

As a device can have zero or more virtqueues for bulk data transport<sup>1</sup>, the driver needs to configure them as part of the device-specific configuration.

The driver typically does this as follows, for each virtqueue a device has:

1. Write the virtqueue index (first queue is 0) to *queue\_select*.
2. Read the virtqueue size from *queue\_size*. This controls how big the virtqueue is (see 2.6 Virtqueues). If this field is 0, the virtqueue does not exist.
3. Optionally, select a smaller virtqueue size and write it to *queue\_size*.

<sup>1</sup>For example, the simplest network device has two virtqueues.

4. Allocate and zero Descriptor Table, Available and Used rings for the virtqueue in contiguous physical memory.
5. Optionally, if MSI-X capability is present and enabled on the device, select a vector to use to request interrupts triggered by virtqueue events. Write the MSI-X Table entry number corresponding to this vector into *queue\_msix\_vector*. Read *queue\_msix\_vector*: on success, previously written value is returned; on failure, NO\_VECTOR value is returned.

#### 4.1.5.1.3.1 Legacy Interface: A Note on Virtqueue Configuration

When using the legacy interface, the queue layout follows [2.7.2 Legacy Interfaces: A Note on Virtqueue Layout](#) with an alignment of 4096. Driver writes the physical address, divided by 4096 to the Queue Address field<sup>2</sup>. There was no mechanism to negotiate the queue size.

#### 4.1.5.2 Available Buffer Notifications

When VIRTIO\_F\_NOTIFICATION\_DATA has not been negotiated, the driver sends an available buffer notification to the device by writing the 16-bit virtqueue index of this virtqueue to the Queue Notify address.

When VIRTIO\_F\_NOTIFICATION\_DATA has been negotiated, the driver sends an available buffer notification to the device by writing the following 32-bit value to the Queue Notify address:

```
le32 {
    vqn : 16;
    next_off : 15;
    next_wrap : 1;
};
```

See [2.9 Driver Notifications](#) for the definition of the components.

See [4.1.4.4](#) for how to calculate the Queue Notify address.

#### 4.1.5.2.1 Driver Requirements: Available Buffer Notifications

If VIRTIO\_F\_NOTIF\_CONFIG\_DATA has been negotiated:

- If VIRTIO\_F\_NOTIFICATION\_DATA has not been negotiated, the driver MUST use the *queue\_notify\_data* value instead of the virtqueue index.
- If VIRTIO\_F\_NOTIFICATION\_DATA has been negotiated, the driver MUST set the *vqn* field to the *queue\_notify\_data* value.

#### 4.1.5.3 Used Buffer Notifications

If a used buffer notification is necessary for a virtqueue, the device would typically act as follows:

- If MSI-X capability is disabled:
  1. Set the lower bit of the ISR Status field for the device.
  2. Send the appropriate PCI interrupt for the device.
- If MSI-X capability is enabled:
  1. If *queue\_msix\_vector* is not NO\_VECTOR, request the appropriate MSI-X interrupt message for the device, *queue\_msix\_vector* sets the MSI-X Table entry number.

#### 4.1.5.3.1 Device Requirements: Used Buffer Notifications

If MSI-X capability is enabled and *queue\_msix\_vector* is NO\_VECTOR for a virtqueue, the device MUST NOT deliver an interrupt for that virtqueue.

<sup>2</sup>The 4096 is based on the x86 page size, but it's also large enough to ensure that the separate parts of the virtqueue are on separate cache lines.

#### 4.1.5.4 Notification of Device Configuration Changes

Some virtio PCI devices can change the device configuration state, as reflected in the device-specific configuration region of the device. In this case:

- If MSI-X capability is disabled:
  1. Set the second lower bit of the ISR Status field for the device.
  2. Send the appropriate PCI interrupt for the device.
- If MSI-X capability is enabled:
  1. If *config\_msix\_vector* is not NO\_VECTOR, request the appropriate MSI-X interrupt message for the device, *config\_msix\_vector* sets the MSI-X Table entry number.

A single interrupt MAY indicate both that one or more virtqueue has been used and that the configuration space has changed.

##### 4.1.5.4.1 Device Requirements: Notification of Device Configuration Changes

If MSI-X capability is enabled and *config\_msix\_vector* is NO\_VECTOR, the device MUST NOT deliver an interrupt for device configuration space changes.

##### 4.1.5.4.2 Driver Requirements: Notification of Device Configuration Changes

A driver MUST handle the case where the same interrupt is used to indicate both device configuration space change and one or more virtqueues being used.

#### 4.1.5.5 Driver Handling Interrupts

The driver interrupt handler would typically:

- If MSI-X capability is disabled:
  - Read the ISR Status field, which will reset it to zero.
  - If the lower bit is set: look through all virtqueues for the device, to see if any progress has been made by the device which requires servicing.
  - If the second lower bit is set: re-examine the configuration space to see what changed.
- If MSI-X capability is enabled:
  - Look through all virtqueues mapped to that MSI-X vector for the device, to see if any progress has been made by the device which requires servicing.
  - If the MSI-X vector is equal to *config\_msix\_vector*, re-examine the configuration space to see what changed.

## 4.2 Virtio Over MMIO

Virtual environments without PCI support (a common situation in embedded devices models) might use simple memory mapped device (“virtio-mmio”) instead of the PCI device.

The memory mapped virtio device behaviour is based on the PCI device specification. Therefore most operations including device initialization, queues configuration and buffer transfers are nearly identical. Existing differences are described in the following sections.

### 4.2.1 MMIO Device Discovery

Unlike PCI, MMIO provides no generic device discovery mechanism. For each device, the guest OS will need to know the location of the registers and interrupt(s) used. The suggested binding for systems using flattened device trees is shown in this example:

```
// EXAMPLE: virtio_block device taking 512 bytes at 0x1e000, interrupt 42.
virtio_block@1e000 {
    compatible = "virtio,mmio";
    reg = <0x1e000 0x200>;
    interrupts = <42>;
}
```

## 4.2.2 MMIO Device Register Layout

MMIO virtio devices provide a set of memory mapped control registers followed by a device-specific configuration space, described in the table 4.1.

All register values are organized as Little Endian.

Table 4.1: MMIO Device Register Layout

<i>Name</i> Offset from base Direction	<b>Function</b> Description
<i>MagicValue</i> 0x000 R	<b>Magic value</b> 0x74726976 (a Little Endian equivalent of the “virt” string).
<i>Version</i> 0x004 R	<b>Device version number</b> 0x2. <b>Note:</b> Legacy devices (see 4.2.4 Legacy interface) used 0x1.
<i>DeviceID</i> 0x008 R	<b>Virtio Subsystem Device ID</b> See 5 Device Types for possible values. Value zero (0x0) is used to define a system memory map with placeholder devices at static, well known addresses, assigning functions to them depending on user’s needs.
<i>VendorID</i> 0x00c R	<b>Virtio Subsystem Vendor ID</b>
<i>DeviceFeatures</i> 0x010 R	<b>Flags representing features the device supports</b> Reading from this register returns 32 consecutive flag bits, the least significant bit depending on the last value written to <i>DeviceFeaturesSel</i> . Access to this register returns bits <i>DeviceFeaturesSel</i> * 32 to ( <i>DeviceFeaturesSel</i> * 32) + 31, eg. feature bits 0 to 31 if <i>DeviceFeaturesSel</i> is set to 0 and features bits 32 to 63 if <i>DeviceFeaturesSel</i> is set to 1. Also see 2.2 Feature Bits.
<i>DeviceFeaturesSel</i> 0x014 W	<b>Device (host) features word selection.</b> Writing to this register selects a set of 32 device feature bits accessible by reading from <i>DeviceFeatures</i> .
<i>DriverFeatures</i> 0x020 W	<b>Flags representing device features understood and activated by the driver</b> Writing to this register sets 32 consecutive flag bits, the least significant bit depending on the last value written to <i>DriverFeaturesSel</i> . Access to this register sets bits <i>DriverFeaturesSel</i> * 32 to ( <i>DriverFeaturesSel</i> * 32) + 31, eg. feature bits 0 to 31 if <i>DriverFeaturesSel</i> is set to 0 and features bits 32 to 63 if <i>DriverFeaturesSel</i> is set to 1. Also see 2.2 Feature Bits.
<i>DriverFeaturesSel</i> 0x024 W	<b>Activated (guest) features word selection</b> Writing to this register selects a set of 32 activated feature bits accessible by writing to <i>DriverFeatures</i> .

Name Offset from the base Direction	Function Description
<i>QueueSel</i> 0x030 W	<b>Virtual queue index</b> Writing to this register selects the virtual queue that the following operations on <i>QueueNumMax</i> , <i>QueueNum</i> , <i>QueueReady</i> , <i>QueueDescLow</i> , <i>QueueDescHigh</i> , <i>QueueDriverLow</i> , <i>QueueDriverHigh</i> , <i>QueueDeviceLow</i> , <i>QueueDeviceHigh</i> and <i>QueueReset</i> apply to. The index number of the first queue is zero (0x0).
<i>QueueNumMax</i> 0x034 R	<b>Maximum virtual queue size</b> Reading from the register returns the maximum size (number of elements) of the queue the device is ready to process or zero (0x0) if the queue is not available. This applies to the queue selected by writing to <i>QueueSel</i> .
<i>QueueNum</i> 0x038 W	<b>Virtual queue size</b> Queue size is the number of elements in the queue. Writing to this register notifies the device what size of the queue the driver will use. This applies to the queue selected by writing to <i>QueueSel</i> .
<i>QueueReady</i> 0x044 RW	<b>Virtual queue ready bit</b> Writing one (0x1) to this register notifies the device that it can execute requests from this virtual queue. Reading from this register returns the last value written to it. Both read and write accesses apply to the queue selected by writing to <i>QueueSel</i> .
<i>QueueNotify</i> 0x050 W	<b>Queue notifier</b> Writing a value to this register notifies the device that there are new buffers to process in a queue. When VIRTIO_F_NOTIFICATION_DATA has not been negotiated, the value written is the queue index. When VIRTIO_F_NOTIFICATION_DATA has been negotiated, the <i>Notification data</i> value has the following format: <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <pre>le32 {     vqn : 16;     next_off : 15;     next_wrap : 1; };</pre> </div> See <a href="#">2.9 Driver Notifications</a> for the definition of the components.
<i>InterruptStatus</i> 0x60 R	<b>Interrupt status</b> Reading from this register returns a bit mask of events that caused the device interrupt to be asserted. The following events are possible: <b>Used Buffer Notification</b> - bit 0 - the interrupt was asserted because the device has used a buffer in at least one of the active virtual queues. <b>Configuration Change Notification</b> - bit 1 - the interrupt was asserted because the configuration of the device has changed.
<i>InterruptACK</i> 0x064 W	<b>Interrupt acknowledge</b> Writing a value with bits set as defined in <i>InterruptStatus</i> to this register notifies the device that events causing the interrupt have been handled.
<i>Status</i> 0x070 RW	<b>Device status</b> Reading from this register returns the current device status flags. Writing non-zero values to this register sets the status flags, indicating the driver progress. Writing zero (0x0) to this register triggers a device reset. See also p. <a href="#">4.2.3.1 Device Initialization</a> .

Name Offset from the base Direction	Function Description
<i>QueueDescLow</i> 0x080 <i>QueueDescHigh</i> 0x084 W	<b>Virtual queue's Descriptor Area 64 bit long physical address</b> Writing to these two registers (lower 32 bits of the address to <i>QueueDescLow</i> , higher 32 bits to <i>QueueDescHigh</i> ) notifies the device about location of the Descriptor Area of the queue selected by writing to <i>QueueSel</i> register.
<i>QueueDriverLow</i> 0x090 <i>QueueDriverHigh</i> 0x094 W	<b>Virtual queue's Driver Area 64 bit long physical address</b> Writing to these two registers (lower 32 bits of the address to <i>QueueDriverLow</i> , higher 32 bits to <i>QueueDriverHigh</i> ) notifies the device about location of the Driver Area of the queue selected by writing to <i>QueueSel</i> .
<i>QueueDeviceLow</i> 0x0a0 <i>QueueDeviceHigh</i> 0x0a4 W	<b>Virtual queue's Device Area 64 bit long physical address</b> Writing to these two registers (lower 32 bits of the address to <i>QueueDeviceLow</i> , higher 32 bits to <i>QueueDeviceHigh</i> ) notifies the device about location of the Device Area of the queue selected by writing to <i>QueueSel</i> .
<i>SHMSel</i> 0x0ac W	<b>Shared memory id</b> Writing to this register selects the shared memory region <a href="#">2.10</a> following operations on <i>SHMLenLow</i> , <i>SHMLenHigh</i> , <i>SHMBaseLow</i> and <i>SHMBaseHigh</i> apply to.
<i>SHMLenLow</i> 0x0b0 <i>SHMLenHigh</i> 0x0b4 R	<b>Shared memory region 64 bit long length</b> These registers return the length of the shared memory region in bytes, as defined by the device for the region selected by the <i>SHMSel</i> register. The lower 32 bits of the length are read from <i>SHMLenLow</i> and the higher 32 bits from <i>SHMLenHigh</i> . Reading from a non-existent region (i.e. where the ID written to <i>SHMSel</i> is unused) results in a length of -1.
<i>SHMBaseLow</i> 0x0b8 <i>SHMBaseHigh</i> 0x0bc R	<b>Shared memory region 64 bit long physical address</b> The driver reads these registers to discover the base address of the region in physical address space. This address is chosen by the device (or other part of the VMM). The lower 32 bits of the address are read from <i>SHMBaseLow</i> with the higher 32 bits from <i>SHMBaseHigh</i> . Reading from a non-existent region (i.e. where the ID written to <i>SHMSel</i> is unused) results in a base address of 0xffffffffffff.
<i>QueueReset</i> 0x0c0 RW	<b>Virtual queue reset bit</b> If VIRTIO_F_RING_RESET has been negotiated, writing one (0x1) to this register selectively resets the queue. Both read and write accesses apply to the queue selected by writing to <i>QueueSel</i> .
<i>ConfigGeneration</i> 0x0fc R	<b>Configuration atomicity value</b> Reading from this register returns a value describing a version of the device-specific configuration space (see <i>Config</i> ). The driver can then access the configuration space and, when finished, read <i>ConfigGeneration</i> again. If no part of the configuration space has changed between these two <i>ConfigGeneration</i> reads, the returned values are identical. If the values are different, the configuration space accesses were not atomic and the driver has to perform the operations again. See also <a href="#">2.5</a> .
<i>Config</i> 0x100+ RW	<b>Configuration space</b> Device-specific configuration space starts at the offset 0x100 and is accessed with byte alignment. Its meaning and size depend on the device and the driver.



#### 4.2.2.1 Device Requirements: MMIO Device Register Layout

The device MUST return 0x74726976 in *MagicValue*.

The device MUST return value 0x2 in *Version*.

The device MUST present each event by setting the corresponding bit in *InterruptStatus* from the moment it takes place, until the driver acknowledges the interrupt by writing a corresponding bit mask to the *InterruptACK* register. Bits which do not represent events which took place MUST be zero.

Upon reset, the device MUST clear all bits in *InterruptStatus* and ready bits in the *QueueReady* register for all queues in the device.

The device MUST change value returned in *ConfigGeneration* if there is any risk of a driver seeing an inconsistent configuration state.

The device MUST NOT access virtual queue contents when *QueueReady* is zero (0x0).

If VIRTIO\_F\_RING\_RESET has been negotiated, the device MUST present a 0 in *QueueReady* after the driver has reset the virtqueue via *QueueReset*.

If VIRTIO\_F\_RING\_RESET has been negotiated, the device MUST present a 0 in *QueueReset* on reset.

If VIRTIO\_F\_RING\_RESET has been negotiated, The device MUST present a 0 in *QueueReset* after the virtqueue is enabled with *QueueReady*.

The device MUST reset the queue when 1 is written to *QueueReset*, and present a 1 in *QueueReset* after the queue has been reset, until the driver re-enables the queue via *QueueReady* or until the device is reset. The device MUST present consistent default values after queue reset. (see 2.6.1).

#### 4.2.2.2 Driver Requirements: MMIO Device Register Layout

The driver MUST NOT access memory locations not described in the table 4.1 (or, in case of the configuration space, described in the device specification), MUST NOT write to the read-only registers (direction R) and MUST NOT read from the write-only registers (direction W).

The driver MUST only use 32 bit wide and aligned reads and writes to access the control registers described in table 4.1. For the device-specific configuration space, the driver MUST use 8 bit wide accesses for 8 bit wide fields, 16 bit wide and aligned accesses for 16 bit wide fields and 32 bit wide and aligned accesses for 32 and 64 bit wide fields.

The driver MUST ignore a device with *MagicValue* which is not 0x74726976, although it MAY report an error.

The driver MUST ignore a device with *Version* which is not 0x2, although it MAY report an error.

The driver MUST ignore a device with *DeviceID* 0x0, but MUST NOT report any error.

Before reading from *DeviceFeatures*, the driver MUST write a value to *DeviceFeaturesSel*.

Before writing to the *DriverFeatures* register, the driver MUST write a value to the *DriverFeaturesSel* register.

The driver MUST write a value to *QueueNum* which is less than or equal to the value presented by the device in *QueueNumMax*.

When *QueueReady* is not zero, the driver MUST NOT access *QueueNum*, *QueueDescLow*, *QueueDescHigh*, *QueueDriverLow*, *QueueDriverHigh*, *QueueDeviceLow*, *QueueDeviceHigh*.

To stop using the queue the driver MUST write zero (0x0) to this *QueueReady* and MUST read the value back to ensure synchronization.

The driver MUST ignore undefined bits in *InterruptStatus*.

The driver MUST write a value with a bit mask describing events it handled into *InterruptACK* when it finishes handling an interrupt and MUST NOT set any of the undefined bits in the value.

If VIRTIO\_F\_RING\_RESET has been negotiated, after the driver writes 1 to *QueueReset* to reset the queue, it MUST verify that the queue has been reset by reading back *QueueReset* and ensuring that it is 1. The driver MAY re-enable the queue by writing a 1 to *QueueReady* after ensuring that the other virtqueue fields

have been set up correctly. The driver MAY set driver-writable queue configuration values to different values than those that were used before the queue reset. (see [2.6.1](#)).

## 4.2.3 MMIO-specific Initialization And Device Operation

### 4.2.3.1 Device Initialization

#### 4.2.3.1.1 Driver Requirements: Device Initialization

The driver MUST start the device initialization by reading and checking values from *MagicValue* and *Version*. If both values are valid, it MUST read *DeviceID* and if its value is zero (0x0) MUST abort initialization and MUST NOT access any other register.

Drivers not expecting shared memory MUST NOT use the shared memory registers.

Further initialization MUST follow the procedure described in [3.1 Device Initialization](#).

### 4.2.3.2 Virtqueue Configuration

The driver will typically initialize the virtual queue in the following way:

1. Select the queue writing its index (first queue is 0) to *QueueSel*.
2. Check if the queue is not already in use: read *QueueReady*, and expect a returned value of zero (0x0).
3. Read maximum queue size (number of elements) from *QueueNumMax*. If the returned value is zero (0x0) the queue is not available.
4. Allocate and zero the queue memory, making sure the memory is physically contiguous.
5. Notify the device about the queue size by writing the size to *QueueNum*.
6. Write physical addresses of the queue's Descriptor Area, Driver Area and Device Area to (respectively) the *QueueDescLow/QueueDescHigh*, *QueueDriverLow/QueueDriverHigh* and *QueueDeviceLow/QueueDeviceHigh* register pairs.
7. Write 0x1 to *QueueReady*.

### 4.2.3.3 Available Buffer Notifications

When *VIRTIO\_F\_NOTIFICATION\_DATA* has not been negotiated, the driver sends an available buffer notification to the device by writing the 16-bit virtqueue index of the queue to be notified to *QueueNotify*.

When *VIRTIO\_F\_NOTIFICATION\_DATA* has been negotiated, the driver sends an available buffer notification to the device by writing the following 32-bit value to *QueueNotify*:

```
le32 {
    vqn : 16;
    next_off : 15;
    next_wrap : 1;
};
```

See [2.9 Driver Notifications](#) for the definition of the components.

### 4.2.3.4 Notifications From The Device

The memory mapped virtio device is using a single, dedicated interrupt signal, which is asserted when at least one of the bits described in the description of *InterruptStatus* is set. This is how the device sends a used buffer notification or a configuration change notification to the device.

#### 4.2.3.4.1 Driver Requirements: Notifications From The Device

After receiving an interrupt, the driver MUST read *InterruptStatus* to check what caused the interrupt (see the register description). The used buffer notification bit being set SHOULD be interpreted as a used buffer

notification for each active virtqueue. After the interrupt is handled, the driver MUST acknowledge it by writing a bit mask corresponding to the handled events to the InterruptACK register.

#### 4.2.4 Legacy interface

The legacy MMIO transport used page-based addressing, resulting in a slightly different control register layout, the device initialization and the virtual queue configuration procedure.

Table 4.2 presents control registers layout, omitting descriptions of registers which did not change their function nor behaviour:

Table 4.2: MMIO Device Legacy Register Layout

<i>Name</i> Offset from base Direction	<b>Function</b> Description
<i>MagicValue</i> 0x000 R	<b>Magic value</b>
<i>Version</i> 0x004 R	<b>Device version number</b> Legacy device returns value 0x1.
<i>DeviceID</i> 0x008 R	<b>Virtio Subsystem Device ID</b>
<i>VendorID</i> 0x00c R	<b>Virtio Subsystem Vendor ID</b>
<i>HostFeatures</i> 0x010 R	<b>Flags representing features the device supports</b>
<i>HostFeaturesSel</i> 0x014 W	<b>Device (host) features word selection.</b>
<i>GuestFeatures</i> 0x020 W	<b>Flags representing device features understood and activated by the driver</b>
<i>GuestFeaturesSel</i> 0x024 W	<b>Activated (guest) features word selection</b>
<i>GuestPageSize</i> 0x028 W	<b>Guest page size</b> The driver writes the guest page size in bytes to the register during initialization, before any queues are used. This value should be a power of 2 and is used by the device to calculate the Guest address of the first queue page (see QueuePFN).
<i>QueueSel</i> 0x030 W	<b>Virtual queue index</b> Writing to this register selects the virtual queue that the following operations on the <i>QueueNumMax</i> , <i>QueueNum</i> , <i>QueueAlign</i> and <i>QueuePFN</i> registers apply to. The index number of the first queue is zero (0x0). .

<i>Name</i> Offset from the base Direction	<b>Function</b> Description
<i>QueueNumMax</i> 0x034 R	<b>Maximum virtual queue size</b> Reading from the register returns the maximum size of the queue the device is ready to process or zero (0x0) if the queue is not available. This applies to the queue selected by writing to <i>QueueSel</i> and is allowed only when <i>QueuePFN</i> is set to zero (0x0), so when the queue is not actively used.
<i>QueueNum</i> 0x038 W	<b>Virtual queue size</b> Queue size is the number of elements in the queue. Writing to this register notifies the device what size of the queue the driver will use. This applies to the queue selected by writing to <i>QueueSel</i> .
<i>QueueAlign</i> 0x03c W	<b>Used Ring alignment in the virtual queue</b> Writing to this register notifies the device about alignment boundary of the Used Ring in bytes. This value should be a power of 2 and applies to the queue selected by writing to <i>QueueSel</i> .
<i>QueuePFN</i> 0x040 RW	<b>Guest physical page number of the virtual queue</b> Writing to this register notifies the device about location of the virtual queue in the Guest's physical address space. This value is the index number of a page starting with the queue Descriptor Table. Value zero (0x0) means physical address zero (0x00000000) and is illegal. When the driver stops using the queue it writes zero (0x0) to this register. Reading from this register returns the currently used page number of the queue, therefore a value other than zero (0x0) means that the queue is in use. Both read and write accesses apply to the queue selected by writing to <i>QueueSel</i> .
<i>QueueNotify</i> 0x050 W	<b>Queue notifier</b>
<i>InterruptStatus</i> 0x60 R	<b>Interrupt status</b>
<i>InterruptACK</i> 0x064 W	<b>Interrupt acknowledge</b>
<i>Status</i> 0x070 RW	<b>Device status</b> Reading from this register returns the current device status flags. Writing non-zero values to this register sets the status flags, indicating the OS/driver progress. Writing zero (0x0) to this register triggers a device reset. The device sets <i>QueuePFN</i> to zero (0x0) for all queues in the device. Also see <a href="#">3.1 Device Initialization</a> .
<i>Config</i> 0x100+ RW	<b>Configuration space</b>

The virtual queue page size is defined by writing to *GuestPageSize*, as written by the guest. The driver does this before the virtual queues are configured.

The virtual queue layout follows p. [2.7.2 Legacy Interfaces: A Note on Virtqueue Layout](#), with the alignment defined in *QueueAlign*.

The virtual queue is configured as follows:

1. Select the queue writing its index (first queue is 0) to *QueueSel*.
2. Check if the queue is not already in use: read *QueuePFN*, expecting a returned value of zero (0x0).

3. Read maximum queue size (number of elements) from *QueueNumMax*. If the returned value is zero (0x0) the queue is not available.
4. Allocate and zero the queue pages in contiguous virtual memory, aligning the Used Ring to an optimal boundary (usually page size). The driver should choose a queue size smaller than or equal to *QueueNumMax*.
5. Notify the device about the queue size by writing the size to *QueueNum*.
6. Notify the device about the used alignment by writing its value in bytes to *QueueAlign*.
7. Write the physical number of the first page of the queue to the *QueuePFN* register.

Notification mechanisms did not change.

## 4.3 Virtio Over Channel I/O

S/390 based virtual machines support neither PCI nor MMIO, so a different transport is needed there.

virtio-ccw uses the standard channel I/O based mechanism used for the majority of devices on S/390. A virtual channel device with a special control unit type acts as proxy to the virtio device (similar to the way virtio-pci uses a PCI device) and configuration and operation of the virtio device is accomplished (mostly) via channel commands. This means virtio devices are discoverable via standard operating system algorithms, and adding virtio support is mainly a question of supporting a new control unit type.

As the S/390 is a big endian machine, the data structures transmitted via channel commands are big-endian: this is made clear by use of the types be16, be32 and be64.

### 4.3.1 Basic Concepts

As a proxy device, virtio-ccw uses a channel-attached I/O control unit with a special control unit type (0x3832) and a control unit model corresponding to the attached virtio device's subsystem device ID, accessed via a virtual I/O subchannel and a virtual channel path of type 0x32. This proxy device is discoverable via normal channel subsystem device discovery (usually a STORE SUBCHANNEL loop) and answers to the basic channel commands:

- NO-OPERATION (0x03)
- BASIC SENSE (0x04)
- TRANSFER IN CHANNEL (0x08)
- SENSE ID (0xe4)

For a virtio-ccw proxy device, SENSE ID will return the following information:

Bytes	Description	Contents
0	reserved	0xff
1-2	control unit type	0x3832
3	control unit model	<virtio device id>
4-5	device type	zeroes (unset)
6	device model	zeroes (unset)
7-255	extended SenseId data	zeroes (unset)

A virtio-ccw proxy device facilitates:

- Discovery and attachment of virtio devices (as described above).
- Initialization of virtqueues and transport-specific facilities (using virtio-specific channel commands).
- Notifications (via hypercall and a combination of I/O interrupts and indicator bits).

#### 4.3.1.1 Channel Commands for Virtio

In addition to the basic channel commands, virtio-ccw defines a set of channel commands related to configuration and operation of virtio:

```
#define CCW_CMD_SET_VQ 0x13
#define CCW_CMD_VDEV_RESET 0x33
#define CCW_CMD_SET_IND 0x43
#define CCW_CMD_SET_CONF_IND 0x53
#define CCW_CMD_SET_IND_ADAPTER 0x73
#define CCW_CMD_READ_FEAT 0x12
#define CCW_CMD_WRITE_FEAT 0x11
#define CCW_CMD_READ_CONF 0x22
#define CCW_CMD_WRITE_CONF 0x21
#define CCW_CMD_WRITE_STATUS 0x31
#define CCW_CMD_READ_VQ_CONF 0x32
#define CCW_CMD_SET_VIRTIO_REV 0x83
#define CCW_CMD_READ_STATUS 0x72
```

#### 4.3.1.2 Notifications

Available buffer notifications are realized as a hypercall. No additional setup by the driver is needed. The operation of available buffer notifications is described in section 4.3.3.2.

Used buffer notifications are realized either as so-called classic or adapter I/O interrupts depending on a transport level negotiation. The initialization is described in sections 4.3.2.6.1 and 4.3.2.6.3 respectively. The operation of each flavor is described in sections 4.3.3.1.1 and 4.3.3.1.2 respectively.

Configuration change notifications are done using so-called classic I/O interrupts. The initialization is described in section 4.3.2.6.2 and the operation in section 4.3.3.1.1.

#### 4.3.1.3 Device Requirements: Basic Concepts

The virtio-ccw device acts like a normal channel device, as specified in [S390 PoP] and [S390 Common I/O]. In particular:

- A device MUST post a unit check with command reject for any command it does not support.
- If a driver did not suppress length checks for a channel command, the device MUST present a sub-channel status as detailed in the architecture when the actual length did not match the expected length.
- If a driver did suppress length checks for a channel command, the device MUST present a check condition if the transmitted data does not contain enough data to process the command. If the driver submitted a buffer that was too long, the device SHOULD accept the command.

#### 4.3.1.4 Driver Requirements: Basic Concepts

A driver for virtio-ccw devices MUST check for a control unit type of 0x3832 and MUST ignore the device type and model.

A driver SHOULD attempt to provide the correct length in a channel command even if it suppresses length checks for that command.

### 4.3.2 Device Initialization

virtio-ccw uses several channel commands to set up a device.

#### 4.3.2.1 Setting the Virtio Revision

CCW\_CMD\_SET\_VIRTIO\_REV is issued by the driver to set the revision of the virtio-ccw transport it intends to drive the device with. It uses the following communication structure:

```
struct virtio_rev_info {
    be16 revision;
    be16 length;
```

```
};      u8 data[];
```

*revision* contains the desired revision id, *length* the length of the data portion and *data* revision-dependent additional desired options.

The following values are supported:

<i>revision</i>	<i>length</i>	<i>data</i>	remarks
0	0	<empty>	legacy interface; transitional devices only
1	0	<empty>	Virtio 1
2	0	<empty>	CCW_CMD_READ_STATUS support
3-n			reserved for later revisions

Note that a change in the virtio standard does not necessarily correspond to a change in the virtio-ccw revision.

#### 4.3.2.1.1 Device Requirements: Setting the Virtio Revision

A device **MUST** post a unit check with command reject for any *revision* it does not support. For any invalid combination of *revision*, *length* and *data*, it **MUST** post a unit check with command reject as well. A non-transitional device **MUST** reject revision id 0.

A device **SHOULD** answer with command reject to any virtio-ccw specific channel command that is not contained in the revision selected by the driver.

A device **MUST** answer with command reject to any attempt to select a different revision after a revision has been successfully selected by the driver.

A device **MUST** treat the revision as unset from the time the associated subchannel has been enabled until a revision has been successfully set by the driver. This implies that revisions are not persistent across disabling and enabling of the associated subchannel.

#### 4.3.2.1.2 Driver Requirements: Setting the Virtio Revision

A driver **SHOULD** start with trying to set the highest revision it supports and continue with lower revisions if it gets a command reject.

A driver **MUST NOT** issue any other virtio-ccw specific channel commands prior to setting the revision.

After a revision has been successfully selected by the driver, it **MUST NOT** attempt to select a different revision.

#### 4.3.2.1.3 Legacy Interfaces: A Note on Setting the Virtio Revision

A legacy device will not support the CCW\_CMD\_SET\_VIRTIO\_REV and answer with a command reject. A non-transitional driver **MUST** stop trying to operate this device in that case. A transitional driver **MUST** operate the device as if it had been able to set revision 0.

A legacy driver will not issue the CCW\_CMD\_SET\_VIRTIO\_REV prior to issuing other virtio-ccw specific channel commands. A non-transitional device therefore **MUST** answer any such attempts with a command reject. A transitional device **MUST** assume in this case that the driver is a legacy driver and continue as if the driver selected revision 0. This implies that the device **MUST** reject any command not valid for revision 0, including a subsequent CCW\_CMD\_SET\_VIRTIO\_REV.

#### 4.3.2.2 Configuring a Virtqueue

CCW\_CMD\_READ\_VQ\_CONF is issued by the driver to obtain information about a queue. It uses the following structure for communicating:



```
struct vq_config_block {
    be16 index;
    be16 max_num;
};
```

The requested number of buffers for queue *index* is returned in *max\_num*.

Afterwards, CCW\_CMD\_SET\_VQ is issued by the driver to inform the device about the location used for its queue. The transmitted structure is

```
struct vq_info_block {
    be64 desc;
    be32 res0;
    be16 index;
    be16 num;
    be64 driver;
    be64 device;
};
```

*desc*, *driver* and *device* contain the guest addresses for the descriptor area, available area and used area for queue *index*, respectively. The actual virtqueue size (number of allocated buffers) is transmitted in *num*.

#### 4.3.2.2.1 Device Requirements: Configuring a Virtqueue

*res0* is reserved and MUST be ignored by the device.

#### 4.3.2.2.2 Legacy Interface: A Note on Configuring a Virtqueue

For a legacy driver or for a driver that selected revision 0, CCW\_CMD\_SET\_VQ uses the following communication block:

```
struct vq_info_block_legacy {
    be64 queue;
    be32 align;
    be16 index;
    be16 num;
};
```

*queue* contains the guest address for queue *index*, *num* the number of buffers and *align* the alignment. The queue layout follows [2.7.2 Legacy Interfaces: A Note on Virtqueue Layout](#).

#### 4.3.2.3 Communicating Status Information

The driver changes the status of a device via the CCW\_CMD\_WRITE\_STATUS command, which transmits an 8 bit status value.

As described in [2.2.2](#), a device sometimes fails to set the *device status* field: For example, it might fail to accept the FEATURES\_OK status bit during device initialization.

With revision 2, CCW\_CMD\_READ\_STATUS is defined: It reads an 8 bit status value from the device and acts as a reverse operation to CCW\_CMD\_WRITE\_STATUS.

##### 4.3.2.3.1 Driver Requirements: Communicating Status Information

If the device posts a unit check with command reject in response to the CCW\_CMD\_WRITE\_STATUS command, the driver MUST assume that the device failed to set the status and the *device status* field retained its previous value.

If at least revision 2 has been negotiated, the driver SHOULD use the CCW\_CMD\_READ\_STATUS command to retrieve the *device status* field after a configuration change has been detected.

If not at least revision 2 has been negotiated, the driver MUST NOT attempt to issue the CCW\_CMD\_READ\_STATUS command.



#### 4.3.2.3.2 Device Requirements: Communicating Status Information

If the device fails to set the *device status* field to the value written by the driver, the device MUST assure that the *device status* field is left unchanged and MUST post a unit check with command reject.

If at least revision 2 has been negotiated, the device MUST return the current *device status* field if the CCW\_CMD\_READ\_STATUS command is issued.

#### 4.3.2.4 Handling Device Features

Feature bits are arranged in an array of 32 bit values, making for a total of 8192 feature bits. Feature bits are in little-endian byte order.

The CCW commands dealing with features use the following communication block:

```
struct virtio_feature_desc {  
    le32 features;  
    u8 index;  
};
```

*features* are the 32 bits of features currently accessed, while *index* describes which of the feature bit values is to be accessed. No padding is added at the end of the structure, it is exactly 5 bytes in length.

The guest obtains the device's device feature set via the CCW\_CMD\_READ\_FEAT command. The device stores the features at *index* to *features*.

For communicating its supported features to the device, the driver uses the CCW\_CMD\_WRITE\_FEAT command, denoting a *features/index* combination.

#### 4.3.2.5 Device Configuration

The device's configuration space is located in host memory.

To obtain information from the configuration space, the driver uses CCW\_CMD\_READ\_CONF, specifying the guest memory for the device to write to.

For changing configuration information, the driver uses CCW\_CMD\_WRITE\_CONF, specifying the guest memory for the device to read from.

In both cases, the complete configuration space is transmitted. This allows the driver to compare the new configuration space with the old version, and keep a generation count internally whenever it changes.

#### 4.3.2.6 Setting Up Indicators

In order to set up the indicator bits for host->guest notification, the driver uses different channel commands depending on whether it wishes to use traditional I/O interrupts tied to a subchannel or adapter I/O interrupts for virtqueue notifications. For any given device, the two mechanisms are mutually exclusive.

For the configuration change indicators, only a mechanism using traditional I/O interrupts is provided, regardless of whether traditional or adapter I/O interrupts are used for virtqueue notifications.

##### 4.3.2.6.1 Setting Up Classic Queue Indicators

Indicators for notification via classic I/O interrupts are contained in a 64 bit value per virtio-ccw proxy device.

To communicate the location of the indicator bits for host->guest notification, the driver uses the CCW\_CMD\_SET\_IND command, pointing to a location containing the guest address of the indicators in a 64 bit value.

If the driver has already set up two-staged queue indicators via the CCW\_CMD\_SET\_IND\_ADAPTER command, the device MUST post a unit check with command reject to any subsequent CCW\_CMD\_SET\_IND command.

#### 4.3.2.6.2 Setting Up Configuration Change Indicators

Indicators for configuration change host->guest notification are contained in a 64 bit value per virtio-ccw proxy device.

To communicate the location of the indicator bits used in the configuration change host->guest notification, the driver issues the CCW\_CMD\_SET\_CONF\_IND command, pointing to a location containing the guest address of the indicators in a 64 bit value.

#### 4.3.2.6.3 Setting Up Two-Stage Queue Indicators

Indicators for notification via adapter I/O interrupts consist of two stages:

- a summary indicator byte covering the virtqueues for one or more virtio-ccw proxy devices
- a set of contiguous indicator bits for the virtqueues for a virtio-ccw proxy device

To communicate the location of the summary and queue indicator bits, the driver uses the CCW\_CMD\_SET\_IND\_ADAPTER command with the following payload:

```
struct virtio_thinint_area {  
    be64 summary_indicator;  
    be64 indicator;  
    be64 bit_nr;  
    u8 isc;  
} __attribute__((packed));
```

*summary\_indicator* contains the guest address of the 8 bit summary indicator. *indicator* contains the guest address of an area wherein the indicators for the devices are contained, starting at *bit\_nr*, one bit per virtqueue of the device. Bit numbers start at the left, i.e. the most significant bit in the first byte is assigned the bit number 0. *isc* contains the I/O interruption subclass to be used for the adapter I/O interrupt. It MAY be different from the isc used by the proxy virtio-ccw device's subchannel. No padding is added at the end of the structure, it is exactly 25 bytes in length.

##### 4.3.2.6.3.1 Device Requirements: Setting Up Two-Stage Queue Indicators

If the driver has already set up classic queue indicators via the CCW\_CMD\_SET\_IND command, the device MUST post a unit check with command reject to any subsequent CCW\_CMD\_SET\_IND\_ADAPTER command.

#### 4.3.2.6.4 Legacy Interfaces: A Note on Setting Up Indicators

In some cases, legacy devices will only support classic queue indicators; in that case, they will reject CCW\_CMD\_SET\_IND\_ADAPTER as they don't know that command. Some legacy devices will support two-stage queue indicators, though, and a driver will be able to successfully use CCW\_CMD\_SET\_IND\_ADAPTER to set them up.

### 4.3.3 Device Operation

#### 4.3.3.1 Host->Guest Notification

There are two modes of operation regarding host->guest notification, classic I/O interrupts and adapter I/O interrupts. The mode to be used is determined by the driver by using CCW\_CMD\_SET\_IND respectively CCW\_CMD\_SET\_IND\_ADAPTER to set up queue indicators.

For configuration changes, the driver always uses classic I/O interrupts.

##### 4.3.3.1.1 Notification via Classic I/O Interrupts

If the driver used the CCW\_CMD\_SET\_IND command to set up queue indicators, the device will use classic I/O interrupts for host->guest notification about virtqueue activity.

For notifying the driver of virtqueue buffers, the device sets the corresponding bit in the guest-provided indicators. If an interrupt is not already pending for the subchannel, the device generates an unsolicited I/O interrupt.

If the device wants to notify the driver about configuration changes, it sets bit 0 in the configuration indicators and generates an unsolicited I/O interrupt, if needed. This also applies if adapter I/O interrupts are used for queue notifications.

#### 4.3.3.1.2 Notification via Adapter I/O Interrupts

If the driver used the `CCW_CMD_SET_IND_ADAPTER` command to set up queue indicators, the device will use adapter I/O interrupts for host->guest notification about virtqueue activity.

For notifying the driver of virtqueue buffers, the device sets the bit in the guest-provided indicator area at the corresponding offset. The guest-provided summary indicator is set to 0x01. An adapter I/O interrupt for the corresponding interruption subclass is generated.

The recommended way to process an adapter I/O interrupt by the driver is as follows:

- Process all queue indicator bits associated with the summary indicator.
- Clear the summary indicator, performing a synchronization (memory barrier) afterwards.
- Process all queue indicator bits associated with the summary indicator again.

##### 4.3.3.1.2.1 Device Requirements: Notification via Adapter I/O Interrupts

The device **SHOULD** only generate an adapter I/O interrupt if the summary indicator had not been set prior to notification.

##### 4.3.3.1.2.2 Driver Requirements: Notification via Adapter I/O Interrupts

The driver **MUST** clear the summary indicator after receiving an adapter I/O interrupt before it processes the queue indicators.

#### 4.3.3.1.3 Legacy Interfaces: A Note on Host->Guest Notification

As legacy devices and drivers support only classic queue indicators, host->guest notification will always be done via classic I/O interrupts.

#### 4.3.3.2 Guest->Host Notification

For notifying the device of virtqueue buffers, the driver unfortunately can't use a channel command (the asynchronous characteristics of channel I/O interact badly with the host block I/O backend). Instead, it uses a diagnose 0x500 call with subcode 3 specifying the queue, as follows:

GPR	Input Value	Output Value
1	0x3	
2	Subchannel ID	Host Cookie
3	Notification data	
4	Host Cookie	

When `VIRTIO_F_NOTIFICATION_DATA` has not been negotiated, the *Notification data* contains the Virtqueue number.

When `VIRTIO_F_NOTIFICATION_DATA` has been negotiated, the value has the following format:

```
be32 {
    vqn : 16;
    next_off : 15;
    next_wrap : 1;
```

```
};
```

See [2.9 Driver Notifications](#) for the definition of the components.

#### 4.3.3.2.1 Device Requirements: Guest->Host Notification

The device MUST ignore bits 0-31 (counting from the left) of GPR2. This aligns passing the subchannel ID with the way it is passed for the existing I/O instructions.

The device MAY return a 64-bit host cookie in GPR2 to speed up the notification execution.

#### 4.3.3.2.2 Driver Requirements: Guest->Host Notification

For each notification, the driver SHOULD use GPR4 to pass the host cookie received in GPR2 from the previous notification.

**Note:** For example:

```
info->cookie = do_notify(schid,
                        virtqueue_get_queue_index(vq),
                        info->cookie);
```

#### 4.3.3.3 Resetting Devices

In order to reset a device, a driver sends the CCW\_CMD\_VDEV\_RESET command. This command does not carry any payload.

The device signals completion of the virtio reset operation through successful conclusion of the CCW\_CMD\_VDEV\_RESET channel command. In particular, the command not only triggers the reset operation, but the reset operation is already completed when the operation concludes successfully.

##### 4.3.3.3.1 Device Requirements: Resetting Devices

The device MUST finish the virtio reset operation and reinitialize *device status* to zero before it concludes the CCW\_CMD\_VDEV\_RESET command successfully.

The device MUST NOT send notifications or interact with the queues after it signaled successful conclusion of the CCW\_CMD\_VDEV\_RESET command.

##### 4.3.3.3.2 Driver Requirements: Resetting Devices

The driver MAY consider the virtio reset operation to be complete already after successful conclusion of the CCW\_CMD\_VDEV\_RESET channel command, although it MAY also choose to verify reset completion by reading *device status* via CCW\_CMD\_READ\_STATUS and checking whether it is 0 afterwards.

---

## 5 Device Types

On top of the queues, config space and feature negotiation facilities built into virtio, several devices are defined.

The following device IDs are used to identify different types of virtio devices. Some device IDs are reserved for devices which are not currently defined in this standard.

Discovering what devices are available and their type is bus-dependent.

Device ID	Virtio Device
0	reserved (invalid)
1	network card
2	block device
3	console
4	entropy source
5	memory ballooning (traditional)
6	ioMemory
7	rpmsg
8	SCSI host
9	9P transport
10	mac80211 wlan
11	rproc serial
12	virtio CAIF
13	memory balloon
16	GPU device
17	Timer/Clock device
18	Input device
19	Socket device
20	Crypto device
21	Signal Distribution Module
22	pstore device
23	IOMMU device
24	Memory device
25	Audio device
26	file system device
27	PMEM device
28	RPMB device
29	mac80211 hwsim wireless simulation device
30	Video encoder device
31	Video decoder device
32	SCMI device
33	NitroSecureModule
34	I2C adapter
35	Watchdog
36	CAN device
38	Parameter Server
39	Audio policy device
40	Bluetooth device
41	GPIO device
42	RDMA device

Some of the devices above are unspecified by this document, because they are seen as immature or especially niche. Be warned that some are only specified by the sole existing implementation; they could become

part of a future specification, be abandoned entirely, or live on outside this standard. We shall speak of them no further.

## 5.1 Network Device

The virtio network device is a virtual ethernet card, and is the most complex of the devices supported so far by virtio. It has enhanced rapidly and demonstrates clearly how support for new features are added to an existing device. Empty buffers are placed in one virtqueue for receiving packets, and outgoing packets are enqueued into another for transmission in that order. A third command queue is used to control advanced filtering features.

### 5.1.1 Device ID

1

### 5.1.2 Virtqueues

0 receiveq1

1 transmitq1

...

2(N-1) receiveqN

2(N-1)+1 transmitqN

2N controlq

N=1 if neither VIRTIO\_NET\_F\_MQ nor VIRTIO\_NET\_F\_RSS are negotiated, otherwise N is set by *max\_virtqueue\_pairs*.

controlq only exists if VIRTIO\_NET\_F\_CTRL\_VQ set.

### 5.1.3 Feature bits

**VIRTIO\_NET\_F\_CSUM (0)** Device handles packets with partial checksum. This “checksum offload” is a common feature on modern network cards.

**VIRTIO\_NET\_F\_GUEST\_CSUM (1)** Driver handles packets with partial checksum.

**VIRTIO\_NET\_F\_CTRL\_GUEST\_OFFLOADS (2)** Control channel offloads reconfiguration support.

**VIRTIO\_NET\_F\_MTU(3)** Device maximum MTU reporting is supported. If offered by the device, device advises driver about the value of its maximum MTU. If negotiated, the driver uses *mtu* as the maximum MTU value.

**VIRTIO\_NET\_F\_MAC (5)** Device has given MAC address.

**VIRTIO\_NET\_F\_GUEST\_TSO4 (7)** Driver can receive TSOv4.

**VIRTIO\_NET\_F\_GUEST\_TSO6 (8)** Driver can receive TSOv6.

**VIRTIO\_NET\_F\_GUEST\_ECN (9)** Driver can receive TSO with ECN.

**VIRTIO\_NET\_F\_GUEST\_UFO (10)** Driver can receive UFO.

**VIRTIO\_NET\_F\_HOST\_TSO4 (11)** Device can receive TSOv4.

**VIRTIO\_NET\_F\_HOST\_TSO6 (12)** Device can receive TSOv6.

**VIRTIO\_NET\_F\_HOST\_ECN (13)** Device can receive TSO with ECN.

**VIRTIO\_NET\_F\_HOST\_UFO (14)** Device can receive UFO.

**VIRTIO\_NET\_F\_MRG\_RXBUF (15)** Driver can merge receive buffers.



**VIRTIO\_NET\_F\_STATUS (16)** Configuration status field is available.

**VIRTIO\_NET\_F\_CTRL\_VQ (17)** Control channel is available.

**VIRTIO\_NET\_F\_CTRL\_RX (18)** Control channel RX mode support.

**VIRTIO\_NET\_F\_CTRL\_VLAN (19)** Control channel VLAN filtering.

**VIRTIO\_NET\_F\_GUEST\_ANNOUNCE(21)** Driver can send gratuitous packets.

**VIRTIO\_NET\_F\_MQ(22)** Device supports multiqueue with automatic receive steering.

**VIRTIO\_NET\_F\_CTRL\_MAC\_ADDR(23)** Set MAC address through control channel.

**VIRTIO\_NET\_F\_HOST\_USO (56)** Device can receive USO packets. Unlike UFO (fragmenting the packet) the USO splits large UDP packet to several segments when each of these smaller packets has UDP header.

**VIRTIO\_NET\_F\_HASH\_REPORT(57)** Device can report per-packet hash value and a type of calculated hash.

**VIRTIO\_NET\_F\_GUEST\_HDRLEN(59)** Driver can provide the exact *hdr\_len* value. Device benefits from knowing the exact header length.

**VIRTIO\_NET\_F\_RSS(60)** Device supports RSS (receive-side scaling) with Toeplitz hash calculation and configurable hash parameters for receive steering.

**VIRTIO\_NET\_F\_RSC\_EXT(61)** Device can process duplicated ACKs and report number of coalesced segments and duplicated ACKs.

**VIRTIO\_NET\_F\_STANDBY(62)** Device may act as a standby for a primary device with the same MAC address.

**VIRTIO\_NET\_F\_SPEED\_DUPLEX(63)** Device reports speed and duplex.

#### 5.1.3.1 Feature bit requirements

Some networking feature bits require other networking feature bits (see [2.2.1](#)):

**VIRTIO\_NET\_F\_GUEST\_TSO4** Requires VIRTIO\_NET\_F\_GUEST\_CSUM.

**VIRTIO\_NET\_F\_GUEST\_TSO6** Requires VIRTIO\_NET\_F\_GUEST\_CSUM.

**VIRTIO\_NET\_F\_GUEST\_ECN** Requires VIRTIO\_NET\_F\_GUEST\_TSO4 or VIRTIO\_NET\_F\_GUEST\_TSO6.

**VIRTIO\_NET\_F\_GUEST\_UFO** Requires VIRTIO\_NET\_F\_GUEST\_CSUM.

**VIRTIO\_NET\_F\_HOST\_TSO4** Requires VIRTIO\_NET\_F\_CSUM.

**VIRTIO\_NET\_F\_HOST\_TSO6** Requires VIRTIO\_NET\_F\_CSUM.

**VIRTIO\_NET\_F\_HOST\_ECN** Requires VIRTIO\_NET\_F\_HOST\_TSO4 or VIRTIO\_NET\_F\_HOST\_TSO6.

**VIRTIO\_NET\_F\_HOST\_UFO** Requires VIRTIO\_NET\_F\_CSUM.

**VIRTIO\_NET\_F\_HOST\_USO** Requires VIRTIO\_NET\_F\_CSUM.

**VIRTIO\_NET\_F\_CTRL\_RX** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

**VIRTIO\_NET\_F\_CTRL\_VLAN** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

**VIRTIO\_NET\_F\_GUEST\_ANNOUNCE** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

**VIRTIO\_NET\_F\_MQ** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

**VIRTIO\_NET\_F\_CTRL\_MAC\_ADDR** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

**VIRTIO\_NET\_F\_RSC\_EXT** Requires VIRTIO\_NET\_F\_HOST\_TSO4 or VIRTIO\_NET\_F\_HOST\_TSO6.

**VIRTIO\_NET\_F\_RSS** Requires VIRTIO\_NET\_F\_CTRL\_VQ.

### 5.1.3.2 Legacy Interface: Feature bits

**VIRTIO\_NET\_F\_GSO (6)** Device handles packets with any GSO type. This was supposed to indicate segmentation offload support, but upon further investigation it became clear that multiple bits were needed.

**VIRTIO\_NET\_F\_GUEST\_RSC4 (41)** Device coalesces TCPIP v4 packets. This was implemented by hypervisor patch for certification purposes and current Windows driver depends on it. It will not function if virtio-net device reports this feature.

**VIRTIO\_NET\_F\_GUEST\_RSC6 (42)** Device coalesces TCPIP v6 packets. Similar to VIRTIO\_NET\_F\_GUEST\_RSC4.

### 5.1.4 Device configuration layout

Device configuration fields are listed below, they are read-only for a driver. The *mac* address field always exists (though is only valid if VIRTIO\_NET\_F\_MAC is set), and *status* only exists if VIRTIO\_NET\_F\_STATUS is set. Two read-only bits (for the driver) are currently defined for the status field: VIRTIO\_NET\_S\_LINK\_UP and VIRTIO\_NET\_S\_ANNOUNCE.

```
#define VIRTIO_NET_S_LINK_UP      1
#define VIRTIO_NET_S_ANNOUNCE    2
```

The following driver-read-only field, *max\_virtqueue\_pairs* only exists if VIRTIO\_NET\_F\_MQ or VIRTIO\_NET\_F\_RSS is set. This field specifies the maximum number of each of transmit and receive virtqueues (receiveq1...receiveqN and transmitq1...transmitqN respectively) that can be configured once at least one of these features is negotiated.

The following driver-read-only field, *mtu* only exists if VIRTIO\_NET\_F\_MTU is set. This field specifies the maximum MTU for the driver to use.

The following two fields, *speed* and *duplex*, only exist if VIRTIO\_NET\_F\_SPEED\_DUPLEX is set.

*speed* contains the device speed, in units of 1 MBit per second, 0 to 0x7ffffff, or 0xffffffff for unknown speed.

*duplex* has the values of 0x01 for full duplex, 0x00 for half duplex and 0xff for unknown duplex state.

Both *speed* and *duplex* can change, thus the driver is expected to re-read these values after receiving a configuration change notification.

```
struct virtio_net_config {
    u8 mac[6];
    le16 status;
    le16 max_virtqueue_pairs;
    le16 mtu;
    le32 speed;
    u8 duplex;
    u8 rss_max_key_size;
    le16 rss_max_indirection_table_length;
    le32 supported_hash_types;
};
```

The following field, *rss\_max\_key\_size* only exists if VIRTIO\_NET\_F\_RSS or VIRTIO\_NET\_F\_HASH\_REPORT is set. It specifies the maximum supported length of RSS key in bytes.

The following field, *rss\_max\_indirection\_table\_length* only exists if VIRTIO\_NET\_F\_RSS is set. It specifies the maximum number of 16-bit entries in RSS indirection table.

The next field, *supported\_hash\_types* only exists if the device supports hash calculation, i.e. if VIRTIO\_NET\_F\_RSS or VIRTIO\_NET\_F\_HASH\_REPORT is set.

Field *supported\_hash\_types* contains the bitmask of supported hash types. See 5.1.6.4.3.1 for details of supported hash types.

#### 5.1.4.1 Device Requirements: Device configuration layout

The device MUST set *max\_virtqueue\_pairs* to between 1 and 0x8000 inclusive, if it offers VIRTIO\_NET\_F\_MQ.

The device MUST set *mtu* to between 68 and 65535 inclusive, if it offers VIRTIO\_NET\_F\_MTU.

The device SHOULD set *mtu* to at least 1280, if it offers VIRTIO\_NET\_F\_MTU.

The device MUST NOT modify *mtu* once it has been set.

The device MUST NOT pass received packets that exceed *mtu* (plus low level ethernet header length) size with *gso\_type* NONE or ECN after VIRTIO\_NET\_F\_MTU has been successfully negotiated.

The device MUST forward transmitted packets of up to *mtu* (plus low level ethernet header length) size with *gso\_type* NONE or ECN, and do so without fragmentation, after VIRTIO\_NET\_F\_MTU has been successfully negotiated.

The device MUST set *rss\_max\_key\_size* to at least 40, if it offers VIRTIO\_NET\_F\_RSS or VIRTIO\_NET\_F\_HASH\_REPORT.

The device MUST set *rss\_max\_indirection\_table\_length* to at least 128, if it offers VIRTIO\_NET\_F\_RSS.

If the driver negotiates the VIRTIO\_NET\_F\_STANDBY feature, the device MAY act as a standby device for a primary device with the same MAC address.

If VIRTIO\_NET\_F\_SPEED\_DUPLEX has been negotiated, *speed* MUST contain the device speed, in units of 1 Mbit per second, 0 to 0x7fffffff, or 0xffffffff for unknown.

If VIRTIO\_NET\_F\_SPEED\_DUPLEX has been negotiated, *duplex* MUST have the values of 0x00 for full duplex, 0x01 for half duplex, or 0xff for unknown.

If VIRTIO\_NET\_F\_SPEED\_DUPLEX and VIRTIO\_NET\_F\_STATUS have both been negotiated, the device SHOULD NOT change the *speed* and *duplex* fields as long as VIRTIO\_NET\_S\_LINK\_UP is set in the *status*.

#### 5.1.4.2 Driver Requirements: Device configuration layout

A driver SHOULD negotiate VIRTIO\_NET\_F\_MAC if the device offers it. If the driver negotiates the VIRTIO\_NET\_F\_MAC feature, the driver MUST set the physical address of the NIC to *mac*. Otherwise, it SHOULD use a locally-administered MAC address (see [IEEE 802](#), “9.2 48-bit universal LAN MAC addresses”).

If the driver does not negotiate the VIRTIO\_NET\_F\_STATUS feature, it SHOULD assume the link is active, otherwise it SHOULD read the link status from the bottom bit of *status*.

A driver SHOULD negotiate VIRTIO\_NET\_F\_MTU if the device offers it.

If the driver negotiates VIRTIO\_NET\_F\_MTU, it MUST supply enough receive buffers to receive at least one receive packet of size *mtu* (plus low level ethernet header length) with *gso\_type* NONE or ECN.

If the driver negotiates VIRTIO\_NET\_F\_MTU, it MUST NOT transmit packets of size exceeding the value of *mtu* (plus low level ethernet header length) with *gso\_type* NONE or ECN.

A driver SHOULD negotiate the VIRTIO\_NET\_F\_STANDBY feature if the device offers it.

If VIRTIO\_NET\_F\_SPEED\_DUPLEX has been negotiated, the driver MUST treat any value of *speed* above 0x7fffffff as well as any value of *duplex* not matching 0x00 or 0x01 as an unknown value.

If VIRTIO\_NET\_F\_SPEED\_DUPLEX has been negotiated, the driver SHOULD re-read *speed* and *duplex* after a configuration change notification.

#### 5.1.4.3 Legacy Interface: Device configuration layout

When using the legacy interface, transitional devices and drivers MUST format *status* and *max\_virtqueue\_pairs* in struct *virtio\_net\_config* according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

When using the legacy interface, *mac* is driver-writable which provided a way for drivers to update the MAC without negotiating VIRTIO\_NET\_F\_CTRL\_MAC\_ADDR.

## 5.1.5 Device Initialization

A driver would perform a typical initialization routine like so:

1. Identify and initialize the receive and transmission virtqueues, up to N of each kind. If VIRTIO\_NET\_F\_MQ feature bit is negotiated,  $N = \text{max\_virtqueue\_pairs}$ , otherwise identify  $N=1$ .
2. If the VIRTIO\_NET\_F\_CTRL\_VQ feature bit is negotiated, identify the control virtqueue.
3. Fill the receive queues with buffers: see [5.1.6.3](#).
4. Even with VIRTIO\_NET\_F\_MQ, only receiveq1, transmitq1 and controlq are used by default. The driver would send the VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command specifying the number of the transmit and receive queues to use.
5. If the VIRTIO\_NET\_F\_MAC feature bit is set, the configuration space *mac* entry indicates the “physical” address of the network card, otherwise the driver would typically generate a random local MAC address.
6. If the VIRTIO\_NET\_F\_STATUS feature bit is negotiated, the link status comes from the bottom bit of *status*. Otherwise, the driver assumes it's active.
7. A performant driver would indicate that it will generate checksumless packets by negotiating the VIRTIO\_NET\_F\_CSUM feature.
8. If that feature is negotiated, a driver can use TCP segmentation or UDP segmentation/fragmentation offload by negotiating the VIRTIO\_NET\_F\_HOST\_TSO4 (IPv4 TCP), VIRTIO\_NET\_F\_HOST\_TSO6 (IPv6 TCP), VIRTIO\_NET\_F\_HOST\_UFO (UDP fragmentation) and VIRTIO\_NET\_F\_HOST\_USO (UDP segmentation) features.
9. The converse features are also available: a driver can save the virtual device some work by negotiating these features.

**Note:** For example, a network packet transported between two guests on the same system might not need checksumming at all, nor segmentation, if both guests are amenable. The VIRTIO\_NET\_F\_GUEST\_CSUM feature indicates that partially checksummed packets can be received, and if it can do that then the VIRTIO\_NET\_F\_GUEST\_TSO4, VIRTIO\_NET\_F\_GUEST\_TSO6, VIRTIO\_NET\_F\_GUEST\_UFO and VIRTIO\_NET\_F\_GUEST\_ECN are the input equivalents of the features described above. See [5.1.6.3 Setting Up Receive Buffers](#) and [5.1.6.4 Processing of Incoming Packets](#) below.

A truly minimal driver would only accept VIRTIO\_NET\_F\_MAC and ignore everything else.

## 5.1.6 Device Operation

Packets are transmitted by placing them in the transmitq1...transmitqN, and buffers for incoming packets are placed in the receiveq1...receiveqN. In each case, the packet itself is preceded by a header:

```
struct virtio_net_hdr {
#define VIRTIO_NET_HDR_F_NEEDS_CSUM    1
#define VIRTIO_NET_HDR_F_DATA_VALID   2
#define VIRTIO_NET_HDR_F_RSC_INFO     4
    u8 flags;
#define VIRTIO_NET_HDR_GSO_NONE        0
#define VIRTIO_NET_HDR_GSO_TCPV4      1
#define VIRTIO_NET_HDR_GSO_UDP        3
#define VIRTIO_NET_HDR_GSO_TCPV6      4
#define VIRTIO_NET_HDR_GSO_UDP_L4     5
#define VIRTIO_NET_HDR_GSO_ECN        0x80
    u8 gso_type;
    le16 hdr_len;
    le16 gso_size;
    le16 csum_start;
    le16 csum_offset;
    le16 num_buffers;
    le32 hash_value;           (Only if VIRTIO_NET_F_HASH_REPORT negotiated)
    le16 hash_report;         (Only if VIRTIO_NET_F_HASH_REPORT negotiated)
    le16 padding_reserved;    (Only if VIRTIO_NET_F_HASH_REPORT negotiated)
```

```
};
```

The controlq is used to control device features such as filtering.

#### 5.1.6.1 Legacy Interface: Device Operation

When using the legacy interface, transitional devices and drivers MUST format the fields in struct virtio\_net\_hdr according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

The legacy driver only presented *num\_buffers* in the struct virtio\_net\_hdr when VIRTIO\_NET\_F\_MRG\_RXBUF was negotiated; without that feature the structure was 2 bytes shorter.

When using the legacy interface, the driver SHOULD ignore the used length for the transmit queues and the controlq queue.

**Note:** Historically, some devices put the total descriptor length there, even though no data was actually written.

#### 5.1.6.2 Packet Transmission

Transmitting a single packet is simple, but varies depending on the different features the driver negotiated.

1. The driver can send a completely checksummed packet. In this case, *flags* will be zero, and *gso\_type* will be VIRTIO\_NET\_HDR\_GSO\_NONE.
2. If the driver negotiated VIRTIO\_NET\_F\_CSUM, it can skip checksumming the packet:
  - *flags* has the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM set,
  - *csum\_start* is set to the offset within the packet to begin checksumming, and
  - *csum\_offset* indicates how many bytes after the *csum\_start* the new (16 bit ones' complement) checksum is placed by the device.
  - The TCP checksum field in the packet is set to the sum of the TCP pseudo header, so that replacing it by the ones' complement checksum of the TCP header and body will give the correct result.

**Note:** For example, consider a partially checksummed TCP (IPv4) packet. It will have a 14 byte ethernet header and 20 byte IP header followed by the TCP header (with the TCP checksum field 16 bytes into that header). *csum\_start* will be 14+20 = 34 (the TCP checksum includes the header), and *csum\_offset* will be 16.

3. If the driver negotiated VIRTIO\_NET\_F\_HOST\_TSO4, TSO6, USO or UFO, and the packet requires TCP segmentation, UDP segmentation or fragmentation, then *gso\_type* is set to VIRTIO\_NET\_HDR\_GSO\_TCPV4, TCPV6, UDP\_L4 or UDP. (Otherwise, it is set to VIRTIO\_NET\_HDR\_GSO\_NONE). In this case, packets larger than 1514 bytes can be transmitted: the metadata indicates how to replicate the packet header to cut it into smaller packets. The other gso fields are set:
  - If the VIRTIO\_NET\_F\_GUEST\_HDRLLEN feature has been negotiated, *hdr\_len* indicates the header length that needs to be replicated for each packet. It's the number of bytes from the beginning of the packet to the beginning of the transport payload. Otherwise, if the VIRTIO\_NET\_F\_GUEST\_HDRLLEN feature has not been negotiated, *hdr\_len* is a hint to the device as to how much of the header needs to be kept to copy into each packet, usually set to the length of the headers, including the transport header<sup>1</sup>.

**Note:** Some devices benefit from knowledge of the exact header length.

  - *gso\_size* is the maximum size of each packet beyond that header (ie. MSS).
  - If the driver negotiated the VIRTIO\_NET\_F\_HOST\_ECN feature, the VIRTIO\_NET\_HDR\_GSO\_ECN bit in *gso\_type* indicates that the TCP packet has the ECN bit set<sup>2</sup>.

<sup>1</sup>Due to various bugs in implementations, this field is not useful as a guarantee of the transport header size.

<sup>2</sup>This case is not handled by some older hardware, so is called out specifically in the protocol.

4. *num\_buffers* is set to zero. This field is unused on transmitted packets.
5. The header and packet are added as one output descriptor to the transmitq, and the device is notified of the new entry (see [5.1.5 Device Initialization](#)).

#### 5.1.6.2.1 Driver Requirements: Packet Transmission

The driver MUST set *num\_buffers* to zero.

If VIRTIO\_NET\_F\_CSUM is not negotiated, the driver MUST set *flags* to zero and SHOULD supply a fully checksummed packet to the device.

If VIRTIO\_NET\_F\_HOST\_TSO4 is negotiated, the driver MAY set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV4 to request TCPv4 segmentation, otherwise the driver MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV4.

If VIRTIO\_NET\_F\_HOST\_TSO6 is negotiated, the driver MAY set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV6 to request TCPv6 segmentation, otherwise the driver MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV6.

If VIRTIO\_NET\_F\_HOST\_UFO is negotiated, the driver MAY set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_UDP to request UDP fragmentation, otherwise the driver MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_UDP.

If VIRTIO\_NET\_F\_HOST\_USO is negotiated, the driver MAY set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_UDP\_L4 to request UDP segmentation, otherwise the driver MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_UDP\_L4.

The driver SHOULD NOT send to the device TCP packets requiring segmentation offload which have the Explicit Congestion Notification bit set, unless the VIRTIO\_NET\_F\_HOST\_ECN feature is negotiated, in which case the driver MUST set the VIRTIO\_NET\_HDR\_GSO\_ECN bit in *gso\_type*.

If the VIRTIO\_NET\_F\_CSUM feature has been negotiated, the driver MAY set the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags*, if so:

1. the driver MUST validate the packet checksum at offset *csum\_offset* from *csum\_start* as well as all preceding offsets;
2. the driver MUST set the packet checksum stored in the buffer to the TCP/UDP pseudo header;
3. the driver MUST set *csum\_start* and *csum\_offset* such that calculating a ones' complement checksum from *csum\_start* up until the end of the packet and storing the result at offset *csum\_offset* from *csum\_start* will result in a fully checksummed packet;

If none of the VIRTIO\_NET\_F\_HOST\_TSO4, TSO6, USO or UFO options have been negotiated, the driver MUST set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_NONE.

If *gso\_type* differs from VIRTIO\_NET\_HDR\_GSO\_NONE, then the driver MUST also set the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags* and MUST set *gso\_size* to indicate the desired MSS.

If one of the VIRTIO\_NET\_F\_HOST\_TSO4, TSO6, USO or UFO options have been negotiated:

- If the VIRTIO\_NET\_F\_GUEST\_HDRLLEN feature has been negotiated, and *gso\_type* differs from VIRTIO\_NET\_HDR\_GSO\_NONE, the driver MUST set *hdr\_len* to a value equal to the length of the headers, including the transport header.
- If the VIRTIO\_NET\_F\_GUEST\_HDRLLEN feature has not been negotiated, or *gso\_type* is VIRTIO\_NET\_HDR\_GSO\_NONE, the driver SHOULD set *hdr\_len* to a value not less than the length of the headers, including the transport header.

The driver SHOULD accept the VIRTIO\_NET\_F\_GUEST\_HDRLLEN feature if it has been offered, and if it's able to provide the exact header length.

The driver MUST NOT set the VIRTIO\_NET\_HDR\_F\_DATA\_VALID and VIRTIO\_NET\_HDR\_F\_RSC\_INFO bits in *flags*.



### 5.1.6.2.2 Device Requirements: Packet Transmission

The device MUST ignore *flag* bits that it does not recognize.

If VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags* is not set, the device MUST NOT use the *csum\_start* and *csum\_offset*.

If one of the VIRTIO\_NET\_F\_HOST\_TSO4, TSO6, USO or UFO options have been negotiated:

- If the VIRTIO\_NET\_F\_GUEST\_HDRLEN feature has been negotiated, and *gso\_type* differs from VIRTIO\_NET\_HDR\_GSO\_NONE, the device MAY use *hdr\_len* as the transport header size.

**Note:** Caution should be taken by the implementation so as to prevent a malicious driver from attacking the device by setting an incorrect *hdr\_len*.

- If the VIRTIO\_NET\_F\_GUEST\_HDRLEN feature has not been negotiated, or *gso\_type* is VIRTIO\_NET\_HDR\_GSO\_NONE, the device MAY use *hdr\_len* only as a hint about the transport header size. The device MUST NOT rely on *hdr\_len* to be correct.

**Note:** This is due to various bugs in implementations.

If VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM is not set, the device MUST NOT rely on the packet checksum being correct.

### 5.1.6.2.3 Packet Transmission Interrupt

Often a driver will suppress transmission virtqueue interrupts and check for used packets in the transmit path of following packets.

The normal behavior in this interrupt handler is to retrieve used buffers from the virtqueue and free the corresponding headers and packets.

### 5.1.6.3 Setting Up Receive Buffers

It is generally a good idea to keep the receive virtqueue as fully populated as possible: if it runs out, network performance will suffer.

If the VIRTIO\_NET\_F\_GUEST\_TSO4, VIRTIO\_NET\_F\_GUEST\_TSO6 or VIRTIO\_NET\_F\_GUEST\_UFO features are used, the maximum incoming packet will be to 65550 bytes long (the maximum size of a TCP or UDP packet, plus the 14 byte ethernet header), otherwise 1514 bytes. The 12-byte struct *virtio\_net\_hdr* is prepended to this, making for 65562 or 1526 bytes.

#### 5.1.6.3.1 Driver Requirements: Setting Up Receive Buffers

- If VIRTIO\_NET\_F\_MRG\_RXBUF is not negotiated:
  - If VIRTIO\_NET\_F\_GUEST\_TSO4, VIRTIO\_NET\_F\_GUEST\_TSO6 or VIRTIO\_NET\_F\_GUEST\_UFO are negotiated, the driver SHOULD populate the receive queue(s) with buffers of at least 65562 bytes.
  - Otherwise, the driver SHOULD populate the receive queue(s) with buffers of at least 1526 bytes.
- If VIRTIO\_NET\_F\_MRG\_RXBUF is negotiated, each buffer MUST be at least the size of the struct *virtio\_net\_hdr*.

**Note:** Obviously each buffer can be split across multiple descriptor elements.

If VIRTIO\_NET\_F\_MQ is negotiated, each of *receiveq1*...*receiveqN* that will be used SHOULD be populated with receive buffers.

#### 5.1.6.3.2 Device Requirements: Setting Up Receive Buffers

The device MUST set *num\_buffers* to the number of descriptors used to hold the incoming packet.

The device MUST use only a single descriptor if VIRTIO\_NET\_F\_MRG\_RXBUF was not negotiated.



**Note:** This means that *num\_buffers* will always be 1 if VIRTIO\_NET\_F\_MRG\_RXBUF is not negotiated.

#### 5.1.6.4 Processing of Incoming Packets

When a packet is copied into a buffer in the receiveq, the optimal path is to disable further used buffer notifications for the receiveq and process packets until no more are found, then re-enable them.

Processing incoming packets involves:

1. *num\_buffers* indicates how many descriptors this packet is spread over (including this one): this will always be 1 if VIRTIO\_NET\_F\_MRG\_RXBUF was not negotiated. This allows receipt of large packets without having to allocate large buffers: a packet that does not fit in a single buffer can flow over to the next buffer, and so on. In this case, there will be at least *num\_buffers* used buffers in the virtqueue, and the device chains them together to form a single packet in a way similar to how it would store it in a single buffer spread over multiple descriptors. The other buffers will not begin with a struct *virtio\_net\_hdr*.
2. If *num\_buffers* is one, then the entire packet will be contained within this buffer, immediately following the struct *virtio\_net\_hdr*.
3. If the VIRTIO\_NET\_F\_GUEST\_CSUM feature was negotiated, the VIRTIO\_NET\_HDR\_F\_DATA\_VALID bit in *flags* can be set: if so, device has validated the packet checksum. In case of multiple encapsulated protocols, one level of checksums has been validated.

Additionally, VIRTIO\_NET\_F\_GUEST\_CSUM, TSO4, TSO6, UDP and ECN features enable receive checksum, large receive offload and ECN support which are the input equivalents of the transmit checksum, transmit segmentation offloading and ECN features, as described in 5.1.6.2:

1. If the VIRTIO\_NET\_F\_GUEST\_TSO4, TSO6 or UFO options were negotiated, then *gso\_type* MAY be something other than VIRTIO\_NET\_HDR\_GSO\_NONE, and *gso\_size* field indicates the desired MSS (see Packet Transmission point 2).
2. If the VIRTIO\_NET\_F\_RSC\_EXT option was negotiated (this implies one of VIRTIO\_NET\_F\_GUEST\_TSO4, TSO6), the device processes also duplicated ACK segments, reports number of coalesced TCP segments in *csum\_start* field and number of duplicated ACK segments in *csum\_offset* field and sets bit VIRTIO\_NET\_HDR\_F\_RSC\_INFO in *flags*.
3. If the VIRTIO\_NET\_F\_GUEST\_CSUM feature was negotiated, the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags* can be set: if so, the packet checksum at offset *csum\_offset* from *csum\_start* and any preceding checksums have been validated. The checksum on the packet is incomplete and if bit VIRTIO\_NET\_HDR\_F\_RSC\_INFO is not set in *flags*, then *csum\_start* and *csum\_offset* indicate how to calculate it (see Packet Transmission point 1).

If applicable, the device calculates per-packet hash for incoming packets as defined in 5.1.6.4.3.

If applicable, the device reports hash information for incoming packets as defined in 5.1.6.4.4.

##### 5.1.6.4.1 Device Requirements: Processing of Incoming Packets

If VIRTIO\_NET\_F\_MRG\_RXBUF has not been negotiated, the device MUST set *num\_buffers* to 1.

If VIRTIO\_NET\_F\_MRG\_RXBUF has been negotiated, the device MUST set *num\_buffers* to indicate the number of buffers the packet (including the header) is spread over.

If a receive packet is spread over multiple buffers, the device MUST use all buffers but the last (i.e. the first *num\_buffers* - 1 buffers) completely up to the full length of each buffer supplied by the driver.

The device MUST use all buffers used by a single receive packet together, such that at least *num\_buffers* are observed by driver as used.

If VIRTIO\_NET\_F\_GUEST\_CSUM is not negotiated, the device MUST set *flags* to zero and SHOULD supply a fully checksummed packet to the driver.

If VIRTIO\_NET\_F\_GUEST\_TSO4 is not negotiated, the device MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV4.

If VIRTIO\_NET\_F\_GUEST\_UDP is not negotiated, the device MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_UDP.

If VIRTIO\_NET\_F\_GUEST\_TSO6 is not negotiated, the device MUST NOT set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_TCPV6.

The device SHOULD NOT send to the driver TCP packets requiring segmentation offload which have the Explicit Congestion Notification bit set, unless the VIRTIO\_NET\_F\_GUEST\_ECN feature is negotiated, in which case the device MUST set the VIRTIO\_NET\_HDR\_GSO\_ECN bit in *gso\_type*.

If the VIRTIO\_NET\_F\_GUEST\_CSUM feature has been negotiated, the device MAY set the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags*, if so:

1. the device MUST validate the packet checksum at offset *csum\_offset* from *csum\_start* as well as all preceding offsets;
2. the device MUST set the packet checksum stored in the receive buffer to the TCP/UDP pseudo header;
3. the device MUST set *csum\_start* and *csum\_offset* such that calculating a ones' complement checksum from *csum\_start* up until the end of the packet and storing the result at offset *csum\_offset* from *csum\_start* will result in a fully checksummed packet;

If none of the VIRTIO\_NET\_F\_GUEST\_TSO4, TSO6 or UFO options have been negotiated, the device MUST set *gso\_type* to VIRTIO\_NET\_HDR\_GSO\_NONE.

If *gso\_type* differs from VIRTIO\_NET\_HDR\_GSO\_NONE, then the device MUST also set the VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags* MUST set *gso\_size* to indicate the desired MSS. If VIRTIO\_NET\_F\_RSC\_EXT was negotiated, the device MUST also set VIRTIO\_NET\_HDR\_F\_RSC\_INFO bit in *flags*, set *csum\_start* to number of coalesced TCP segments and set *csum\_offset* to number of received duplicated ACK segments.

If VIRTIO\_NET\_F\_RSC\_EXT was not negotiated, the device MUST not set VIRTIO\_NET\_HDR\_F\_RSC\_INFO bit in *flags*.

If one of the VIRTIO\_NET\_F\_GUEST\_TSO4, TSO6 or UFO options have been negotiated, the device SHOULD set *hdr\_len* to a value not less than the length of the headers, including the transport header.

If the VIRTIO\_NET\_F\_GUEST\_CSUM feature has been negotiated, the device MAY set the VIRTIO\_NET\_HDR\_F\_DATA\_VALID bit in *flags*, if so, the device MUST validate the packet checksum (in case of multiple encapsulated protocols, one level of checksums is validated).

#### 5.1.6.4.2 Driver Requirements: Processing of Incoming Packets

The driver MUST ignore *flag* bits that it does not recognize.

If VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM bit in *flags* is not set or if VIRTIO\_NET\_HDR\_F\_RSC\_INFO bit in *flags* is set, the driver MUST NOT use the *csum\_start* and *csum\_offset*.

If one of the VIRTIO\_NET\_F\_GUEST\_TSO4, TSO6 or UFO options have been negotiated, the driver MAY use *hdr\_len* only as a hint about the transport header size. The driver MUST NOT rely on *hdr\_len* to be correct.

**Note:** This is due to various bugs in implementations.

If neither VIRTIO\_NET\_HDR\_F\_NEEDS\_CSUM nor VIRTIO\_NET\_HDR\_F\_DATA\_VALID is set, the driver MUST NOT rely on the packet checksum being correct.

#### 5.1.6.4.3 Hash calculation for incoming packets

A device attempts to calculate a per-packet hash in the following cases:

- The feature VIRTIO\_NET\_F\_RSS was negotiated. The device uses the hash to determine the receive virtqueue to place incoming packets.

- The feature VIRTIO\_NET\_F\_HASH\_REPORT was negotiated. The device reports the hash value and the hash type with the packet.

If the feature VIRTIO\_NET\_F\_RSS was negotiated:

- The device uses *hash\_types* of the *virtio\_net\_rss\_config* structure as 'Enabled hash types' bitmask.
- The device uses a key as defined in *hash\_key\_data* and *hash\_key\_length* of the *virtio\_net\_rss\_config* structure (see 5.1.6.5.7.1).

If the feature VIRTIO\_NET\_F\_RSS was not negotiated:

- The device uses *hash\_types* of the *virtio\_net\_hash\_config* structure as 'Enabled hash types' bitmask.
- The device uses a key as defined in *hash\_key\_data* and *hash\_key\_length* of the *virtio\_net\_hash\_config* structure (see 5.1.6.5.6.4).

Note that if the device offers VIRTIO\_NET\_F\_HASH\_REPORT, even if it supports only one pair of virtqueues, it MUST support at least one of commands of VIRTIO\_NET\_CTRL\_MQ class to configure reported hash parameters:

- If the device offers VIRTIO\_NET\_F\_RSS, it MUST support VIRTIO\_NET\_CTRL\_MQ\_RSS\_CONFIG command per 5.1.6.5.7.1.
- Otherwise the device MUST support VIRTIO\_NET\_CTRL\_MQ\_HASH\_CONFIG command per 5.1.6.5.6.4.

#### 5.1.6.4.3.1 Supported/enabled hash types

Hash types applicable for IPv4 packets:

```
#define VIRTIO_NET_HASH_TYPE_IPV4      (1 << 0)
#define VIRTIO_NET_HASH_TYPE_TCPV4     (1 << 1)
#define VIRTIO_NET_HASH_TYPE_UDPV4     (1 << 2)
```

Hash types applicable for IPv6 packets without extension headers

```
#define VIRTIO_NET_HASH_TYPE_IPV6      (1 << 3)
#define VIRTIO_NET_HASH_TYPE_TCPV6     (1 << 4)
#define VIRTIO_NET_HASH_TYPE_UDPV6     (1 << 5)
```

Hash types applicable for IPv6 packets with extension headers

```
#define VIRTIO_NET_HASH_TYPE_IP_EX      (1 << 6)
#define VIRTIO_NET_HASH_TYPE_TCP_EX     (1 << 7)
#define VIRTIO_NET_HASH_TYPE_UDP_EX     (1 << 8)
```

#### 5.1.6.4.3.2 IPv4 packets

The device calculates the hash on IPv4 packets according to 'Enabled hash types' bitmask as follows:

- If VIRTIO\_NET\_HASH\_TYPE\_TCPV4 is set and the packet has a TCP header, the hash is calculated over the following fields:
  - Source IP address
  - Destination IP address
  - Source TCP port
  - Destination TCP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_UDPV4 is set and the packet has a UDP header, the hash is calculated over the following fields:
  - Source IP address
  - Destination IP address

- Source UDP port
- Destination UDP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_IPV4 is set, the hash is calculated over the following fields:
  - Source IP address
  - Destination IP address
- Else the device does not calculate the hash

#### 5.1.6.4.3.3 IPv6 packets without extension header

The device calculates the hash on IPv6 packets without extension headers according to 'Enabled hash types' bitmask as follows:

- If VIRTIO\_NET\_HASH\_TYPE\_TCPv6 is set and the packet has a TCPv6 header, the hash is calculated over the following fields:
  - Source IPv6 address
  - Destination IPv6 address
  - Source TCP port
  - Destination TCP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_UDPv6 is set and the packet has a UDPv6 header, the hash is calculated over the following fields:
  - Source IPv6 address
  - Destination IPv6 address
  - Source UDP port
  - Destination UDP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_IPV6 is set, the hash is calculated over the following fields:
  - Source IPv6 address
  - Destination IPv6 address
- Else the device does not calculate the hash

#### 5.1.6.4.3.4 IPv6 packets with extension header

The device calculates the hash on IPv6 packets with extension headers according to 'Enabled hash types' bitmask as follows:

- If VIRTIO\_NET\_HASH\_TYPE\_TCP\_EX is set and the packet has a TCPv6 header, the hash is calculated over the following fields:
  - Home address from the home address option in the IPv6 destination options header. If the extension header is not present, use the Source IPv6 address.
  - IPv6 address that is contained in the Routing-Header-Type-2 from the associated extension header. If the extension header is not present, use the Destination IPv6 address.
  - Source TCP port
  - Destination TCP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_UDP\_EX is set and the packet has a UDPv6 header, the hash is calculated over the following fields:
  - Home address from the home address option in the IPv6 destination options header. If the extension header is not present, use the Source IPv6 address.

- IPv6 address that is contained in the Routing-Header-Type-2 from the associated extension header. If the extension header is not present, use the Destination IPv6 address.
- Source UDP port
- Destination UDP port
- Else if VIRTIO\_NET\_HASH\_TYPE\_IP\_EX is set, the hash is calculated over the following fields:
  - Home address from the home address option in the IPv6 destination options header. If the extension header is not present, use the Source IPv6 address.
  - IPv6 address that is contained in the Routing-Header-Type-2 from the associated extension header. If the extension header is not present, use the Destination IPv6 address.
- Else skip IPv6 extension headers and calculate the hash as defined for an IPv6 packet without extension headers (see 5.1.6.4.3.3).

#### 5.1.6.4.4 Hash reporting for incoming packets

If VIRTIO\_NET\_F\_HASH\_REPORT was negotiated and the device has calculated the hash for the packet, the device fills *hash\_report* with the report type of calculated hash and *hash\_value* with the value of calculated hash.

If VIRTIO\_NET\_F\_HASH\_REPORT was negotiated but due to any reason the hash was not calculated, the device sets *hash\_report* to VIRTIO\_NET\_HASH\_REPORT\_NONE.

Possible values that the device can report in *hash\_report* are defined below. They correspond to supported hash types defined in 5.1.6.4.3.1 as follows:

VIRTIO\_NET\_HASH\_TYPE\_XXX = 1 « (VIRTIO\_NET\_HASH\_REPORT\_XXX - 1)

#define VIRTIO_NET_HASH_REPORT_NONE	0
#define VIRTIO_NET_HASH_REPORT_IPv4	1
#define VIRTIO_NET_HASH_REPORT_TCPv4	2
#define VIRTIO_NET_HASH_REPORT_UDPv4	3
#define VIRTIO_NET_HASH_REPORT_IPv6	4
#define VIRTIO_NET_HASH_REPORT_TCPv6	5
#define VIRTIO_NET_HASH_REPORT_UDPv6	6
#define VIRTIO_NET_HASH_REPORT_IPv6_EX	7
#define VIRTIO_NET_HASH_REPORT_TCPv6_EX	8
#define VIRTIO_NET_HASH_REPORT_UDPv6_EX	9

#### 5.1.6.5 Control Virtqueue

The driver uses the control virtqueue (if VIRTIO\_NET\_F\_CTRL\_VQ is negotiated) to send commands to manipulate various features of the device which would not easily map into the configuration space.

All commands are of the following form:

```
struct virtio_net_ctrl {
    u8 class;
    u8 command;
    u8 command-specific-data[];
    u8 ack;
};

/* ack values */
#define VIRTIO_NET_OK      0
#define VIRTIO_NET_ERR    1
```

The *class*, *command* and *command-specific-data* are set by the driver, and the device sets the *ack* byte. There is little it can do except issue a diagnostic if *ack* is not VIRTIO\_NET\_OK.

### 5.1.6.5.1 Packet Receive Filtering

If the VIRTIO\_NET\_F\_CTRL\_RX and VIRTIO\_NET\_F\_CTRL\_RX\_EXTRA features are negotiated, the driver can send control commands for promiscuous mode, multicast, unicast and broadcast receiving.

**Note:** In general, these commands are best-effort: unwanted packets could still arrive.

```
#define VIRTIO_NET_CTRL_RX      0
#define VIRTIO_NET_CTRL_RX_PROMISC    0
#define VIRTIO_NET_CTRL_RX_ALLMULTI  1
#define VIRTIO_NET_CTRL_RX_ALLUNI    2
#define VIRTIO_NET_CTRL_RX_NOMULTI   3
#define VIRTIO_NET_CTRL_RX_NOUNI     4
#define VIRTIO_NET_CTRL_RX_NOBCAST   5
```

#### 5.1.6.5.1.1 Device Requirements: Packet Receive Filtering

If the VIRTIO\_NET\_F\_CTRL\_RX feature has been negotiated, the device MUST support the following VIRTIO\_NET\_CTRL\_RX class commands:

- VIRTIO\_NET\_CTRL\_RX\_PROMISC turns promiscuous mode on and off. The command-specific-data is one byte containing 0 (off) or 1 (on). If promiscuous mode is on, the device SHOULD receive all incoming packets. This SHOULD take effect even if one of the other modes set by a VIRTIO\_NET\_CTRL\_RX class command is on.
- VIRTIO\_NET\_CTRL\_RX\_ALLMULTI turns all-multicast receive on and off. The command-specific-data is one byte containing 0 (off) or 1 (on). When all-multicast receive is on the device SHOULD allow all incoming multicast packets.

If the VIRTIO\_NET\_F\_CTRL\_RX\_EXTRA feature has been negotiated, the device MUST support the following VIRTIO\_NET\_CTRL\_RX class commands:

- VIRTIO\_NET\_CTRL\_RX\_ALLUNI turns all-unicast receive on and off. The command-specific-data is one byte containing 0 (off) or 1 (on). When all-unicast receive is on the device SHOULD allow all incoming unicast packets.
- VIRTIO\_NET\_CTRL\_RX\_NOMULTI suppresses multicast receive. The command-specific-data is one byte containing 0 (multicast receive allowed) or 1 (multicast receive suppressed). When multicast receive is suppressed, the device SHOULD NOT send multicast packets to the driver. This SHOULD take effect even if VIRTIO\_NET\_CTRL\_RX\_ALLMULTI is on. This filter SHOULD NOT apply to broadcast packets.
- VIRTIO\_NET\_CTRL\_RX\_NOUNI suppresses unicast receive. The command-specific-data is one byte containing 0 (unicast receive allowed) or 1 (unicast receive suppressed). When unicast receive is suppressed, the device SHOULD NOT send unicast packets to the driver. This SHOULD take effect even if VIRTIO\_NET\_CTRL\_RX\_ALLUNI is on.
- VIRTIO\_NET\_CTRL\_RX\_NOBCAST suppresses broadcast receive. The command-specific-data is one byte containing 0 (broadcast receive allowed) or 1 (broadcast receive suppressed). When broadcast receive is suppressed, the device SHOULD NOT send broadcast packets to the driver. This SHOULD take effect even if VIRTIO\_NET\_CTRL\_RX\_ALLMULTI is on.

#### 5.1.6.5.1.2 Driver Requirements: Packet Receive Filtering

If the VIRTIO\_NET\_F\_CTRL\_RX feature has not been negotiated, the driver MUST NOT issue commands VIRTIO\_NET\_CTRL\_RX\_PROMISC or VIRTIO\_NET\_CTRL\_RX\_ALLMULTI.

If the VIRTIO\_NET\_F\_CTRL\_RX\_EXTRA feature has not been negotiated, the driver MUST NOT issue commands VIRTIO\_NET\_CTRL\_RX\_ALLUNI, VIRTIO\_NET\_CTRL\_RX\_NOMULTI, VIRTIO\_NET\_CTRL\_RX\_NOUNI or VIRTIO\_NET\_CTRL\_RX\_NOBCAST.

### 5.1.6.5.2 Setting MAC Address Filtering

If the VIRTIO\_NET\_F\_CTRL\_RX feature is negotiated, the driver can send control commands for MAC address filtering.

```
struct virtio_net_ctrl_mac {
    le32 entries;
    u8 macs[entries][6];
};

#define VIRTIO_NET_CTRL_MAC      1
#define VIRTIO_NET_CTRL_MAC_TABLE_SET 0
#define VIRTIO_NET_CTRL_MAC_ADDR_SET 1
```

The device can filter incoming packets by any number of destination MAC addresses<sup>3</sup>. This table is set using the class VIRTIO\_NET\_CTRL\_MAC and the command VIRTIO\_NET\_CTRL\_MAC\_TABLE\_SET. The command-specific-data is two variable length tables of 6-byte MAC addresses (as described in struct virtio\_net\_ctrl\_mac). The first table contains unicast addresses, and the second contains multicast addresses.

The VIRTIO\_NET\_CTRL\_MAC\_ADDR\_SET command is used to set the default MAC address which rx filtering accepts (and if VIRTIO\_NET\_F\_MAC has been negotiated, this will be reflected in *mac* in config space).

The command-specific-data for VIRTIO\_NET\_CTRL\_MAC\_ADDR\_SET is the 6-byte MAC address.

#### 5.1.6.5.2.1 Device Requirements: Setting MAC Address Filtering

The device MUST have an empty MAC filtering table on reset.

The device MUST update the MAC filtering table before it consumes the VIRTIO\_NET\_CTRL\_MAC\_TABLE\_SET command.

The device MUST update *mac* in config space before it consumes the VIRTIO\_NET\_CTRL\_MAC\_ADDR\_SET command, if VIRTIO\_NET\_F\_MAC has been negotiated.

The device SHOULD drop incoming packets which have a destination MAC which matches neither the *mac* (or that set with VIRTIO\_NET\_CTRL\_MAC\_ADDR\_SET) nor the MAC filtering table.

#### 5.1.6.5.2.2 Driver Requirements: Setting MAC Address Filtering

If VIRTIO\_NET\_F\_CTRL\_RX has not been negotiated, the driver MUST NOT issue VIRTIO\_NET\_CTRL\_MAC class commands.

If VIRTIO\_NET\_F\_CTRL\_RX has been negotiated, the driver SHOULD issue VIRTIO\_NET\_CTRL\_MAC\_ADDR\_SET to set the default mac if it is different from *mac*.

The driver MUST follow the VIRTIO\_NET\_CTRL\_MAC\_TABLE\_SET command by a le32 number, followed by that number of non-multicast MAC addresses, followed by another le32 number, followed by that number of multicast addresses. Either number MAY be 0.

#### 5.1.6.5.2.3 Legacy Interface: Setting MAC Address Filtering

When using the legacy interface, transitional devices and drivers MUST format *entries* in struct virtio\_net\_ctrl\_mac according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

Legacy drivers that didn't negotiate VIRTIO\_NET\_F\_CTRL\_MAC\_ADDR changed *mac* in config space when NIC is accepting incoming packets. These drivers always wrote the mac value from first to last byte, therefore after detecting such drivers, a transitional device MAY defer MAC update, or MAY defer processing incoming packets until driver writes the last byte of *mac* in the config space.

<sup>3</sup>Since there are no guarantees, it can use a hash filter or silently switch to allmulti or promiscuous mode if it is given too many addresses.



### 5.1.6.5.3 VLAN Filtering

If the driver negotiates the `VIRTIO_NET_F_CTRL_VLAN` feature, it can control a VLAN filter table in the device.

**Note:** Similar to the MAC address based filtering, the VLAN filtering is also best-effort: unwanted packets could still arrive.

<code>#define VIRTIO_NET_CTRL_VLAN</code>	<code>2</code>	
<code>#define VIRTIO_NET_CTRL_VLAN_ADD</code>		<code>0</code>
<code>#define VIRTIO_NET_CTRL_VLAN_DEL</code>		<code>1</code>

Both the `VIRTIO_NET_CTRL_VLAN_ADD` and `VIRTIO_NET_CTRL_VLAN_DEL` command take a little-endian 16-bit VLAN id as the command-specific-data.

#### 5.1.6.5.3.1 Legacy Interface: VLAN Filtering

When using the legacy interface, transitional devices and drivers **MUST** format the VLAN id according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

### 5.1.6.5.4 Gratuitous Packet Sending

If the driver negotiates the `VIRTIO_NET_F_GUEST_ANNOUNCE` (depends on `VIRTIO_NET_F_CTRL_VQ`), the device can ask the driver to send gratuitous packets; this is usually done after the guest has been physically migrated, and needs to announce its presence on the new network links. (As hypervisor does not have the knowledge of guest network configuration (eg. tagged vlan) it is simplest to prod the guest in this way).

<code>#define VIRTIO_NET_CTRL_ANNOUNCE</code>	<code>3</code>	
<code>#define VIRTIO_NET_CTRL_ANNOUNCE_ACK</code>		<code>0</code>

The driver checks `VIRTIO_NET_S_ANNOUNCE` bit in the device configuration *status* field when it notices the changes of device configuration. The command `VIRTIO_NET_CTRL_ANNOUNCE_ACK` is used to indicate that driver has received the notification and device clears the `VIRTIO_NET_S_ANNOUNCE` bit in *status*.

Processing this notification involves:

1. Sending the gratuitous packets (eg. ARP) or marking there are pending gratuitous packets to be sent and letting deferred routine to send them.
2. Sending `VIRTIO_NET_CTRL_ANNOUNCE_ACK` command through control vq.

#### 5.1.6.5.4.1 Driver Requirements: Gratuitous Packet Sending

If the driver negotiates `VIRTIO_NET_F_GUEST_ANNOUNCE`, it **SHOULD** notify network peers of its new location after it sees the `VIRTIO_NET_S_ANNOUNCE` bit in *status*. The driver **MUST** send a command on the command queue with class `VIRTIO_NET_CTRL_ANNOUNCE` and command `VIRTIO_NET_CTRL_ANNOUNCE_ACK`.

#### 5.1.6.5.4.2 Device Requirements: Gratuitous Packet Sending

If `VIRTIO_NET_F_GUEST_ANNOUNCE` is negotiated, the device **MUST** clear the `VIRTIO_NET_S_ANNOUNCE` bit in *status* upon receipt of a command buffer with class `VIRTIO_NET_CTRL_ANNOUNCE` and command `VIRTIO_NET_CTRL_ANNOUNCE_ACK` before marking the buffer as used.

### 5.1.6.5.5 Device operation in multiqueue mode

This specification defines the following modes that a device **MAY** implement for operation with multiple transmit/receive virtqueues:

- Automatic receive steering as defined in 5.1.6.5.6. If a device supports this mode, it offers the VIRTIO\_ - NET\_F\_MQ feature bit.
- Receive-side scaling as defined in 5.1.6.5.7.3. If a device supports this mode, it offers the VIRTIO\_ - NET\_F\_RSS feature bit.

A device MAY support one of these features or both. The driver MAY negotiate any set of these features that the device supports.

Multiqueue is disabled by default.

The driver enables multiqueue by sending a command using *class* VIRTIO\_NET\_CTRL\_MQ. The *command* selects the mode of multiqueue operation, as follows:

```
#define VIRTIO_NET_CTRL_MQ      4
#define VIRTIO_NET_CTRL_MQ_VQ_PAIRS_SET    0 (for automatic receive steering)
#define VIRTIO_NET_CTRL_MQ_RSS_CONFIG      1 (for configurable receive steering)
#define VIRTIO_NET_CTRL_MQ_HASH_CONFIG     2 (for configurable hash calculation)
```

If more than one multiqueue mode is negotiated, the resulting device configuration is defined by the last command sent by the driver.

#### 5.1.6.5.6 Automatic receive steering in multiqueue mode

If the driver negotiates the VIRTIO\_NET\_F\_MQ feature bit (depends on VIRTIO\_NET\_F\_CTRL\_VQ), it MAY transmit outgoing packets on one of the multiple transmitq1...transmitqN and ask the device to queue incoming packets into one of the multiple receiveq1...receiveqN depending on the packet flow.

The driver enables multiqueue by sending the VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command, specifying the number of the transmit and receive queues to be used up to *max\_virtqueue\_pairs*; subsequently, transmitq1...transmitqn and receiveq1...receiveqn where *n=virtqueue\_pairs* MAY be used.

```
struct virtio_net_ctrl_mq_pairs_set {
    le16 virtqueue_pairs;
};
#define VIRTIO_NET_CTRL_MQ_VQ_PAIRS_MIN      1
#define VIRTIO_NET_CTRL_MQ_VQ_PAIRS_MAX     0x8000
```

When multiqueue is enabled by VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command, the device MUST use automatic receive steering based on packet flow. Programming of the receive steering classifier is implicit. After the driver transmitted a packet of a flow on transmitqX, the device SHOULD cause incoming packets for that flow to be steered to receiveqX. For uni-directional protocols, or where no packets have been transmitted yet, the device MAY steer a packet to a random queue out of the specified receiveq1...receiveqn.

Multiqueue is disabled by VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET with *virtqueue\_pairs* to 1 (this is the default) and waiting for the device to use the command buffer.

##### 5.1.6.5.6.1 Driver Requirements: Automatic receive steering in multiqueue mode

The driver MUST configure the virtqueues before enabling them with the VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command.

The driver MUST NOT request a *virtqueue\_pairs* of 0 or greater than *max\_virtqueue\_pairs* in the device configuration space.

The driver MUST queue packets only on any transmitq1 before the VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command.

The driver MUST NOT queue packets on transmit queues greater than *virtqueue\_pairs* once it has placed the VIRTIO\_NET\_CTRL\_MQ\_VQ\_PAIRS\_SET command in the available ring.

#### 5.1.6.5.6.2 Device Requirements: Automatic receive steering in multiqueue mode

After initialization of reset, the device MUST queue packets only on `receiveq1`.

The device MUST NOT queue packets on receive queues greater than `virtqueue_pairs` once it has placed the `VIRTIO_NET_CTRL_MQ_VQ_PAIRS_SET` command in a used buffer.

#### 5.1.6.5.6.3 Legacy Interface: Automatic receive steering in multiqueue mode

When using the legacy interface, transitional devices and drivers MUST format `virtqueue_pairs` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.1.6.5.6.4 Hash calculation

If `VIRTIO_NET_F_HASH_REPORT` was negotiated and the device uses automatic receive steering, the device MUST support a command to configure hash calculation parameters.

The driver provides parameters for hash calculation as follows:

`class VIRTIO_NET_CTRL_MQ`, `command VIRTIO_NET_CTRL_MQ_HASH_CONFIG`.

The *command-specific-data* has following format:

```
struct virtio_net_hash_config {
    le32 hash_types;
    le16 reserved[4];
    u8 hash_key_length;
    u8 hash_key_data[hash_key_length];
};
```

Field `hash_types` contains a bitmask of allowed hash types as defined in 5.1.6.4.3.1. Initially the device has all hash types disabled and reports only `VIRTIO_NET_HASH_REPORT_NONE`.

Field `reserved` MUST contain zeroes. It is defined to make the structure to match the layout of `virtio_net_rss_config` structure, defined in 5.1.6.5.7.

Fields `hash_key_length` and `hash_key_data` define the key to be used in hash calculation.

#### 5.1.6.5.7 Receive-side scaling (RSS)

A device offers the feature `VIRTIO_NET_F_RSS` if it supports RSS receive steering with Toeplitz hash calculation and configurable parameters.

A driver queries RSS capabilities of the device by reading device configuration as defined in 5.1.4

##### 5.1.6.5.7.1 Setting RSS parameters

Driver sends a `VIRTIO_NET_CTRL_MQ_RSS_CONFIG` command using the following format for *command-specific-data*:

```
struct virtio_net_rss_config {
    le32 hash_types;
    le16 indirection_table_mask;
    le16 unclassified_queue;
    le16 indirection_table[indirection_table_length];
    le16 max_tx_vq;
    u8 hash_key_length;
    u8 hash_key_data[hash_key_length];
};
```

Field `hash_types` contains a bitmask of allowed hash types as defined in 5.1.6.4.3.1.

Field `indirection_table_mask` is a mask to be applied to the calculated hash to produce an index in the `indirection_table` array. Number of entries in `indirection_table` is  $(indirection\_table\_mask + 1)$ .

Field *unclassified\_queue* contains the 0-based index of the receive virtqueue to place unclassified packets in. Index 0 corresponds to receiveq1.

Field *indirection\_table* contains an array of 0-based indices of receive virtqueues. Index 0 corresponds to receiveq1.

A driver sets *max\_tx\_vq* to inform a device how many transmit virtqueues it may use (transmitq1...transmitq *max\_tx\_vq*).

Fields *hash\_key\_length* and *hash\_key\_data* define the key to be used in hash calculation.

#### 5.1.6.5.7.2 Driver Requirements: Setting RSS parameters

A driver MUST NOT send the VIRTIO\_NET\_CTRL\_MQ\_RSS\_CONFIG command if the feature VIRTIO\_NET\_F\_RSS has not been negotiated.

A driver MUST fill the *indirection\_table* array only with indices of enabled queues. Index 0 corresponds to receiveq1.

The number of entries in *indirection\_table* (*indirection\_table\_mask* + 1) MUST be a power of two.

A driver MUST use *indirection\_table\_mask* values that are less than *rss\_max\_indirection\_table\_length* reported by a device.

A driver MUST NOT set any VIRTIO\_NET\_HASH\_TYPE\_ flags that are not supported by a device.

#### 5.1.6.5.7.3 Device Requirements: RSS processing

The device MUST determine the destination queue for a network packet as follows:

- Calculate the hash of the packet as defined in 5.1.6.4.3.
- If the device did not calculate the hash for the specific packet, the device directs the packet to the receiveq specified by *unclassified\_queue* of *virtio\_net\_rss\_config* structure (value of 0 corresponds to receiveq1).
- Apply *indirection\_table\_mask* to the calculated hash and use the result as the index in the *indirection\_table* to get 0-based number of destination receiveq (value of 0 corresponds to receiveq1).

#### 5.1.6.5.8 Offloads State Configuration

If the VIRTIO\_NET\_F\_CTRL\_GUEST\_OFFLOADS feature is negotiated, the driver can send control commands for dynamic offloads state configuration.

##### 5.1.6.5.8.1 Setting Offloads State

To configure the offloads, the following layout structure and definitions are used:

```
le64 offloads;

#define VIRTIO_NET_F_GUEST_CSUM      1
#define VIRTIO_NET_F_GUEST_TS04     7
#define VIRTIO_NET_F_GUEST_TS06     8
#define VIRTIO_NET_F_GUEST_ECN      9
#define VIRTIO_NET_F_GUEST_UFO     10

#define VIRTIO_NET_CTRL_GUEST_OFFLOADS 5
#define VIRTIO_NET_CTRL_GUEST_OFFLOADS_SET 0
```

The class VIRTIO\_NET\_CTRL\_GUEST\_OFFLOADS has one command: VIRTIO\_NET\_CTRL\_GUEST\_OFFLOADS\_SET applies the new offloads configuration.

le64 value passed as command data is a bitmask, bits set define offloads to be enabled, bits cleared - offloads to be disabled.

There is a corresponding device feature for each offload. Upon feature negotiation corresponding offload gets enabled to preserve backward compatibility.

#### 5.1.6.5.8.2 Driver Requirements: Setting Offloads State

A driver MUST NOT enable an offload for which the appropriate feature has not been negotiated.

#### 5.1.6.5.8.3 Legacy Interface: Setting Offloads State

When using the legacy interface, transitional devices and drivers MUST format *offloads* according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.1.6.6 Legacy Interface: Framing Requirements

When using legacy interfaces, transitional drivers which have not negotiated `VIRTIO_F_ANY_LAYOUT` MUST use a single descriptor for the struct `virtio_net_hdr` on both transmit and receive, with the network data in the following descriptors.

Additionally, when using the control virtqueue (see 5.1.6.5) , transitional drivers which have not negotiated `VIRTIO_F_ANY_LAYOUT` MUST:

- for all commands, use a single 2-byte descriptor including the first two fields: *class* and *command*
- for all commands except `VIRTIO_NET_CTRL_MAC_TABLE_SET` use a single descriptor including command-specific-data with no padding.
- for the `VIRTIO_NET_CTRL_MAC_TABLE_SET` command use exactly two descriptors including command-specific-data with no padding: the first of these descriptors MUST include the `virtio_net_ctrl_mac` table structure for the unicast addresses with no padding, the second of these descriptors MUST include the `virtio_net_ctrl_mac` table structure for the multicast addresses with no padding.
- for all commands, use a single 1-byte descriptor for the *ack* field

See 2.7.4.

## 5.2 Block Device

The virtio block device is a simple virtual block device (ie. disk). Read and write requests (and other exotic requests) are placed in one of its queues, and serviced (probably out of order) by the device except where noted.

### 5.2.1 Device ID

2

### 5.2.2 Virtqueues

0 requestq1

...

N-1 requestqN

N=1 if `VIRTIO_BLK_F_MQ` is not negotiated, otherwise N is set by *num\_queues*.

### 5.2.3 Feature bits

**VIRTIO\_BLK\_F\_SIZE\_MAX (1)** Maximum size of any single segment is in *size\_max*.

**VIRTIO\_BLK\_F\_SEG\_MAX (2)** Maximum number of segments in a request is in *seg\_max*.

**VIRTIO\_BLK\_F\_GEOMETRY (4)** Disk-style geometry specified in *geometry*.

**VIRTIO\_BLK\_F\_RO (5)** Device is read-only.

**VIRTIO\_BLK\_F\_BLK\_SIZE (6)** Block size of disk is in *blk\_size*.

**VIRTIO\_BLK\_F\_FLUSH (9)** Cache flush command support.

**VIRTIO\_BLK\_F\_TOPOLOGY (10)** Device exports information on optimal I/O alignment.

**VIRTIO\_BLK\_F\_CONFIG\_WCE (11)** Device can toggle its cache between writeback and writethrough modes.

**VIRTIO\_BLK\_F\_MQ (12)** Device supports multiqueue.

**VIRTIO\_BLK\_F\_DISCARD (13)** Device can support discard command, maximum discard sectors size in *max\_discard\_sectors* and maximum discard segment number in *max\_discard\_seg*.

**VIRTIO\_BLK\_F\_WRITE\_ZEROES (14)** Device can support write zeroes command, maximum write zeroes sectors size in *max\_write\_zeroes\_sectors* and maximum write zeroes segment number in *max\_write\_zeroes\_seg*.

**VIRTIO\_BLK\_F\_LIFETIME (15)** Device supports providing storage lifetime information.

**VIRTIO\_BLK\_F\_SECURE\_ERASE (16)** Device supports secure erase command, maximum erase sectors count in *max\_secure\_erase\_sectors* and maximum erase segment number in *max\_secure\_erase\_seg*.

### 5.2.3.1 Legacy Interface: Feature bits

**VIRTIO\_BLK\_F\_BARRIER (0)** Device supports request barriers.

**VIRTIO\_BLK\_F\_SCSI (7)** Device supports scsi packet commands.

**Note:** In the legacy interface, VIRTIO\_BLK\_F\_FLUSH was also called VIRTIO\_BLK\_F\_WCE.

## 5.2.4 Device configuration layout

The *capacity* of the device (expressed in 512-byte sectors) is always present. The availability of the others all depend on various feature bits as indicated above.

The field *num\_queues* only exists if VIRTIO\_BLK\_F\_MQ is set. This field specifies the number of queues.

The parameters in the configuration space of the device *max\_discard\_sectors* *discard\_sector\_alignment* are expressed in 512-byte units if the VIRTIO\_BLK\_F\_DISCARD feature bit is negotiated. The *max\_write\_zeroes\_sectors* is expressed in 512-byte units if the VIRTIO\_BLK\_F\_WRITE\_ZEROES feature bit is negotiated. The parameters in the configuration space of the device *max\_secure\_erase\_sectors* *secure\_erase\_sector\_alignment* are expressed in 512-byte units if the VIRTIO\_BLK\_F\_SECURE\_ERASE feature bit is negotiated.

```
struct virtio_blk_config {
    le64 capacity;
    le32 size_max;
    le32 seg_max;
    struct virtio_blk_geometry {
        le16 cylinders;
        u8 heads;
        u8 sectors;
    } geometry;
    le32 blk_size;
    struct virtio_blk_topology {
        // # of logical blocks per physical block (log2)
        u8 physical_block_exp;
        // offset of first aligned logical block
        u8 alignment_offset;
        // suggested minimum I/O size in blocks
        le16 min_io_size;
        // optimal (suggested maximum) I/O size in blocks
        le32 opt_io_size;
    } topology;
    u8 writeback;
    u8 unused0;
    u16 num_queues;
}
```

```

le32 max_discard_sectors;
le32 max_discard_seg;
le32 discard_sector_alignment;
le32 max_write_zeroes_sectors;
le32 max_write_zeroes_seg;
u8 write_zeroes_may_unmap;
u8 unused1[3];
le32 max_secure_erase_sectors;
le32 max_secure_erase_seg;
le32 secure_erase_sector_alignment;
};

```

#### 5.2.4.1 Legacy Interface: Device configuration layout

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_blk_config` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

### 5.2.5 Device Initialization

1. The device size can be read from *capacity*.
2. If the `VIRTIO_BLK_F_BLK_SIZE` feature is negotiated, *blk\_size* can be read to determine the optimal sector size for the driver to use. This does not affect the units used in the protocol (always 512 bytes), but awareness of the correct value can affect performance.
3. If the `VIRTIO_BLK_F_RO` feature is set by the device, any write requests will fail.
4. If the `VIRTIO_BLK_F_TOPOLOGY` feature is negotiated, the fields in the *topology* struct can be read to determine the physical block size and optimal I/O lengths for the driver to use. This also does not affect the units in the protocol, only performance.
5. If the `VIRTIO_BLK_F_CONFIG_WCE` feature is negotiated, the cache mode can be read or set through the *writeback* field. 0 corresponds to a writethrough cache, 1 to a writeback cache<sup>4</sup>. The cache mode after reset can be either writeback or writethrough. The actual mode can be determined by reading *writeback* after feature negotiation.
6. If the `VIRTIO_BLK_F_DISCARD` feature is negotiated, *max\_discard\_sectors* and *max\_discard\_seg* can be read to determine the maximum discard sectors and maximum number of discard segments for the block driver to use. *discard\_sector\_alignment* can be used by OS when splitting a request based on alignment.
7. If the `VIRTIO_BLK_F_WRITE_ZEROES` feature is negotiated, *max\_write\_zeroes\_sectors* and *max\_write\_zeroes\_seg* can be read to determine the maximum write zeroes sectors and maximum number of write zeroes segments for the block driver to use.
8. If the `VIRTIO_BLK_F_MQ` feature is negotiated, *num\_queues* field can be read to determine the number of queues.
9. If the `VIRTIO_BLK_F_SECURE_ERASE` feature is negotiated, *max\_secure\_erase\_sectors* and *max\_secure\_erase\_seg* can be read to determine the maximum secure erase sectors and maximum number of secure erase segments for the block driver to use. *secure\_erase\_sector\_alignment* can be used by OS when splitting a request based on alignment.

#### 5.2.5.1 Driver Requirements: Device Initialization

Drivers SHOULD NOT negotiate `VIRTIO_BLK_F_FLUSH` if they are incapable of sending `VIRTIO_BLK_T_FLUSH` commands.

<sup>4</sup>Consistent with 5.2.6.2, a writethrough cache can be defined broadly as a cache that commits writes to persistent device backend storage before reporting their completion. For example, a battery-backed writeback cache actually counts as writethrough according to this definition.



If neither `VIRTIO_BLK_F_CONFIG_WCE` nor `VIRTIO_BLK_F_FLUSH` are negotiated, the driver MAY deduce the presence of a writethrough cache. If `VIRTIO_BLK_F_CONFIG_WCE` was not negotiated but `VIRTIO_BLK_F_FLUSH` was, the driver SHOULD assume presence of a writeback cache.

The driver MUST NOT read *writeback* before setting the `FEATURES_OK` device status bit.

### 5.2.5.2 Device Requirements: Device Initialization

Devices SHOULD always offer `VIRTIO_BLK_F_FLUSH`, and MUST offer it if they offer `VIRTIO_BLK_F_CONFIG_WCE`.

If `VIRTIO_BLK_F_CONFIG_WCE` is negotiated but `VIRTIO_BLK_F_FLUSH` is not, the device MUST initialize *writeback* to 0.

The device MUST initialize padding bytes *unused0* and *unused1* to 0.

### 5.2.5.3 Legacy Interface: Device Initialization

Because legacy devices do not have `FEATURES_OK`, transitional devices MUST implement slightly different behavior around feature negotiation when used through the legacy interface. In particular, when using the legacy interface:

- the driver MAY read or write *writeback* before setting the `DRIVER` or `DRIVER_OK` device status bit
- the device MUST NOT modify the cache mode (and *writeback*) as a result of a driver setting a status bit, unless the `DRIVER_OK` bit is being set and the driver has not set the `VIRTIO_BLK_F_CONFIG_WCE` driver feature bit.
- the device MUST NOT modify the cache mode (and *writeback*) as a result of a driver modifying the driver feature bits, for example if the driver sets the `VIRTIO_BLK_F_CONFIG_WCE` driver feature bit but does not set the `VIRTIO_BLK_F_FLUSH` bit.

## 5.2.6 Device Operation

The driver queues requests to the virtqueues, and they are used by the device (not necessarily in order). Each request is of form:

```
struct virtio_blk_req {
    le32 type;
    le32 reserved;
    le64 sector;
    u8 data[];
    u8 status;
};
```

The type of the request is either a read (`VIRTIO_BLK_T_IN`), a write (`VIRTIO_BLK_T_OUT`), a discard (`VIRTIO_BLK_T_DISCARD`), a write zeroes (`VIRTIO_BLK_T_WRITE_ZEROES`), a flush (`VIRTIO_BLK_T_FLUSH`), a get device ID string command (`VIRTIO_BLK_T_GET_ID`), a secure erase (`VIRTIO_BLK_T_SECURE_ERASE`), or a get device lifetime command (`VIRTIO_BLK_T_GET_LIFETIME`).

```
#define VIRTIO_BLK_T_IN          0
#define VIRTIO_BLK_T_OUT        1
#define VIRTIO_BLK_T_FLUSH      4
#define VIRTIO_BLK_T_GET_ID     8
#define VIRTIO_BLK_T_GET_LIFETIME 10
#define VIRTIO_BLK_T_DISCARD    11
#define VIRTIO_BLK_T_WRITE_ZEROES 13
#define VIRTIO_BLK_T_SECURE_ERASE 14
```

The *sector* number indicates the offset (multiplied by 512) where the read or write is to occur. This field is unused and set to 0 for commands other than read or write.

`VIRTIO_BLK_T_IN` requests populate *data* with the contents of sectors read from the block device (in multiples of 512 bytes). `VIRTIO_BLK_T_OUT` requests write the contents of *data* to the block device (in multiples of 512 bytes).



The *data* used for discard, secure erase or write zeroes commands consists of one or more segments. The maximum number of segments is *max\_discard\_seg* for discard commands, *max\_secure\_erase\_seg* for secure erase commands and *max\_write\_zeroes\_seg* for write zeroes commands. Each segment is of form:

```
struct virtio_blk_discard_write_zeroes {
    le64 sector;
    le32 num_sectors;
    struct {
        le32 unmap:1;
        le32 reserved:31;
    } flags;
};
```

*sector* indicates the starting offset (in 512-byte units) of the segment, while *num\_sectors* indicates the number of sectors in each discarded range. *unmap* is only used in write zeroes commands and allows the device to discard the specified range, provided that following reads return zeroes.

VIRTIO\_BLK\_T\_GET\_ID requests fetch the device ID string from the device into *data*. The device ID string is a NUL-padded ASCII string up to 20 bytes long. If the string is 20 bytes long then there is no NUL terminator.

The *data* used for VIRTIO\_BLK\_T\_GET\_LIFETIME requests is populated by the device, and is of the form

```
struct virtio_blk_lifetime {
    le16 pre_eol_info;
    le16 device_lifetime_est_typ_a;
    le16 device_lifetime_est_typ_b;
};
```

The *pre\_eol\_info* specifies the percentage of reserved blocks that are consumed and will have one of these values:

```
/* Value not available */
#define VIRTIO_BLK_PRE_EOL_INFO_UNDEFINED    0
/* < 80% of reserved blocks are consumed */
#define VIRTIO_BLK_PRE_EOL_INFO_NORMAL      1
/* 80% of reserved blocks are consumed */
#define VIRTIO_BLK_PRE_EOL_INFO_WARNING     2
/* 90% of reserved blocks are consumed */
#define VIRTIO_BLK_PRE_EOL_INFO_URGENT      3
/* All others values are reserved */
```

The *device\_lifetime\_est\_typ\_a* refers to wear of SLC cells and is provided in increments of 10 used, and so on, thru to 11 meaning estimated lifetime exceeded. All values above 11 are reserved.

The *device\_lifetime\_est\_typ\_b* refers to wear of MLC cells and is provided with the same semantics as *device\_lifetime\_est\_typ\_a*.

The final *status* byte is written by the device: either VIRTIO\_BLK\_S\_OK for success, VIRTIO\_BLK\_S\_IOERR for device or driver error or VIRTIO\_BLK\_S\_UNSUPP for a request unsupported by device:

```
#define VIRTIO_BLK_S_OK          0
#define VIRTIO_BLK_S_IOERR      1
#define VIRTIO_BLK_S_UNSUPP     2
```

The status of individual segments is indeterminate when a discard or write zero command produces VIRTIO\_BLK\_S\_IOERR. A segment may have completed successfully, failed, or not been processed by the device.

### 5.2.6.1 Driver Requirements: Device Operation

A driver **MUST NOT** submit a request which would cause a read or write beyond *capacity*.

A driver **SHOULD** accept the VIRTIO\_BLK\_F\_RO feature if offered.

A driver **MUST** set *sector* to 0 for a VIRTIO\_BLK\_T\_FLUSH request. A driver **SHOULD NOT** include any data in a VIRTIO\_BLK\_T\_FLUSH request.

The length of *data* **MUST** be a multiple of 512 bytes for VIRTIO\_BLK\_T\_IN and VIRTIO\_BLK\_T\_OUT requests.

The length of *data* **MUST** be a multiple of the size of struct `virtio_blk_discard_write_zeroes` for VIRTIO\_BLK\_T\_DISCARD, VIRTIO\_BLK\_T\_SECURE\_ERASE and VIRTIO\_BLK\_T\_WRITE\_ZEROES requests.

The length of *data* **MUST** be 20 bytes for VIRTIO\_BLK\_T\_GET\_ID requests.

VIRTIO\_BLK\_T\_DISCARD requests **MUST NOT** contain more than *max\_discard\_seg* struct `virtio_blk_discard_write_zeroes` segments in *data*.

VIRTIO\_BLK\_T\_SECURE\_ERASE requests **MUST NOT** contain more than *max\_secure\_erase\_seg* struct `virtio_blk_discard_write_zeroes` segments in *data*.

VIRTIO\_BLK\_T\_WRITE\_ZEROES requests **MUST NOT** contain more than *max\_write\_zeroes\_seg* struct `virtio_blk_discard_write_zeroes` segments in *data*.

If the VIRTIO\_BLK\_F\_CONFIG\_WCE feature is negotiated, the driver **MAY** switch to writethrough or write-back mode by writing respectively 0 and 1 to the *writeback* field. After writing a 0 to *writeback*, the driver **MUST NOT** assume that any volatile writes have been committed to persistent device backend storage.

The *unmap* bit **MUST** be zero for discard commands. The driver **MUST NOT** assume anything about the data returned by read requests after a range of sectors has been discarded.

A driver **MUST NOT** assume that individual segments in a multi-segment VIRTIO\_BLK\_T\_DISCARD or VIRTIO\_BLK\_T\_WRITE\_ZEROES request completed successfully, failed, or were processed by the device at all if the request failed with VIRTIO\_BLK\_S\_IOERR.

### 5.2.6.2 Device Requirements: Device Operation

A device **MUST** set the *status* byte to VIRTIO\_BLK\_S\_IOERR for a write request if the VIRTIO\_BLK\_F\_RO feature is offered, and **MUST NOT** write any data.

The device **MUST** set the *status* byte to VIRTIO\_BLK\_S\_UNSUPP for discard, secure erase and write zeroes commands if any unknown flag is set. Furthermore, the device **MUST** set the *status* byte to VIRTIO\_BLK\_S\_UNSUPP for discard commands if the *unmap* flag is set.

For discard commands, the device **MAY** deallocate the specified range of sectors in the device backend storage.

For write zeroes commands, if the *unmap* is set, the device **MAY** deallocate the specified range of sectors in the device backend storage, as if the discard command had been sent. After a write zeroes command is completed, reads of the specified ranges of sectors **MUST** return zeroes. This is true independent of whether *unmap* was set or clear.

The device **SHOULD** clear the *write\_zeroes\_may\_unmap* field of the virtio configuration space if and only if a write zeroes request cannot result in deallocating one or more sectors. The device **MAY** change the content of the field during operation of the device; when this happens, the device **SHOULD** trigger a configuration change notification.

A write is considered volatile when it is submitted; the contents of sectors covered by a volatile write are undefined in persistent device backend storage until the write becomes stable. A write becomes stable once it is completed and one or more of the following conditions is true:

1. neither VIRTIO\_BLK\_F\_CONFIG\_WCE nor VIRTIO\_BLK\_F\_FLUSH feature were negotiated, but VIRTIO\_BLK\_F\_FLUSH was offered by the device;
2. the VIRTIO\_BLK\_F\_CONFIG\_WCE feature was negotiated and the *writeback* field in configuration space was 0 **all the time between the submission of the write and its completion**;
3. a VIRTIO\_BLK\_T\_FLUSH request is sent **after the write is completed** and is completed itself.

If the device is backed by persistent storage, the device **MUST** ensure that stable writes are committed to it, before reporting completion of the write (cases 1 and 2) or the flush (case 3). Failure to do so can cause data loss in case of a crash.

If the driver changes *writeback* between the submission of the write and its completion, the write could be either volatile or stable when its completion is reported; in other words, the exact behavior is undefined.

If `VIRTIO_BLK_F_FLUSH` was not offered by the device<sup>5</sup>, the device **MAY** also commit writes to persistent device backend storage before reporting their completion. Unlike case 1, however, this is not an absolute requirement of the specification.

**Note:** An implementation that does not offer `VIRTIO_BLK_F_FLUSH` and does not commit completed writes will not be resilient to data loss in case of crashes. Not offering `VIRTIO_BLK_F_FLUSH` is an absolute requirement for implementations that do not wish to be safe against such data losses.

If the device is backed by storage providing lifetime metrics (such as eMMC or UFS persistent storage), the device **SHOULD** offer the `VIRTIO_BLK_F_LIFETIME` flag. The flag **MUST NOT** be offered if the device is backed by storage for which the lifetime metrics described in this document cannot be obtained or for which such metrics have no useful meaning. If the metrics are offered, the device **MUST NOT** send any reserved values, as defined in this specification.

**Note:** The device lifetime metrics *pre\_eol\_info*, *device\_lifetime\_est\_a* and *device\_lifetime\_est\_b* are discussed in the JESD84-B50 specification.

The complete JESD84-B50 is available at the JEDEC website (<https://www.jedec.org>) pursuant to JEDEC's licensing terms and conditions. This information is provided to simplify passthrough implementations from eMMC devices.

### 5.2.6.3 Legacy Interface: Device Operation

When using the legacy interface, transitional devices and drivers **MUST** format the fields in struct `virtio_blk_req` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

When using the legacy interface, transitional drivers **SHOULD** ignore the used length values.

**Note:** Historically, some devices put the total descriptor length, or the total length of device-writable buffers there, even when only the status byte was actually written.

The *reserved* field was previously called *ioprio*. *ioprio* is a hint about the relative priorities of requests to the device: higher numbers indicate more important requests.

```
#define VIRTIO_BLK_T_FLUSH_OUT    5
```

The command `VIRTIO_BLK_T_FLUSH_OUT` was a synonym for `VIRTIO_BLK_T_FLUSH`; a driver **MUST** treat it as a `VIRTIO_BLK_T_FLUSH` command.

```
#define VIRTIO_BLK_T_BARRIER    0x80000000
```

If the device has `VIRTIO_BLK_F_BARRIER` feature the high bit (`VIRTIO_BLK_T_BARRIER`) indicates that this request acts as a barrier and that all preceding requests **SHOULD** be complete before this one, and all following requests **SHOULD NOT** be started until this is complete.

**Note:** A barrier does not flush caches in the underlying backend device in host, and thus does not serve as data consistency guarantee. Only a `VIRTIO_BLK_T_FLUSH` request does that.

Some older legacy devices did not commit completed writes to persistent device backend storage when `VIRTIO_BLK_F_FLUSH` was offered but not negotiated. In order to work around this, the driver **MAY** set the *writeback* to 0 (if available) or it **MAY** send an explicit flush request after every completed write.

If the device has `VIRTIO_BLK_F_SCSI` feature, it can also support scsi packet command requests, each of these requests is of form:

<sup>5</sup>Note that in this case, according to 5.2.5.2, the device will not have offered `VIRTIO_BLK_F_CONFIG_WCE` either.

```

/* All fields are in guest's native endian. */
struct virtio_scsi_pc_req {
    u32 type;
    u32 ioprio;
    u64 sector;
    u8 cmd[];
    u8 data[][512];
#define SCSI_SENSE_BUFFERSIZE 96
    u8 sense[SCSI_SENSE_BUFFERSIZE];
    u32 errors;
    u32 data_len;
    u32 sense_len;
    u32 residual;
    u8 status;
};

```

A request type can also be a scsi packet command (VIRTIO\_BLK\_T SCSI\_CMD or VIRTIO\_BLK\_T SCSI\_CMD\_OUT). The two types are equivalent, the device does not distinguish between them:

```

#define VIRTIO_BLK_T SCSI_CMD 2
#define VIRTIO_BLK_T SCSI_CMD_OUT 3

```

The *cmd* field is only present for scsi packet command requests, and indicates the command to perform. This field **MUST** reside in a single, separate device-readable buffer; command length can be derived from the length of this buffer.

Note that these first three (four for scsi packet commands) fields are always device-readable: *data* is either device-readable or device-writable, depending on the request. The size of the read or write can be derived from the total size of the request buffers.

*sense* is only present for scsi packet command requests, and indicates the buffer for scsi sense data.

*data\_len* is only present for scsi packet command requests, this field is deprecated, and **SHOULD** be ignored by the driver. Historically, devices copied data length there.

*sense\_len* is only present for scsi packet command requests and indicates the number of bytes actually written to the *sense* buffer.

*residual* field is only present for scsi packet command requests and indicates the residual size, calculated as data length - number of bytes actually transferred.

#### 5.2.6.4 Legacy Interface: Framing Requirements

When using legacy interfaces, transitional drivers which have not negotiated VIRTIO\_F\_ANY\_LAYOUT:

- **MUST** use a single 8-byte descriptor containing *type*, *reserved* and *sector*, followed by descriptors for *data*, then finally a separate 1-byte descriptor for *status*.
- For SCSI commands there are additional constraints. *sense* **MUST** reside in a single separate device-writable descriptor of size 96 bytes, and *errors*, *data\_len*, *sense\_len* and *residual* **MUST** reside a single separate device-writable descriptor.

See [2.7.4](#).

## 5.3 Console Device

The virtio console device is a simple device for data input and output. A device **MAY** have one or more ports. Each port has a pair of input and output virtqueues. Moreover, a device has a pair of control IO virtqueues. The control virtqueues are used to communicate information between the device and the driver about ports being opened and closed on either side of the connection, indication from the device about whether a particular port is a console port, adding new ports, port hot-plug/unplug, etc., and indication from the driver about whether a port or a device was successfully added, port open/close, etc. For data IO, one or more empty buffers are placed in the receive queue for incoming data and outgoing characters are placed in the transmit queue.

### 5.3.1 Device ID

3

### 5.3.2 Virtqueues

- 0 receiveq(port0)
- 1 transmitq(port0)
- 2 control receiveq
- 3 control transmitq
- 4 receiveq(port1)
- 5 transmitq(port1)

...

The port 0 receive and transmit queues always exist: other queues only exist if VIRTIO\_CONSOLE\_F\_MULTIPORT is set.

### 5.3.3 Feature bits

**VIRTIO\_CONSOLE\_F\_SIZE (0)** Configuration *cols* and *rows* are valid.

**VIRTIO\_CONSOLE\_F\_MULTIPORT (1)** Device has support for multiple ports; *max\_nr\_ports* is valid and control virtqueues will be used.

**VIRTIO\_CONSOLE\_F\_EMERG\_WRITE (2)** Device has support for emergency write. Configuration field *emerg\_wr* is valid.

### 5.3.4 Device configuration layout

The size of the console is supplied in the configuration space if the VIRTIO\_CONSOLE\_F\_SIZE feature is set. Furthermore, if the VIRTIO\_CONSOLE\_F\_MULTIPORT feature is set, the maximum number of ports supported by the device can be fetched.

If VIRTIO\_CONSOLE\_F\_EMERG\_WRITE is set then the driver can use emergency write to output a single character without initializing virtio queues, or even acknowledging the feature.

```
struct virtio_console_config {  
    le16 cols;  
    le16 rows;  
    le32 max_nr_ports;  
    le32 emerg_wr;  
};
```

#### 5.3.4.1 Legacy Interface: Device configuration layout

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_console_config` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

### 5.3.5 Device Initialization

1. If the VIRTIO\_CONSOLE\_F\_EMERG\_WRITE feature is offered, *emerg\_wr* field of the configuration can be written at any time. Thus it works for very early boot debugging output as well as catastrophic OS failures (eg. virtio ring corruption).
2. If the VIRTIO\_CONSOLE\_F\_SIZE feature is negotiated, the driver can read the console dimensions from *cols* and *rows*.

3. If the `VIRTIO_CONSOLE_F_MULTIPORT` feature is negotiated, the driver can spawn multiple ports, not all of which are necessarily attached to a console. Some could be generic ports. In this case, the control virtqueues are enabled and according to *max\_nr\_ports*, the appropriate number of virtqueues are created. A control message indicating the driver is ready is sent to the device. The device can then send control messages for adding new ports to the device. After creating and initializing each port, a `VIRTIO_CONSOLE_PORT_READY` control message is sent to the device for that port so the device can let the driver know of any additional configuration options set for that port.
4. The receiveq for each port is populated with one or more receive buffers.

#### 5.3.5.1 Device Requirements: Device Initialization

The device **MUST** allow a write to *emerg\_wr*, even on an unconfigured device.

The device **SHOULD** transmit the lower byte written to *emerg\_wr* to an appropriate log or output method.

### 5.3.6 Device Operation

1. For output, a buffer containing the characters is placed in the port's transmitq<sup>6</sup>.
2. When a buffer is used in the receiveq (signalled by a used buffer notification), the contents is the input to the port associated with the virtqueue for which the notification was received.
3. If the driver negotiated the `VIRTIO_CONSOLE_F_SIZE` feature, a configuration change notification indicates that the updated size can be read from the configuration fields. This size applies to port 0 only.
4. If the driver negotiated the `VIRTIO_CONSOLE_F_MULTIPORT` feature, active ports are announced by the device using the `VIRTIO_CONSOLE_PORT_ADD` control message. The same message is used for port hot-plug as well.

#### 5.3.6.1 Driver Requirements: Device Operation

The driver **MUST NOT** put a device-readable buffer in a receiveq. The driver **MUST NOT** put a device-writable buffer in a transmitq.

#### 5.3.6.2 Multiport Device Operation

If the driver negotiated the `VIRTIO_CONSOLE_F_MULTIPORT`, the two control queues are used to manipulate the different console ports: the control receiveq for messages from the device to the driver, and the control sendq for driver-to-device messages. The layout of the control messages is:

```
struct virtio_console_control {
    le32 id; /* Port number */
    le16 event; /* The kind of control event */
    le16 value; /* Extra information for the event */
};
```

The values for *event* are:

**VIRTIO\_CONSOLE\_DEVICE\_READY (0)** Sent by the driver at initialization to indicate that it is ready to receive control messages. A value of 1 indicates success, and 0 indicates failure. The port number *id* is unused.

**VIRTIO\_CONSOLE\_DEVICE\_ADD (1)** Sent by the device, to create a new port. *value* is unused.

**VIRTIO\_CONSOLE\_DEVICE\_REMOVE (2)** Sent by the device, to remove an existing port. *value* is unused.

<sup>6</sup>Because this is high importance and low bandwidth, the current Linux implementation polls for the buffer to become used, rather than waiting for a used buffer notification, simplifying the implementation significantly. However, for generic serial ports with the `O_NONBLOCK` flag set, the polling limitation is relaxed and the consumed buffers are freed upon the next write or poll call or when a port is closed or hot-unplugged.

**VIRTIO\_CONSOLE\_PORT\_READY (3)** Sent by the driver in response to the device's VIRTIO\_CONSOLE\_PORT\_ADD message, to indicate that the port is ready to be used. A *value* of 1 indicates success, and 0 indicates failure.

**VIRTIO\_CONSOLE\_CONSOLE\_PORT (4)** Sent by the device to nominate a port as a console port. There MAY be more than one console port.

**VIRTIO\_CONSOLE\_RESIZE (5)** Sent by the device to indicate a console size change. *value* is unused. The buffer is followed by the number of columns and rows:

```
struct virtio_console_resize {
    le16 cols;
    le16 rows;
};
```

**VIRTIO\_CONSOLE\_PORT\_OPEN (6)** This message is sent by both the device and the driver. *value* indicates the state: 0 (port closed) or 1 (port open). This allows for ports to be used directly by guest and host processes to communicate in an application-defined manner.

**VIRTIO\_CONSOLE\_PORT\_NAME (7)** Sent by the device to give a tag to the port. This control command is immediately followed by the UTF-8 name of the port for identification within the guest (without a NUL terminator).

#### 5.3.6.2.1 Device Requirements: Multiport Device Operation

The device MUST NOT specify a port which exists in a VIRTIO\_CONSOLE\_DEVICE\_ADD message, nor a port which is equal or greater than *max\_nr\_ports*.

The device MUST NOT specify a port in VIRTIO\_CONSOLE\_DEVICE\_REMOVE which has not been created with a previous VIRTIO\_CONSOLE\_DEVICE\_ADD.

#### 5.3.6.2.2 Driver Requirements: Multiport Device Operation

The driver MUST send a VIRTIO\_CONSOLE\_DEVICE\_READY message if VIRTIO\_CONSOLE\_F\_MULTIPORT is negotiated.

Upon receipt of a VIRTIO\_CONSOLE\_CONSOLE\_PORT message, the driver SHOULD treat the port in a manner suitable for text console access and MUST respond with a VIRTIO\_CONSOLE\_PORT\_OPEN message, which MUST have *value* set to 1.

#### 5.3.6.3 Legacy Interface: Device Operation

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_console_control` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

When using the legacy interface, the driver SHOULD ignore the used length values for the transmit queues and the control transmitq.

**Note:** Historically, some devices put the total descriptor length there, even though no data was actually written.

#### 5.3.6.4 Legacy Interface: Framing Requirements

When using legacy interfaces, transitional drivers which have not negotiated VIRTIO\_F\_ANY\_LAYOUT MUST use only a single descriptor for all buffers in the control receiveq and control transmitq.

## 5.4 Entropy Device

The virtio entropy device supplies high-quality randomness for guest use.



### 5.4.1 Device ID

4

### 5.4.2 Virtqueues

0 requestq

### 5.4.3 Feature bits

None currently defined

### 5.4.4 Device configuration layout

None currently defined.

### 5.4.5 Device Initialization

1. The virtqueue is initialized

### 5.4.6 Device Operation

When the driver requires random bytes, it places the descriptor of one or more buffers in the queue. It will be completely filled by random data by the device.

#### 5.4.6.1 Driver Requirements: Device Operation

The driver **MUST NOT** place device-readable buffers into the queue.

The driver **MUST** examine the length written by the device to determine how many random bytes were received.

#### 5.4.6.2 Device Requirements: Device Operation

The device **MUST** place one or more random bytes into the buffer, but it **MAY** use less than the entire buffer length.

## 5.5 Traditional Memory Balloon Device

This is the traditional balloon device. The device number 13 is reserved for a new memory balloon interface, with different semantics, which is expected in a future version of the standard.

The traditional virtio memory balloon device is a primitive device for managing guest memory: the device asks for a certain amount of memory, and the driver supplies it (or withdraws it, if the device has more than it asks for). This allows the guest to adapt to changes in allowance of underlying physical memory. If the feature is negotiated, the device can also be used to communicate guest memory statistics to the host.

### 5.5.1 Device ID

5

### 5.5.2 Virtqueues

0 inflateq

1 deflateq

2 statsq

3 free\_page\_vq

#### 4 reporting\_vq

statsq only exists if VIRTIO\_BALLOON\_F\_STATS\_VQ is set.

free\_page\_vq only exists if VIRTIO\_BALLOON\_F\_FREE\_PAGE\_HINT is set.

reporting\_vq only exists if VIRTIO\_BALLOON\_F\_PAGE\_REPORTING is set.

### 5.5.3 Feature bits

**VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST (0)** Host has to be told before pages from the balloon are used.

**VIRTIO\_BALLOON\_F\_STATS\_VQ (1)** A virtqueue for reporting guest memory statistics is present.

**VIRTIO\_BALLOON\_F\_DEFLATE\_ON\_OOM (2)** Deflate balloon on guest out of memory condition.

**VIRTIO\_BALLOON\_F\_FREE\_PAGE\_HINT(3)** The device has support for free page hinting. A virtqueue for providing hints as to what memory is currently free is present. Configuration field *free\_page\_hint\_cmd\_id* is valid.

**VIRTIO\_BALLOON\_F\_PAGE\_POISON(4)** A hint to the device, that the driver will immediately write *poison\_val* to pages after deflating them. Configuration field *poison\_val* is valid.

**VIRTIO\_BALLOON\_F\_PAGE\_REPORTING(5)** The device has support for free page reporting. A virtqueue for reporting free guest memory is present.

#### 5.5.3.1 Driver Requirements: Feature bits

The driver SHOULD accept the VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature if offered by the device.

The driver SHOULD clear the VIRTIO\_BALLOON\_F\_PAGE\_POISON flag if it will not immediately write *poison\_val* to deflated pages (e.g., to initialize them, or fill them with a poison value).

If the driver is expecting the pages to retain some initialized value, it MUST NOT accept VIRTIO\_BALLOON\_F\_PAGE\_REPORTING unless it also negotiates VIRTIO\_BALLOON\_F\_PAGE\_POISON.

#### 5.5.3.2 Device Requirements: Feature bits

If the device offers the VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature bit, and if the driver did not accept this feature bit, the device MAY signal failure by failing to set FEATURES\_OK *device status* bit when the driver writes it.

##### 5.5.3.2.0.1 Legacy Interface: Feature bits

As the legacy interface does not have a way to gracefully report feature negotiation failure, when using the legacy interface, transitional devices MUST support guests which do not negotiate VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature, and SHOULD allow guest to use memory before notifying host if VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST is not negotiated.

### 5.5.4 Device configuration layout

*num\_pages* and *actual* are always available.

*free\_page\_hint\_cmd\_id* is available if VIRTIO\_BALLOON\_F\_FREE\_PAGE\_HINT has been negotiated. The field is read-only by the driver. *poison\_val* is available if VIRTIO\_BALLOON\_F\_PAGE\_POISON has been negotiated.

```
struct virtio_balloon_config {
    le32 num_pages;
    le32 actual;
    le32 free_page_hint_cmd_id;
    le32 poison_val;
};
```

#### 5.5.4.0.0.1 Legacy Interface: Device configuration layout

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_balloon_config` according to the little-endian format.

**Note:** This is unlike the usual convention that legacy device fields are guest endian.

### 5.5.5 Device Initialization

The device initialization process is outlined below:

1. The inflate and deflate virtqueues are identified.
2. If the `VIRTIO_BALLOON_F_STATS_VQ` feature bit is negotiated:
  - (a) Identify the stats virtqueue.
  - (b) Add one empty buffer to the stats virtqueue.
3. If the `VIRTIO_BALLOON_F_FREE_PAGE_HINT` feature bit is negotiated, identify the `free_page_vq`.
4. If the `VIRTIO_BALLOON_F_PAGE_POISON` feature bit is negotiated, update the *poison\_val* configuration field.
5. If the `VIRTIO_BALLOON_F_PAGE_REPORTING` feature bit is negotiated, identify the `reporting_vq`.
6. `DRIVER_OK` is set: device operation begins.
7. If the `VIRTIO_BALLOON_F_STATS_VQ` feature bit is negotiated, then notify the device about the stats virtqueue buffer.
8. If the `VIRTIO_BALLOON_F_PAGE_REPORTING` feature bit is negotiated, then begin reporting free pages to the device.

### 5.5.6 Device Operation

The device is driven either by the receipt of a configuration change notification, or by changing guest memory needs, such as performing memory compaction or responding to out of memory conditions.

1. *num\_pages* configuration field is examined. If this is greater than the *actual* number of pages, the balloon wants more memory from the guest. If it is less than *actual*, the balloon doesn't need it all.
2. To supply memory to the balloon (aka. inflate):
  - (a) The driver constructs an array of addresses of unused memory pages. These addresses are divided by 4096<sup>7</sup> and the descriptor describing the resulting 32-bit array is added to the `inflateq`.
3. To remove memory from the balloon (aka. deflate):
  - (a) The driver constructs an array of addresses of memory pages it has previously given to the balloon, as described above. This descriptor is added to the `deflateq`.
  - (b) If the `VIRTIO_BALLOON_F_MUST_TELL_HOST` feature is negotiated, the guest informs the device of pages before it uses them.
  - (c) Otherwise, the guest is allowed to re-use pages previously given to the balloon before the device has acknowledged their withdrawal<sup>8</sup>.
4. In either case, the device acknowledges inflate and deflate requests by using the descriptor.
5. Once the device has acknowledged the inflation or deflation, the driver updates *actual* to reflect the new number of pages in the balloon.

---

<sup>7</sup>This is historical, and independent of the guest page size.

<sup>8</sup>In this case, deflation advice is merely a courtesy.

### 5.5.6.1 Driver Requirements: Device Operation

The driver SHOULD supply pages to the balloon when *num\_pages* is greater than the actual number of pages in the balloon.

The driver MAY use pages from the balloon when *num\_pages* is less than the actual number of pages in the balloon.

The driver MAY supply pages to the balloon when *num\_pages* is greater than or equal to the actual number of pages in the balloon.

If VIRTIO\_BALLOON\_F\_DEFLATE\_ON\_OOM has not been negotiated, the driver MUST NOT use pages from the balloon when *num\_pages* is less than or equal to the actual number of pages in the balloon.

If VIRTIO\_BALLOON\_F\_DEFLATE\_ON\_OOM has been negotiated, the driver MAY use pages from the balloon when *num\_pages* is less than or equal to the actual number of pages in the balloon if this is required for system stability (e.g. if memory is required by applications running within the guest).

The driver MUST use the *deflateq* to inform the device of pages that it wants to use from the balloon.

If the VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature is negotiated, the driver MUST NOT use pages from the balloon until the device has acknowledged the deflate request.

Otherwise, if the VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature is not negotiated, the driver MAY begin to re-use pages previously given to the balloon before the device has acknowledged the deflate request.

In any case, the driver MUST NOT use pages from the balloon after adding the pages to the balloon, but before the device has acknowledged the inflate request.

The driver MUST NOT request deflation of pages in the balloon before the device has acknowledged the inflate request.

The driver MUST update *actual* after changing the number of pages in the balloon.

The driver MAY update *actual* once after multiple inflate and deflate operations.

### 5.5.6.2 Device Requirements: Device Operation

The device MAY modify the contents of a page in the balloon after detecting its physical number in an inflate request and before acknowledging the inflate request by using the *inflateq* descriptor.

If the VIRTIO\_BALLOON\_F\_MUST\_TELL\_HOST feature is negotiated, the device MAY modify the contents of a page in the balloon after detecting its physical number in an inflate request and before detecting its physical number in a deflate request and acknowledging the deflate request.

#### 5.5.6.2.1 Legacy Interface: Device Operation

When using the legacy interface, the driver SHOULD ignore the used length values.

**Note:** Historically, some devices put the total descriptor length there, even though no data was actually written.

When using the legacy interface, the driver MUST write out all 4 bytes each time it updates the *actual* value in the configuration space, using a single atomic operation.

When using the legacy interface, the device SHOULD NOT use the *actual* value written by the driver in the configuration space, until the last, most-significant byte of the value has been written.

**Note:** Historically, devices used the *actual* value, even though when using Virtio Over PCI Bus the device-specific configuration space was not guaranteed to be atomic. Using intermediate values during update by driver is best avoided, except for debugging.

Historically, drivers using Virtio Over PCI Bus wrote the *actual* value by using multiple single-byte writes in order, from the least-significant to the most-significant value.

### 5.5.6.3 Memory Statistics

The stats virtqueue is atypical because communication is driven by the device (not the driver). The channel becomes active at driver initialization time when the driver adds an empty buffer and notifies the device. A request for memory statistics proceeds as follows:

1. The device uses the buffer and sends a used buffer notification.
2. The driver pops the used buffer and discards it.
3. The driver collects memory statistics and writes them into a new buffer.
4. The driver adds the buffer to the virtqueue and notifies the device.
5. The device pops the buffer (retaining it to initiate a subsequent request) and consumes the statistics.

Within the buffer, statistics are an array of 10-byte entries. Each statistic consists of a 16 bit tag and a 64 bit value. All statistics are optional and the driver chooses which ones to supply. To guarantee backwards compatibility, devices omit unsupported statistics.

```
struct virtio_balloon_stat {  
#define VIRTIO_BALLOON_S_SWAP_IN  0  
#define VIRTIO_BALLOON_S_SWAP_OUT 1  
#define VIRTIO_BALLOON_S_MAJFLT   2  
#define VIRTIO_BALLOON_S_MINFLT   3  
#define VIRTIO_BALLOON_S_MEMFREE  4  
#define VIRTIO_BALLOON_S_MEMTOT   5  
#define VIRTIO_BALLOON_S_AVAIL    6  
#define VIRTIO_BALLOON_S_CACHES   7  
#define VIRTIO_BALLOON_S_HTLB_PGALLOC 8  
#define VIRTIO_BALLOON_S_HTLB_PGFAIL 9  
    le16 tag;  
    le64 val;  
} __attribute__((packed));
```

#### 5.5.6.3.1 Driver Requirements: Memory Statistics

Normative statements in this section apply if and only if the VIRTIO\_BALLOON\_F\_STATS\_VQ feature has been negotiated.

The driver **MUST** make at most one buffer available to the device in the statsq, at all times.

After initializing the device, the driver **MUST** make an output buffer available in the statsq.

Upon detecting that device has used a buffer in the statsq, the driver **MUST** make an output buffer available in the statsq.

Before making an output buffer available in the statsq, the driver **MUST** initialize it, including one struct virtio\_balloon\_stat entry for each statistic that it supports.

Driver **MUST** use an output buffer size which is a multiple of 6 bytes for all buffers submitted to the statsq.

Driver **MAY** supply struct virtio\_balloon\_stat entries in the output buffer submitted to the statsq in any order, without regard to tag values.

Driver **MAY** supply a subset of all statistics in the output buffer submitted to the statsq.

Driver **MUST** supply the same subset of statistics in all buffers submitted to the statsq.

#### 5.5.6.3.2 Device Requirements: Memory Statistics

Normative statements in this section apply if and only if the VIRTIO\_BALLOON\_F\_STATS\_VQ feature has been negotiated.

Within an output buffer submitted to the statsq, the device **MUST** ignore entries with tag values that it does not recognize.

Within an output buffer submitted to the statsq, the device MUST accept struct `virtio_balloon_stat` entries in any order without regard to *tag* values.

#### 5.5.6.3.3 Legacy Interface: Memory Statistics

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_balloon_stat` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

When using the legacy interface, the device SHOULD ignore all values in the first buffer in the statsq supplied by the driver after device initialization.

**Note:** Historically, drivers supplied an uninitialized buffer in the first buffer.

#### 5.5.6.4 Memory Statistics Tags

**VIRTIO\_BALLOON\_S\_SWAP\_IN (0)** The amount of memory that has been swapped in (in bytes).

**VIRTIO\_BALLOON\_S\_SWAP\_OUT (1)** The amount of memory that has been swapped out to disk (in bytes).

**VIRTIO\_BALLOON\_S\_MAJFLT (2)** The number of major page faults that have occurred.

**VIRTIO\_BALLOON\_S\_MINFLT (3)** The number of minor page faults that have occurred.

**VIRTIO\_BALLOON\_S\_MEMFREE (4)** The amount of memory not being used for any purpose (in bytes).

**VIRTIO\_BALLOON\_S\_MEMTOT (5)** The total amount of memory available (in bytes).

**VIRTIO\_BALLOON\_S\_AVAIL (6)** An estimate of how much memory is available (in bytes) for starting new applications, without pushing the system to swap.

**VIRTIO\_BALLOON\_S\_CACHES (7)** The amount of memory, in bytes, that can be quickly reclaimed without additional I/O. Typically these pages are used for caching files from disk.

**VIRTIO\_BALLOON\_S\_HTLB\_PGALLOC (8)** The number of successful hugetlb page allocations in the guest.

**VIRTIO\_BALLOON\_S\_HTLB\_PGFAIL (9)** The number of failed hugetlb page allocations in the guest.

#### 5.5.6.5 Free Page Hinting

Free page hinting is designed to be used during migration to determine what pages within the guest are currently unused so that they can be skipped over while migrating the guest. The device will indicate that it is ready to start performing hinting by setting the *free\_page\_hint\_cmd\_id* to one of the non-reserved values that can be used as a command ID. The following values are reserved:

**VIRTIO\_BALLOON\_CMD\_ID\_STOP (0)** Any command ID previously supplied by the device is invalid. The driver should stop hinting free pages until a new command ID is supplied, but should not release any hinted pages for use by the guest.

**VIRTIO\_BALLOON\_CMD\_ID\_DONE (1)** Any command ID previously supplied by the device is invalid. The driver should stop hinting free pages, and should release all hinted pages for use by the guest.

When a hint is provided by the driver it indicates that the data contained in the given page is no longer needed and can be discarded. If the driver writes to the page this overrides the hint and the data will be retained. The contents of any stale pages that have not been written to since the page was hinted may be lost, and if read the contents of such pages will be uninitialized memory.

A request for free page hinting proceeds as follows:

1. The driver examines the *free\_page\_hint\_cmd\_id* configuration field. If it contains a non-reserved value then free page hinting will begin.
2. To supply free page hints:

- (a) The driver constructs an output buffer containing the new value from the *free\_page\_hint\_cmd\_id* configuration field and adds it to the *free\_page\_vq*.
  - (b) The driver maps a series of pages and adds them to the *free\_page\_vq* as individual scatter-gather input buffer entries.
  - (c) When the driver is no longer able to fetch additional pages to add to the *free\_page\_vq*, it will construct an output buffer containing the command ID *VIRTIO\_BALLOON\_CMD\_ID\_STOP*.
3. A round of hinting ends either when the driver is no longer able to supply more pages for hinting as described above, or when the device updates *free\_page\_hint\_cmd\_id* configuration field to contain either *VIRTIO\_BALLOON\_CMD\_ID\_STOP* or *VIRTIO\_BALLOON\_CMD\_ID\_DONE*.
  4. The device may follow *VIRTIO\_BALLOON\_CMD\_ID\_STOP* with a new non-reserved value for the *free\_page\_hint\_cmd\_id* configuration field in which case it will resume supplying free page hints.
  5. Otherwise, if the device provides *VIRTIO\_BALLOON\_CMD\_ID\_DONE* then hinting is complete and the driver may release all previously hinted pages for use by the guest.

#### 5.5.6.5.1 Driver Requirements: Free Page Hinting

Normative statements in this section apply if the *VIRTIO\_BALLOON\_F\_FREE\_PAGE\_HINT* feature has been negotiated.

The driver **MUST** use an output buffer size of 4 bytes for all output buffers submitted to the *free\_page\_vq*.

The driver **MUST** start hinting by providing an output buffer containing the current command ID for the given block of pages.

The driver **MUST NOT** provide more than one output buffer containing the current command ID.

The driver **SHOULD** supply pages to the *free\_page\_vq* as input buffers when *free\_page\_hint\_cmd\_id* specifies a value of 2 or greater.

The driver **SHOULD** stop supplying pages for hinting when *free\_page\_hint\_cmd\_id* specifies a value of *VIRTIO\_BALLOON\_CMD\_ID\_STOP* or *VIRTIO\_BALLOON\_CMD\_ID\_DONE*.

If the driver is unable to supply pages, it **MUST** complete hinting by adding an output buffer containing the command ID *VIRTIO\_BALLOON\_CMD\_ID\_STOP*.

The driver **MAY** release hinted pages for use by the guest including when the device has not yet used the descriptor containing the hinting request.

The driver **MUST** treat the content of all hinted pages as uninitialized memory.

The driver **MUST** initialize the contents of any previously hinted page released before *free\_page\_hint\_cmd\_id* specifies a value of *VIRTIO\_BALLOON\_CMD\_ID\_DONE*.

The driver **SHOULD** release all previously hinted pages once *free\_page\_hint\_cmd\_id* specifies a value of *VIRTIO\_BALLOON\_CMD\_ID\_DONE*.

#### 5.5.6.5.2 Device Requirements: Free Page Hinting

Normative statements in this section apply if the *VIRTIO\_BALLOON\_F\_FREE\_PAGE\_HINT* feature has been negotiated.

The device **SHOULD** set *free\_page\_hint\_cmd\_id* to *VIRTIO\_BALLOON\_CMD\_ID\_STOP* any time that it will not be able to make use of the hints provided by the driver.

The device **MUST NOT** reuse a command ID until it has received an output buffer containing *VIRTIO\_BALLOON\_CMD\_ID\_STOP* from the driver.

The device **MUST** ignore pages that are provided with a command ID that does not match the current value in *free\_page\_hint\_cmd\_id*.



If the content of a previously hinted page has not been modified by the guest since the device issued the *free\_page\_hint\_cmd\_id* associated with the hint, the device MAY modify the contents of the page.

The device MUST NOT modify the content of a previously hinted page after *free\_page\_hint\_cmd\_id* is set to `VIRTIO_BALLOON_CMD_ID_DONE`.

The device MUST report a value of `VIRTIO_BALLOON_CMD_ID_DONE` in *free\_page\_hint\_cmd\_id* when it no longer has need for the previously hinted pages.

#### 5.5.6.5.3 Legacy Interface: Free Page Hinting

When using the legacy interface, transitional devices and drivers MUST format the command ID field in output buffers according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.5.6.6 Page Poison

Page Poison provides a way to notify the host that the guest is initializing free pages with *poison\_val*. When the feature is enabled, pages will be immediately written to by the driver after deflating, and pages reported by free page reporting will retain the value indicated by *poison\_val*.

If the guest is not initializing freed pages, the driver should reject the `VIRTIO_BALLOON_F_PAGE_POISON` feature.

If `VIRTIO_BALLOON_F_PAGE_POISON` feature has been negotiated, the driver will place the initialization value into the *poison\_val* configuration field data.

##### 5.5.6.6.1 Driver Requirements: Page Poison

Normative statements in this section apply if the `VIRTIO_BALLOON_F_PAGE_POISON` feature has been negotiated.

The driver MUST initialize the deflated pages with *poison\_val* when they are reused by the driver.

The driver MUST populate the *poison\_val* configuration data before setting the `DRIVER_OK` bit.

The driver MUST NOT modify *poison\_val* while the `DRIVER_OK` bit is set.

##### 5.5.6.6.2 Device Requirements: Page Poison

Normative statements in this section apply if the `VIRTIO_BALLOON_F_PAGE_POISON` feature has been negotiated.

The device MAY use the content of *poison\_val* as a hint to guest behavior.

#### 5.5.6.7 Free Page Reporting

Free Page Reporting provides a mechanism similar to balloon inflation, however it does not provide a deflation queue. Reported free pages can be reused by the driver after the reporting request has been acknowledged without notifying the device.

The driver will begin reporting free pages. When exactly and which free pages are reported is up to the driver.

1. The driver determines it has enough pages available to begin reporting free pages.
2. The driver gathers free pages into a scatter-gather list and adds them to the `reporting_vq`.
3. The device acknowledges the reporting request by using the `reporting_vq` descriptor.
4. Once the device has acknowledged the report, the driver can reuse the reported free pages when needed (e.g., by putting them back to free page lists in the guest operating system).
5. The driver can then continue to gather and report free pages until it has determined it has reported a sufficient quantity of pages.

#### 5.5.6.7.1 Driver Requirements: Free Page Reporting

Normative statements in this section apply if the `VIRTIO_BALLOON_F_PAGE_REPORTING` feature has been negotiated.

If the `VIRTIO_BALLOON_F_PAGE_POISON` feature has not been negotiated, then the driver **MUST** treat all reported pages as uninitialized memory.

If the `VIRTIO_BALLOON_F_PAGE_POISON` feature has been negotiated, the driver **MUST** initialize all free pages with *poison\_val* before reporting them.

The driver **MUST NOT** use the reported pages until the device has acknowledged the reporting request.

The driver **MAY** report free pages any time after `DRIVER_OK` is set.

The driver **SHOULD** attempt to report large pages rather than smaller ones.

The driver **SHOULD** avoid reading/writing reported pages if not strictly necessary.

#### 5.5.6.7.2 Device Requirements: Free Page Reporting

Normative statements in this section apply if the `VIRTIO_BALLOON_F_PAGE_REPORTING` feature has been negotiated.

If the `VIRTIO_BALLOON_F_PAGE_POISON` feature has not been negotiated, the device **MAY** modify the contents of any page supplied in a report request before acknowledging that request by using the `reporting_vq` descriptor.

If the `VIRTIO_BALLOON_F_PAGE_POISON` feature has been negotiated, the device **MUST NOT** modify the the content of a reported page to a value other than *poison\_val*.

## 5.6 SCSI Host Device

The virtio SCSI host device groups together one or more virtual logical units (such as disks), and allows communicating to them using the SCSI protocol. An instance of the device represents a SCSI host to which many targets and LUNs are attached.

The virtio SCSI device services two kinds of requests:

- command requests for a logical unit;
- task management functions related to a logical unit, target or command.

The device is also able to send out notifications about added and removed logical units. Together, these capabilities provide a SCSI transport protocol that uses virtqueues as the transfer medium. In the transport protocol, the virtio driver acts as the initiator, while the virtio SCSI host provides one or more targets that receive and process the requests.

This section relies on definitions from [SAM](#).

### 5.6.1 Device ID

8

### 5.6.2 Virtqueues

0 controlq

1 eventq

2...n request queues

### 5.6.3 Feature bits

**VIRTIO\_SCSI\_F\_INOUT (0)** A single request can include both device-readable and device-writable data buffers.

**VIRTIO\_SCSI\_F\_HOTPLUG (1)** The host SHOULD enable reporting of hot-plug and hot-unplug events for LUNs and targets on the SCSI bus. The guest SHOULD handle hot-plug and hot-unplug events.

**VIRTIO\_SCSI\_F\_CHANGE (2)** The host will report changes to LUN parameters via a VIRTIO\_SCSI\_T\_-PARAM\_CHANGE event; the guest SHOULD handle them.

**VIRTIO\_SCSI\_F\_T10\_PI (3)** The extended fields for T10 protection information (DIF/DIX) are included in the SCSI request header.

### 5.6.4 Device configuration layout

All fields of this configuration are always available.

```
struct virtio_scsi_config {
    le32 num_queues;
    le32 seg_max;
    le32 max_sectors;
    le32 cmd_per_lun;
    le32 event_info_size;
    le32 sense_size;
    le32 cdb_size;
    le16 max_channel;
    le16 max_target;
    le32 max_lun;
};
```

**num\_queues** is the total number of request virtqueues exposed by the device. The driver MAY use only one request queue, or it can use more to achieve better performance.

**seg\_max** is the maximum number of segments that can be in a command. A bidirectional command can include *seg\_max* input segments and *seg\_max* output segments.

**max\_sectors** is a hint to the driver about the maximum transfer size to use.

**cmd\_per\_lun** tells the driver the maximum number of linked commands it can send to one LUN.

**event\_info\_size** is the maximum size that the device will fill for buffers that the driver places in the eventq. It is written by the device depending on the set of negotiated features.

**sense\_size** is the maximum size of the sense data that the device will write. The default value is written by the device and MUST be 96, but the driver can modify it. It is restored to the default when the device is reset.

**cdb\_size** is the maximum size of the CDB that the driver will write. The default value is written by the device and MUST be 32, but the driver can likewise modify it. It is restored to the default when the device is reset.

**max\_channel, max\_target and max\_lun** can be used by the driver as hints to constrain scanning the logical units on the host to channel/target/logical unit numbers that are less than or equal to the value of the fields. *max\_channel* SHOULD be zero. *max\_target* SHOULD be less than or equal to 255. *max\_lun* SHOULD be less than or equal to 16383.

#### 5.6.4.1 Driver Requirements: Device configuration layout

The driver MUST NOT write to device configuration fields other than *sense\_size* and *cdb\_size*.

The driver MUST NOT send more than *cmd\_per\_lun* linked commands to one LUN, and MUST NOT send more than the virtqueue size number of linked commands to one LUN.

#### 5.6.4.2 Device Requirements: Device configuration layout

On reset, the device MUST set *sense\_size* to 96 and *cdb\_size* to 32.

### 5.6.4.3 Legacy Interface: Device configuration layout

When using the legacy interface, transitional devices and drivers MUST format the fields in struct virtio\_scsi\_config according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

## 5.6.5 Device Requirements: Device Initialization

On initialization the driver SHOULD first discover the device's virtqueues.

If the driver uses the eventq, the driver SHOULD place at least one buffer in the eventq.

The driver MAY immediately issue requests<sup>9</sup> or task management functions<sup>10</sup>.

## 5.6.6 Device Operation

Device operation consists of operating request queues, the control queue and the event queue.

### 5.6.6.0.1 Legacy Interface: Device Operation

When using the legacy interface, the driver SHOULD ignore the used length values.

**Note:** Historically, devices put the total descriptor length, or the total length of device-writable buffers there, even when only part of the buffers were actually written.

### 5.6.6.1 Device Operation: Request Queues

The driver queues requests to an arbitrary request queue, and they are used by the device on that same queue. It is the responsibility of the driver to ensure strict request ordering for commands placed on different queues, because they will be consumed with no order constraints.

Requests have the following format:

```
struct virtio_scsi_req_cmd {
    // Device-readable part
    u8 lun[8];
    le64 id;
    u8 task_attr;
    u8 prio;
    u8 crn;
    u8 cdb[cdb_size];
    // The next three fields are only present if VIRTIO SCSI_F_T10_PI
    // is negotiated.
    le32 pi_bytesout;
    le32 pi_bytesin;
    u8 pi_out[pi_bytesout];
    u8 dataout[];

    // Device-writable part
    le32 sense_len;
    le32 residual;
    le16 status_qualifier;
    u8 status;
    u8 response;
    u8 sense[sense_size];
    // The next field is only present if VIRTIO SCSI_F_T10_PI
    // is negotiated
    u8 pi_in[pi_bytesin];
    u8 datain[];
};

/* command-specific response values */
#define VIRTIO SCSI_S_OK 0
#define VIRTIO SCSI_S_OVERRUN 1
```

<sup>9</sup>For example, INQUIRY or REPORT LUNS.

<sup>10</sup>For example, I\_T RESET.

```

#define VIRTIO SCSI_S_ABORTED          2
#define VIRTIO SCSI_S_BAD_TARGET       3
#define VIRTIO SCSI_S_RESET            4
#define VIRTIO SCSI_S_BUSY             5
#define VIRTIO SCSI_S_TRANSPORT_FAILURE 6
#define VIRTIO SCSI_S_TARGET_FAILURE   7
#define VIRTIO SCSI_S_NEXUS_FAILURE    8
#define VIRTIO SCSI_S_FAILURE          9

/* task_attr */
#define VIRTIO SCSI_S_SIMPLE            0
#define VIRTIO SCSI_S_ORDERED          1
#define VIRTIO SCSI_S_HEAD              2
#define VIRTIO SCSI_S_ACA               3

```

*lun* addresses the REPORT LUNS well-known logical unit, or a target and logical unit in the virtio-scsi device's SCSI domain. When used to address the REPORT LUNS logical unit, *lun* is 0xC1, 0x01 and six zero bytes. The virtio-scsi device SHOULD implement the REPORT LUNS well-known logical unit.

When used to address a target and logical unit, the only supported format for *lun* is: first byte set to 1, second byte set to target, third and fourth byte representing a single level LUN structure, followed by four zero bytes. With this representation, a virtio-scsi device can serve up to 256 targets and 16384 LUNs per target. The device MAY also support having a well-known logical units in the third and fourth byte.

*id* is the command identifier ("tag").

*task\_attr* defines the task attribute as in the table above, but all task attributes MAY be mapped to SIMPLE by the device. Some commands are defined by SCSI standards as "implicit head of queue"; for such commands, all task attributes MAY also be mapped to HEAD OF QUEUE. Drivers and applications SHOULD NOT send a command with the ORDERED task attribute if the command has an implicit HEAD OF QUEUE attribute, because whether the ORDERED task attribute is honored is vendor-specific.

*crn* may also be provided by clients, but is generally expected to be 0. The maximum CRN value defined by the protocol is 255, since CRN is stored in an 8-bit integer.

The CDB is included in *cdb* and its size, *cdb\_size*, is taken from the configuration space.

All of these fields are defined in [SAM](#) and are always device-readable.

*pi\_bytesout* determines the size of the *pi\_out* field in bytes. If it is nonzero, the *pi\_out* field contains outgoing protection information for write operations. *pi\_bytesin* determines the size of the *pi\_in* field in the device-writable section, in bytes. All three fields are only present if VIRTIO SCSI\_F\_T10\_PI has been negotiated.

The remainder of the device-readable part is the data output buffer, *dataout*.

*sense* and subsequent fields are always device-writable. *sense\_len* indicates the number of bytes actually written to the sense buffer.

*residual* indicates the residual size, calculated as "data\_length - number\_of\_transferred\_bytes", for read or write operations. For bidirectional commands, the number\_of\_transferred\_bytes includes both read and written bytes. A *residual* that is less than the size of *datain* means that *dataout* was processed entirely. A *residual* that exceeds the size of *datain* means that *dataout* was processed partially and *datain* was not processed at all.

If the *pi\_bytesin* is nonzero, the *pi\_in* field contains incoming protection information for read operations. *pi\_in* is only present if VIRTIO SCSI\_F\_T10\_PI has been negotiated<sup>11</sup>.

The remainder of the device-writable part is the data input buffer, *datain*.

#### 5.6.6.1.1 Device Requirements: Device Operation: Request Queues

The device MUST write the *status* byte as the status code as defined in [SAM](#).

The device MUST write the *response* byte as one of the following:

<sup>11</sup>There is no separate residual size for *pi\_bytesout* and *pi\_bytesin*. It can be computed from the *residual* field, the size of the data integrity information per sector, and the sizes of *pi\_out*, *pi\_in*, *dataout* and *datain*.

**VIRTIO\_SCSI\_S\_OK** when the request was completed and the *status* byte is filled with a SCSI status code (not necessarily “GOOD”).

**VIRTIO\_SCSI\_S\_OVERRUN** if the content of the CDB (such as the allocation length, parameter length or transfer size) requires more data than is available in the *datain* and *dataout* buffers.

**VIRTIO\_SCSI\_S\_ABORTED** if the request was cancelled due to an ABORT TASK or ABORT TASK SET task management function.

**VIRTIO\_SCSI\_S\_BAD\_TARGET** if the request was never processed because the target indicated by *lun* does not exist.

**VIRTIO\_SCSI\_S\_RESET** if the request was cancelled due to a bus or device reset (including a task management function).

**VIRTIO\_SCSI\_S\_TRANSPORT\_FAILURE** if the request failed due to a problem in the connection between the host and the target (severed link).

**VIRTIO\_SCSI\_S\_TARGET\_FAILURE** if the target is suffering a failure and to tell the driver not to retry on other paths.

**VIRTIO\_SCSI\_S\_NEXUS\_FAILURE** if the nexus is suffering a failure but retrying on other paths might yield a different result.

**VIRTIO\_SCSI\_S\_BUSY** if the request failed but retrying on the same path is likely to work.

**VIRTIO\_SCSI\_S\_FAILURE** for other host or driver error. In particular, if neither *dataout* nor *datain* is empty, and the **VIRTIO\_SCSI\_F\_INOUT** feature has not been negotiated, the request will be immediately returned with a response equal to **VIRTIO\_SCSI\_S\_FAILURE**.

All commands must be completed before the virtio-scsi device is reset or unplugged. The device MAY choose to abort them, or if it does not do so MUST pick the **VIRTIO\_SCSI\_S\_FAILURE** response.

#### 5.6.6.1.2 Driver Requirements: Device Operation: Request Queues

*task\_attr*, *prio* and *crn* SHOULD be zero.

Upon receiving a **VIRTIO\_SCSI\_S\_TARGET\_FAILURE** response, the driver SHOULD NOT retry the request on other paths.

#### 5.6.6.1.3 Legacy Interface: Device Operation: Request Queues

When using the legacy interface, transitional devices and drivers MUST format the fields in struct *virtio\_scsi\_req\_cmd* according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.6.6.2 Device Operation: controlq

The controlq is used for other SCSI transport operations. Requests have the following format:

```
struct virtio_scsi_ctrl {
    le32 type;
    ...
    u8 response;
};

/* response values valid for all commands */
#define VIRTIO_SCSI_S_OK 0
#define VIRTIO_SCSI_S_BAD_TARGET 3
#define VIRTIO_SCSI_S_BUSY 5
#define VIRTIO_SCSI_S_TRANSPORT_FAILURE 6
#define VIRTIO_SCSI_S_TARGET_FAILURE 7
#define VIRTIO_SCSI_S_NEXUS_FAILURE 8
#define VIRTIO_SCSI_S_FAILURE 9
#define VIRTIO_SCSI_S_INCORRECT_LUN 12
```

The *type* identifies the remaining fields.

The following commands are defined:

- Task management function.

```
#define VIRTIO SCSI_T_TMF 0

#define VIRTIO SCSI_T_TMF_ABORT_TASK 0
#define VIRTIO SCSI_T_TMF_ABORT_TASK_SET 1
#define VIRTIO SCSI_T_TMF_CLEAR_ACA 2
#define VIRTIO SCSI_T_TMF_CLEAR_TASK_SET 3
#define VIRTIO SCSI_T_TMF_I_T_NEXUS_RESET 4
#define VIRTIO SCSI_T_TMF_LOGICAL_UNIT_RESET 5
#define VIRTIO SCSI_T_TMF_QUERY_TASK 6
#define VIRTIO SCSI_T_TMF_QUERY_TASK_SET 7

struct virtio_scsi_ctrl_tmf {
    // Device-readable part
    le32 type;
    le32 subtype;
    u8 lun[8];
    le64 id;
    // Device-writable part
    u8 response;
};

/* command-specific response values */
#define VIRTIO SCSI_S_FUNCTION_COMPLETE 0
#define VIRTIO SCSI_S_FUNCTION_SUCCEEDED 10
#define VIRTIO SCSI_S_FUNCTION_REJECTED 11
```

The *type* is `VIRTIO SCSI_T_TMF`; *subtype* defines which task management function. All fields except *response* are filled by the driver.

Other fields which are irrelevant for the requested TMF are ignored but they are still present. *lun* is in the same format specified for request queues; the single level LUN is ignored when the task management function addresses a whole I\_T nexus. When relevant, the value of *id* is matched against the id values passed on the requestq.

The outcome of the task management function is written by the device in *response*. The command-specific response values map 1-to-1 with those defined in [SAM](#).

Task management function can affect the response value for commands that are in the request queue and have not been completed yet. For example, the device MUST complete all active commands on a logical unit or target (possibly with a `VIRTIO SCSI_S_RESET` response code) upon receiving a "logical unit reset" or "I\_T nexus reset" TMF. Similarly, the device MUST complete the selected commands (possibly with a `VIRTIO SCSI_S_ABORTED` response code) upon receiving an "abort task" or "abort task set" TMF. Such effects MUST take place before the TMF itself is successfully completed, and the device MUST use memory barriers appropriately in order to ensure that the driver sees these writes in the correct order.

- Asynchronous notification query.

```
#define VIRTIO SCSI_T_AN_QUERY 1

struct virtio_scsi_ctrl_an {
    // Device-readable part
    le32 type;
    u8 lun[8];
    le32 event_requested;
    // Device-writable part
    le32 event_actual;
    u8 response;
};

#define VIRTIO SCSI_EVT_ASYNC_OPERATIONAL_CHANGE 2
#define VIRTIO SCSI_EVT_ASYNC_POWER_MGMT 4
#define VIRTIO SCSI_EVT_ASYNC_EXTERNAL_REQUEST 8
```



```
#define VIRTIO_SCSI_EVT_ASYNC_MEDIA_CHANGE    16
#define VIRTIO_SCSI_EVT_ASYNC_MULTI_HOST     32
#define VIRTIO_SCSI_EVT_ASYNC_DEVICE_BUSY    64
```

By sending this command, the driver asks the device which events the given LUN can report, as described in paragraphs 6.6 and A.6 of [SCSI MMC](#). The driver writes the events it is interested in into *event\_requested*; the device responds by writing the events that it supports into *event\_actual*.

The *type* is VIRTIO\_SCSI\_T\_AN\_QUERY. *lun* and *event\_requested* are written by the driver. *event\_actual* and *response* fields are written by the device.

No command-specific values are defined for the *response* byte.

- Asynchronous notification subscription.

```
#define VIRTIO_SCSI_T_AN_SUBSCRIBE            2

struct virtio_scsi_ctrl_an {
    // Device-readable part
    le32 type;
    u8 lun[8];
    le32 event_requested;
    // Device-writable part
    le32 event_actual;
    u8 response;
};
```

By sending this command, the driver asks the specified LUN to report events for its physical interface, again as described in [SCSI MMC](#). The driver writes the events it is interested in into *event\_requested*; the device responds by writing the events that it supports into *event\_actual*.

Event types are the same as for the asynchronous notification query message.

The *type* is VIRTIO\_SCSI\_T\_AN\_SUBSCRIBE. *lun* and *event\_requested* are written by the driver. *event\_actual* and *response* are written by the device.

No command-specific values are defined for the response byte.

#### 5.6.6.2.1 Legacy Interface: Device Operation: controlq

When using the legacy interface, transitional devices and drivers MUST format the fields in struct *virtio\_scsi\_ctrl*, struct *virtio\_scsi\_ctrl\_tmf*, struct *virtio\_scsi\_ctrl\_an* and struct *virtio\_scsi\_ctrl\_an* according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.6.6.3 Device Operation: eventq

The eventq is populated by the driver for the device to report information on logical units that are attached to it. In general, the device will not queue events to cope with an empty eventq, and will end up dropping events if it finds no buffer ready. However, when reporting events for many LUNs (e.g. when a whole target disappears), the device can throttle events to avoid dropping them. For this reason, placing 10-15 buffers on the event queue is sufficient.

Buffers returned by the device on the eventq will be referred to as “events” in the rest of this section. Events have the following format:

```
#define VIRTIO_SCSI_T_EVENTS_MISSED    0x80000000

struct virtio_scsi_event {
    // Device-writable part
    le32 event;
    u8 lun[8];
    le32 reason;
};
```

The device sets bit 31 in *event* to report lost events due to missing buffers.

The meaning of *reason* depends on the contents of *event*. The following events are defined:

- No event.

```
#define VIRTIO_SCSI_T_NO_EVENT 0
```

This event is fired in the following cases:

- When the device detects in the eventq a buffer that is shorter than what is indicated in the configuration field, it MAY use it immediately and put this dummy value in *event*. A well-written driver will never observe this situation.
- When events are dropped, the device MAY signal this event as soon as the drivers makes a buffer available, in order to request action from the driver. In this case, of course, this event will be reported with the VIRTIO\_SCSI\_T\_EVENTS\_MISSED flag.

- Transport reset

```
#define VIRTIO_SCSI_T_TRANSPORT_RESET 1
#define VIRTIO_SCSI_EVT_RESET_HARD 0
#define VIRTIO_SCSI_EVT_RESET_RESCAN 1
#define VIRTIO_SCSI_EVT_RESET_REMOVED 2
```

By sending this event, the device signals that a logical unit on a target has been reset, including the case of a new device appearing or disappearing on the bus. The device fills in all fields. *event* is set to VIRTIO\_SCSI\_T\_TRANSPORT\_RESET. *lun* addresses a logical unit in the SCSI host.

The *reason* value is one of the three #define values appearing above:

**VIRTIO\_SCSI\_EVT\_RESET\_REMOVED** (“LUN/target removed”) is used if the target or logical unit is no longer able to receive commands.

**VIRTIO\_SCSI\_EVT\_RESET\_HARD** (“LUN hard reset”) is used if the logical unit has been reset, but is still present.

**VIRTIO\_SCSI\_EVT\_RESET\_RESCAN** (“rescan LUN/target”) is used if a target or logical unit has just appeared on the device.

The “removed” and “rescan” events can happen when VIRTIO\_SCSI\_F\_HOTPLUG feature was negotiated; when sent for LUN 0, they MAY apply to the entire target so the driver can ask the initiator to rescan the target to detect this.

Events will also be reported via sense codes (this obviously does not apply to newly appeared buses or targets, since the application has never discovered them):

- “LUN/target removed” maps to sense key ILLEGAL REQUEST, asc 0x25, ascq 0x00 (LOGICAL UNIT NOT SUPPORTED)
- “LUN hard reset” maps to sense key UNIT ATTENTION, asc 0x29 (POWER ON, RESET OR BUS DEVICE RESET OCCURRED)
- “rescan LUN/target” maps to sense key UNIT ATTENTION, asc 0x3f, ascq 0x0e (REPORTED LUNS DATA HAS CHANGED)

The preferred way to detect transport reset is always to use events, because sense codes are only seen by the driver when it sends a SCSI command to the logical unit or target. However, in case events are dropped, the initiator will still be able to synchronize with the actual state of the controller if the driver asks the initiator to rescan of the SCSI bus. During the rescan, the initiator will be able to observe the above sense codes, and it will process them as if it the driver had received the equivalent event.

- Asynchronous notification

```
#define VIRTIO_SCSI_T_ASYNC_NOTIFY 2
```

By sending this event, the device signals that an asynchronous event was fired from a physical interface.

All fields are written by the device. *event* is set to `VIRTIO_SCSI_T_ASYNC_NOTIFY`. *lun* addresses a logical unit in the SCSI host. *reason* is a subset of the events that the driver has subscribed to via the “Asynchronous notification subscription” command.

- LUN parameter change

```
#define VIRTIO_SCSI_T_PARAM_CHANGE 3
```

By sending this event, the device signals a change in the configuration parameters of a logical unit, for example the capacity or cache mode. *event* is set to `VIRTIO_SCSI_T_PARAM_CHANGE`. *lun* addresses a logical unit in the SCSI host.

The same event SHOULD also be reported as a unit attention condition. *reason* contains the additional sense code and additional sense code qualifier, respectively in bits 0...7 and 8...15.

**Note:** For example, a change in capacity will be reported as asc 0x2a, ascq 0x09 (CAPACITY DATA HAS CHANGED).

For MMC devices (inquiry type 5) there would be some overlap between this event and the asynchronous notification event, so for simplicity the host never reports this event for MMC devices.

#### 5.6.6.3.1 Driver Requirements: Device Operation: eventq

The driver SHOULD keep the eventq populated with buffers. These buffers MUST be device-writable, and SHOULD be at least *event\_info\_size* bytes long, and MUST be at least the size of struct `virtio_scsi_event`.

If *event* has bit 31 set, the driver SHOULD poll the logical units for unit attention conditions, and/or do whatever form of bus scan is appropriate for the guest operating system and SHOULD poll for asynchronous events manually using SCSI commands.

When receiving a `VIRTIO_SCSI_T_TRANSPORT_RESET` message with *reason* set to `VIRTIO_SCSI_EVT_RESET_REMOVED` or `VIRTIO_SCSI_EVT_RESET_RESCAN` for LUN 0, the driver SHOULD ask the initiator to rescan the target, in order to detect the case when an entire target has appeared or disappeared.

#### 5.6.6.3.2 Device Requirements: Device Operation: eventq

The device MUST set bit 31 in *event* if events were lost due to missing buffers, and it MAY use a `VIRTIO_SCSI_T_NO_EVENT` event to report this.

The device MUST NOT send `VIRTIO_SCSI_T_TRANSPORT_RESET` messages with *reason* set to `VIRTIO_SCSI_EVT_RESET_REMOVED` or `VIRTIO_SCSI_EVT_RESET_RESCAN` unless `VIRTIO_SCSI_F_HOTPLUG` was negotiated.

The device MUST NOT report `VIRTIO_SCSI_T_PARAM_CHANGE` for MMC devices.

#### 5.6.6.3.3 Legacy Interface: Device Operation: eventq

When using the legacy interface, transitional devices and drivers MUST format the fields in struct `virtio_scsi_event` according to the native endian of the guest rather than (necessarily when not using the legacy interface) little-endian.

#### 5.6.6.4 Legacy Interface: Framing Requirements

When using legacy interfaces, transitional drivers which have not negotiated `VIRTIO_F_ANY_LAYOUT` MUST use a single descriptor for the *lun*, *id*, *task\_attr*, *prio*, *crn* and *cdb* fields, and MUST only use a single descriptor for the *sense\_len*, *residual*, *status\_qualifier*, *status*, *response* and *sense* fields.

## 5.7 GPU Device

virtio-gpu is a virtio based graphics adapter. It can operate in 2D mode and in 3D mode. 3D mode will offload rendering ops to the host gpu and therefore requires a gpu with 3D support on the host machine.

In 2D mode the virtio-gpu device provides support for ARGB Hardware cursors and multiple scanouts (aka heads).

### 5.7.1 Device ID

16

### 5.7.2 Virtqueues

**0** controlq - queue for sending control commands

**1** cursorq - queue for sending cursor updates

Both queues have the same format. Each request and each response have a fixed header, followed by command specific data fields. The separate cursor queue is the "fast track" for cursor commands (VIRTIO\_GPU\_CMD\_UPDATE\_CURSOR and VIRTIO\_GPU\_CMD\_MOVE\_CURSOR), so they go through without being delayed by time-consuming commands in the control queue.

### 5.7.3 Feature bits

**VIRTIO\_GPU\_F\_VIRGL (0)** virgl 3D mode is supported.

**VIRTIO\_GPU\_F\_EDID (1)** EDID is supported.

**VIRTIO\_GPU\_F\_RESOURCE\_UUID (2)** assigning resources UUIDs for export to other virtio devices is supported.

**VIRTIO\_GPU\_F\_RESOURCE\_BLOB (3)** creating and using size-based blob resources is supported.

**VIRTIO\_GPU\_F\_CONTEXT\_INIT (4)** multiple context types and synchronization timelines supported. Requires VIRTIO\_GPU\_F\_VIRGL.

### 5.7.4 Device configuration layout

GPU device configuration uses the following layout structure and definitions:

```
#define VIRTIO_GPU_EVENT_DISPLAY (1 << 0)

struct virtio_gpu_config {
    le32 events_read;
    le32 events_clear;
    le32 num_scanouts;
    le32 num_capsets;
};
```

#### 5.7.4.1 Device configuration fields

**events\_read** signals pending events to the driver. The driver MUST NOT write to this field.

**events\_clear** clears pending events in the device. Writing a '1' into a bit will clear the corresponding bit in **events\_read**, mimicking write-to-clear behavior.

**num\_scanouts** specifies the maximum number of scanouts supported by the device. Minimum value is 1, maximum value is 16.

**num\_capsets** specifies the maximum number of capability sets supported by the device. The minimum value is zero.

### 5.7.4.2 Events

**VIRTIO\_GPU\_EVENT\_DISPLAY** Display configuration has changed. The driver SHOULD use the `VIRTIO_GPU_CMD_GET_DISPLAY_INFO` command to fetch the information from the device. In case EDID support is negotiated (`VIRTIO_GPU_F_EDID` feature flag) the device SHOULD also fetch the updated EDID blobs using the `VIRTIO_GPU_CMD_GET_EDID` command.

### 5.7.5 Device Requirements: Device Initialization

The driver SHOULD query the display information from the device using the `VIRTIO_GPU_CMD_GET_DISPLAY_INFO` command and use that information for the initial scanout setup. In case EDID support is negotiated (`VIRTIO_GPU_F_EDID` feature flag) the device SHOULD also fetch the EDID information using the `VIRTIO_GPU_CMD_GET_EDID` command. If no information is available or all displays are disabled the driver MAY choose to use a fallback, such as 1024x768 at display 0.

The driver SHOULD query all shared memory regions supported by the device. If the device supports shared memory, the *shm*id of a region MUST (see [2.10 Shared Memory Regions](#)) be one of the following:

```
enum virtio_gpu_shm_id {
    VIRTIO_GPU_SHM_ID_UNDEFINED = 0,
    VIRTIO_GPU_SHM_ID_HOST_VISIBLE = 1,
};
```

The shared memory region with `VIRTIO_GPU_SHM_ID_HOST_VISIBLE` is referred as the "host visible memory region". The device MUST support the `VIRTIO_GPU_CMD_RESOURCE_MAP_BLOB` and `VIRTIO_GPU_CMD_RESOURCE_UNMAP_BLOB` if the host visible memory region is available.

### 5.7.6 Device Operation

The virtio-gpu is based around the concept of resources private to the host. The guest must DMA transfer into these resources, unless shared memory regions are supported. This is a design requirement in order to interface with future 3D rendering. In the unaccelerated 2D mode there is no support for DMA transfers from resources, just to them.

Resources are initially simple 2D resources, consisting of a width, height and format along with an identifier. The guest must then attach backing store to the resources in order for DMA transfers to work. This is like a GART in a real GPU.

#### 5.7.6.1 Device Operation: Create a framebuffer and configure scanout

- Create a host resource using `VIRTIO_GPU_CMD_RESOURCE_CREATE_2D`.
- Allocate a framebuffer from guest ram, and attach it as backing storage to the resource just created, using `VIRTIO_GPU_CMD_RESOURCE_ATTACH_BACKING`. Scatter lists are supported, so the framebuffer doesn't need to be contiguous in guest physical memory.
- Use `VIRTIO_GPU_CMD_SET_SCANOUT` to link the framebuffer to a display scanout.

#### 5.7.6.2 Device Operation: Update a framebuffer and scanout

- Render to your framebuffer memory.
- Use `VIRTIO_GPU_CMD_TRANSFER_TO_HOST_2D` to update the host resource from guest memory.
- Use `VIRTIO_GPU_CMD_RESOURCE_FLUSH` to flush the updated resource to the display.

#### 5.7.6.3 Device Operation: Using pageflip

It is possible to create multiple framebuffers, flip between them using `VIRTIO_GPU_CMD_SET_SCANOUT` and `VIRTIO_GPU_CMD_RESOURCE_FLUSH`, and update the invisible framebuffer using `VIRTIO_GPU_CMD_TRANSFER_TO_HOST_2D`.

#### 5.7.6.4 Device Operation: Multihead setup

In case two or more displays are present there are different ways to configure things:

- Create a single framebuffer, link it to all displays (mirroring).
- Create an framebuffer for each display.
- Create one big framebuffer, configure scanouts to display a different rectangle of that framebuffer each.

#### 5.7.6.5 Device Requirements: Device Operation: Command lifecycle and fencing

The device MAY process controlq commands asynchronously and return them to the driver before the processing is complete. If the driver needs to know when the processing is finished it can set the VIRTIO\_GPU\_FLAG\_FENCE flag in the request. The device MUST finish the processing before returning the command then.

Note: current qemu implementation does asynchronous processing only in 3d mode, when offloading the processing to the host gpu.

#### 5.7.6.6 Device Operation: Configure mouse cursor

The mouse cursor image is a normal resource, except that it must be 64x64 in size. The driver MUST create and populate the resource (using the usual VIRTIO\_GPU\_CMD\_RESOURCE\_CREATE\_2D, VIRTIO\_GPU\_CMD\_RESOURCE\_ATTACH\_BACKING and VIRTIO\_GPU\_CMD\_TRANSFER\_TO\_HOST\_2D controlq commands) and make sure they are completed (using VIRTIO\_GPU\_FLAG\_FENCE).

Then VIRTIO\_GPU\_CMD\_UPDATE\_CURSOR can be sent to the cursorq to set the pointer shape and position. To move the pointer without updating the shape use VIRTIO\_GPU\_CMD\_MOVE\_CURSOR instead.

#### 5.7.6.7 Device Operation: Request header

All requests and responses on the virt queues have a fixed header using the following layout structure and definitions:

```
enum virtio_gpu_ctrl_type {  
  
    /* 2d commands */  
    VIRTIO_GPU_CMD_GET_DISPLAY_INFO = 0x0100,  
    VIRTIO_GPU_CMD_RESOURCE_CREATE_2D,  
    VIRTIO_GPU_CMD_RESOURCE_UNREF,  
    VIRTIO_GPU_CMD_SET_SCANOUT,  
    VIRTIO_GPU_CMD_RESOURCE_FLUSH,  
    VIRTIO_GPU_CMD_TRANSFER_TO_HOST_2D,  
    VIRTIO_GPU_CMD_RESOURCE_ATTACH_BACKING,  
    VIRTIO_GPU_CMD_RESOURCE_DETACH_BACKING,  
    VIRTIO_GPU_CMD_GET_CAPSET_INFO,  
    VIRTIO_GPU_CMD_GET_CAPSET,  
    VIRTIO_GPU_CMD_GET_EDID,  
    VIRTIO_GPU_CMD_RESOURCE_ASSIGN_UUID,  
    VIRTIO_GPU_CMD_RESOURCE_CREATE_BLOB,  
    VIRTIO_GPU_CMD_SET_SCANOUT_BLOB,  
  
    /* 3d commands */  
    VIRTIO_GPU_CMD_CTX_CREATE = 0x0200,  
    VIRTIO_GPU_CMD_CTX_DESTROY,  
    VIRTIO_GPU_CMD_CTX_ATTACH_RESOURCE,  
    VIRTIO_GPU_CMD_CTX_DETACH_RESOURCE,  
    VIRTIO_GPU_CMD_RESOURCE_CREATE_3D,  
    VIRTIO_GPU_CMD_TRANSFER_TO_HOST_3D,  
    VIRTIO_GPU_CMD_TRANSFER_FROM_HOST_3D,  
    VIRTIO_GPU_CMD_SUBMIT_3D,  
    VIRTIO_GPU_CMD_RESOURCE_MAP_BLOB,  
    VIRTIO_GPU_CMD_RESOURCE_UNMAP_BLOB,  
  
    /* cursor commands */  
    VIRTIO_GPU_CMD_UPDATE_CURSOR = 0x0300,  
    VIRTIO_GPU_CMD_MOVE_CURSOR,  
  
    /* success responses */  
}
```

```

    VIRTIO_GPU_RESP_OK_NODATA = 0x1100,
    VIRTIO_GPU_RESP_OK_DISPLAY_INFO,
    VIRTIO_GPU_RESP_OK_CAPSET_INFO,
    VIRTIO_GPU_RESP_OK_CAPSET,
    VIRTIO_GPU_RESP_OK_EDID,
    VIRTIO_GPU_RESP_OK_RESOURCE_UUID,
    VIRTIO_GPU_RESP_OK_MAP_INFO,

    /* error responses */
    VIRTIO_GPU_RESP_ERR_UNSPEC = 0x1200,
    VIRTIO_GPU_RESP_ERR_OUT_OF_MEMORY,
    VIRTIO_GPU_RESP_ERR_INVALID_SCANOUT_ID,
    VIRTIO_GPU_RESP_ERR_INVALID_RESOURCE_ID,
    VIRTIO_GPU_RESP_ERR_INVALID_CONTEXT_ID,
    VIRTIO_GPU_RESP_ERR_INVALID_PARAMETER,
};

#define VIRTIO_GPU_FLAG_FENCE (1 << 0)
#define VIRTIO_GPU_FLAG_INFO_RING_IDX (1 << 1)

struct virtio_gpu_ctrl_hdr {
    le32 type;
    le32 flags;
    le64 fence_id;
    le32 ctx_id;
    u8 ring_idx;
    u8 padding[3];
};

```

The fixed header *struct virtio\_gpu\_ctrl\_hdr* in each request includes the following fields:

**type** specifies the type of the driver request (VIRTIO\_GPU\_CMD\_\*) or device response (VIRTIO\_GPU\_RESP\_\*).

**flags** request / response flags.

**fence\_id** If the driver sets the VIRTIO\_GPU\_FLAG\_FENCE bit in the request *flags* field the device MUST:

- set VIRTIO\_GPU\_FLAG\_FENCE bit in the response,
- copy the content of the *fence\_id* field from the request to the response, and
- send the response only after command processing is complete.

**ctx\_id** Rendering context (used in 3D mode only).

**ring\_idx** If VIRTIO\_GPU\_F\_CONTEXT\_INIT is supported, then the driver MAY set VIRTIO\_GPU\_FLAG\_INFO\_RING\_IDX bit in the request *flags*. In that case:

- *ring\_idx* indicates the value of a context-specific ring index. The minimum value is 0 and maximum value is 63 (inclusive).
- If VIRTIO\_GPU\_FLAG\_FENCE is set, *fence\_id* acts as a sequence number on the synchronization timeline defined by *ctx\_id* and the ring index.
- If VIRTIO\_GPU\_FLAG\_FENCE is set and when the command associated with *fence\_id* is complete, the device MUST send a response for all outstanding commands with a sequence number less than or equal to *fence\_id* on the same synchronization timeline.

On success the device will return VIRTIO\_GPU\_RESP\_OK\_NODATA in case there is no payload. Otherwise the *type* field will indicate the kind of payload.

On error the device will return one of the VIRTIO\_GPU\_RESP\_ERR\_\* error codes.

#### 5.7.6.8 Device Operation: controlq

For any coordinates given 0,0 is top left, larger x moves right, larger y moves down.

**VIRTIO\_GPU\_CMD\_GET\_DISPLAY\_INFO** Retrieve the current output configuration. No request data (just bare *struct virtio\_gpu\_ctrl\_hdr*). Response type is VIRTIO\_GPU\_RESP\_OK\_DISPLAY\_INFO, response data is *struct virtio\_gpu\_resp\_display\_info*.



```

#define VIRTIO_GPU_MAX_SCANOUTS 16

struct virtio_gpu_rect {
    le32 x;
    le32 y;
    le32 width;
    le32 height;
};

struct virtio_gpu_resp_display_info {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_display_one {
        struct virtio_gpu_rect r;
        le32 enabled;
        le32 flags;
    } pmodes[VIRTIO_GPU_MAX_SCANOUTS];
};

```

The response contains a list of per-scanout information. The info contains whether the scanout is enabled and what its preferred position and size is.

The size (fields *width* and *height*) is similar to the native panel resolution in EDID display information, except that in the virtual machine case the size can change when the host window representing the guest display is gets resized.

The position (fields *x* and *y*) describe how the displays are arranged (i.e. which is – for example – the left display).

The *enabled* field is set when the user enabled the display. It is roughly the same as the connected state of a physical display connector.

**VIRTIO\_GPU\_CMD\_GET\_EDID** Retrieve the EDID data for a given scanout. Request data is *struct virtio\_gpu\_get\_edid*. Response type is VIRTIO\_GPU\_RESP\_OK\_EDID, response data is *struct virtio\_gpu\_resp\_edid*. Support is optional and negotiated using the VIRTIO\_GPU\_F\_EDID feature flag.

```

struct virtio_gpu_get_edid {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 scanout;
    le32 padding;
};

struct virtio_gpu_resp_edid {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 size;
    le32 padding;
    u8 edid[1024];
};

```

The response contains the EDID display data blob (as specified by VESA) for the scanout.

**VIRTIO\_GPU\_CMD\_RESOURCE\_CREATE\_2D** Create a 2D resource on the host. Request data is *struct virtio\_gpu\_resource\_create\_2d*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

enum virtio_gpu_formats {
    VIRTIO_GPU_FORMAT_B8G8R8A8_UNORM = 1,
    VIRTIO_GPU_FORMAT_B8G8R8X8_UNORM = 2,
    VIRTIO_GPU_FORMAT_A8R8G8B8_UNORM = 3,
    VIRTIO_GPU_FORMAT_X8R8G8B8_UNORM = 4,

    VIRTIO_GPU_FORMAT_R8G8B8A8_UNORM = 67,
    VIRTIO_GPU_FORMAT_X8B8G8R8_UNORM = 68,

    VIRTIO_GPU_FORMAT_A8B8G8R8_UNORM = 121,
    VIRTIO_GPU_FORMAT_R8G8B8X8_UNORM = 134,
};

struct virtio_gpu_resource_create_2d {
    struct virtio_gpu_ctrl_hdr hdr;

```

```

    le32 resource_id;
    le32 format;
    le32 width;
    le32 height;
};

```

This creates a 2D resource on the host with the specified width, height and format. The resource ids are generated by the guest.

**VIRTIO\_GPU\_CMD\_RESOURCE\_UNREF** Destroy a resource. Request data is *struct virtio\_gpu\_resource\_unref*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

struct virtio_gpu_resource_unref {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 padding;
};

```

This informs the host that a resource is no longer required by the guest.

**VIRTIO\_GPU\_CMD\_SET\_SCANOUT** Set the scanout parameters for a single output. Request data is *struct virtio\_gpu\_set\_scanout*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

struct virtio_gpu_set_scanout {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_rect r;
    le32 scanout_id;
    le32 resource_id;
};

```

This sets the scanout parameters for a single scanout. The resource\_id is the resource to be scanned out from, along with a rectangle.

Scanout rectangles must be completely covered by the underlying resource. Overlapping (or identical) scanouts are allowed, typical use case is screen mirroring.

The driver can use resource\_id = 0 to disable a scanout.

**VIRTIO\_GPU\_CMD\_RESOURCE\_FLUSH** Flush a scanout resource Request data is *struct virtio\_gpu\_resource\_flush*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

struct virtio_gpu_resource_flush {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_rect r;
    le32 resource_id;
    le32 padding;
};

```

This flushes a resource to screen. It takes a rectangle and a resource id, and flushes any scanouts the resource is being used on.

**VIRTIO\_GPU\_CMD\_TRANSFER\_TO\_HOST\_2D** Transfer from guest memory to host resource. Request data is *struct virtio\_gpu\_transfer\_to\_host\_2d*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

struct virtio_gpu_transfer_to_host_2d {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_rect r;
    le64 offset;
    le32 resource_id;
    le32 padding;
};

```

This takes a resource id along with an destination offset into the resource, and a box to transfer to the host backing for the resource.

**VIRTIO\_GPU\_CMD\_RESOURCE\_ATTACH\_BACKING** Assign backing pages to a resource. Request data is *struct virtio\_gpu\_resource\_attach\_backing*, followed by *struct virtio\_gpu\_mem\_entry* entries. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

```

struct virtio_gpu_resource_attach_backing {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 nr_entries;
};

struct virtio_gpu_mem_entry {
    le64 addr;
    le32 length;
    le32 padding;
};

```

This assigns an array of guest pages as the backing store for a resource. These pages are then used for the transfer operations for that resource from that point on.

**VIRTIO\_GPU\_CMD\_RESOURCE\_DETACH\_BACKING** Detach backing pages from a resource. Request data is *struct virtio\_gpu\_resource\_detach\_backing*. Response type is VIRTIO\_GPU\_RESP\_OK\_NO\_DATA.

```

struct virtio_gpu_resource_detach_backing {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 padding;
};

```

This detaches any backing pages from a resource, to be used in case of guest swapping or object destruction.

**VIRTIO\_GPU\_CMD\_GET\_CAPSET\_INFO** Gets the information associated with a particular *capset\_index*, which MUST be less than *num\_capsets* defined in the device configuration. Request data is *struct virtio\_gpu\_get\_capset\_info*. Response type is VIRTIO\_GPU\_RESP\_OK\_CAPSET\_INFO.

On success, *struct virtio\_gpu\_resp\_capset\_info* contains the *capset\_id*, *capset\_max\_version*, *capset\_max\_size* associated with capset at the specified *capset\_idx*. *capset\_id* MUST be one of the following (see listing for values):

- [VIRTIO\\_GPU\\_CAPSET\\_VIRGL](#) – the first edition of Virgl (Gallium OpenGL) protocol.
- [VIRTIO\\_GPU\\_CAPSET\\_VIRGL2](#) – the second edition of Virgl (Gallium OpenGL) protocol after the capset fix.
- [VIRTIO\\_GPU\\_CAPSET\\_GFXSTREAM](#) – gfxstream's (mostly) autogenerated GLES and Vulkan streaming protocols.
- [VIRTIO\\_GPU\\_CAPSET\\_VENUS](#) – Mesa's (mostly) autogenerated Vulkan protocol.
- [VIRTIO\\_GPU\\_CAPSET\\_CROSS\\_DOMAIN](#) – protocol for display virtualization via Wayland proxying.

```

struct virtio_gpu_get_capset_info {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 capset_index;
    le32 padding;
};

#define VIRTIO_GPU_CAPSET_VIRGL 1
#define VIRTIO_GPU_CAPSET_VIRGL2 2
#define VIRTIO_GPU_CAPSET_GFXSTREAM 3
#define VIRTIO_GPU_CAPSET_VENUS 4
#define VIRTIO_GPU_CAPSET_CROSS_DOMAIN 5
struct virtio_gpu_resp_capset_info {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 capset_id;
    le32 capset_max_version;
    le32 capset_max_size;
    le32 padding;
};

```

**VIRTIO\_GPU\_CMD\_GET\_CAPSET** Gets the capset associated with a particular *capset\_id* and *capset\_version*. Request data is *struct virtio\_gpu\_get\_capset*. Response type is VIRTIO\_GPU\_RESP\_OK\_CAPSET.

```
struct virtio_gpu_get_capset {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 capset_id;
    le32 capset_version;
};

struct virtio_gpu_resp_capset {
    struct virtio_gpu_ctrl_hdr hdr;
    u8 capset_data[];
};
```

**VIRTIO\_GPU\_CMD\_RESOURCE\_ASSIGN\_UUID** Creates an exported object from a resource. Request data is *struct virtio\_gpu\_resource\_assign\_uuid*. Response type is VIRTIO\_GPU\_RESP\_OK\_RESOURCE\_UUID, response data is *struct virtio\_gpu\_resp\_resource\_uuid*. Support is optional and negotiated using the VIRTIO\_GPU\_F\_RESOURCE\_UUID feature flag.

```
struct virtio_gpu_resource_assign_uuid {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 padding;
};

struct virtio_gpu_resp_resource_uuid {
    struct virtio_gpu_ctrl_hdr hdr;
    u8 uuid[16];
};
```

The response contains a UUID which identifies the exported object created from the host private resource. Note that if the resource has an attached backing, modifications made to the host private resource through the exported object by other devices are not visible in the attached backing until they are transferred into the backing.

**VIRTIO\_GPU\_CMD\_RESOURCE\_CREATE\_BLOB** Creates a virtio-gpu blob resource. Request data is *struct virtio\_gpu\_resource\_create\_blob*, followed by *struct virtio\_gpu\_mem\_entry* entries. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA. Support is optional and negotiated using the VIRTIO\_GPU\_F\_RESOURCE\_BLOB feature flag.

```
#define VIRTIO_GPU_BLOB_MEM_GUEST      0x0001
#define VIRTIO_GPU_BLOB_MEM_HOST3D     0x0002
#define VIRTIO_GPU_BLOB_MEM_HOST3D_GUEST 0x0003

#define VIRTIO_GPU_BLOB_FLAG_USE_MAPPABLE 0x0001
#define VIRTIO_GPU_BLOB_FLAG_USE_SHAREABLE 0x0002
#define VIRTIO_GPU_BLOB_FLAG_USE_CROSS_DEVICE 0x0004

struct virtio_gpu_resource_create_blob {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 blob_mem;
    le32 blob_flags;
    le32 nr_entries;
    le64 blob_id;
    le64 size;
};
```

A blob resource is a container for:

- a guest memory allocation (referred to as a "guest-only blob resource").
- a host memory allocation (referred to as a "host-only blob resource").
- a guest memory and host memory allocation (referred to as a "default blob resource").

The memory properties of the blob resource MUST be described by *blob\_mem*, which MUST be non-zero.

For default and guest-only blob resources, *nr\_entries* guest memory entries may be assigned to the resource. For default blob resources (i.e, when *blob\_mem* is `VIRTIO_GPU_BLOB_MEM_HOST3D_GUEST`), these memory entries are used as a shadow buffer for the host memory. To facilitate drivers that support swap-in and swap-out, *nr\_entries* may be zero and `VIRTIO_GPU_CMD_RESOURCE_ATTACH_BACKING` may be subsequently used. `VIRTIO_GPU_CMD_RESOURCE_DETACH_BACKING` may be used to unassign memory entries.

*blob\_mem* can only be `VIRTIO_GPU_BLOB_MEM_HOST3D` and `VIRTIO_GPU_BLOB_MEM_HOST3D_GUEST` if `VIRTIO_GPU_F_VIRGL` is supported. `VIRTIO_GPU_BLOB_MEM_GUEST` is valid regardless whether `VIRTIO_GPU_F_VIRGL` is supported or not.

For `VIRTIO_GPU_BLOB_MEM_HOST3D` and `VIRTIO_GPU_BLOB_MEM_HOST3D_GUEST`, the virtio-gpu resource MUST be created from the rendering context local object identified by the *blob\_id*. The actual allocation is done via `VIRTIO_GPU_CMD_SUBMIT_3D`.

The driver MUST inform the device if the blob resource is used for memory access, sharing between driver instances and/or sharing with other devices. This is done via the *blob\_flags* field.

If `VIRTIO_GPU_F_VIRGL` is set, both `VIRTIO_GPU_CMD_TRANSFER_TO_HOST_3D` and `VIRTIO_GPU_CMD_TRANSFER_FROM_HOST_3D` may be used to update the resource. There is no restriction on the image/buffer view the driver has on the blob resource.

**VIRTIO\_GPU\_CMD\_SET\_SCANOUT\_BLOB** sets scanout parameters for a blob resource. Request data is *struct virtio\_gpu\_set\_scanout\_blob*. Response type is `VIRTIO_GPU_RESP_OK_NODATA`. Support is optional and negotiated using the `VIRTIO_GPU_F_RESOURCE_BLOB` feature flag.

```
struct virtio_gpu_set_scanout_blob {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_rect r;
    le32 scanout_id;
    le32 resource_id;
    le32 width;
    le32 height;
    le32 format;
    le32 padding;
    le32 strides[4];
    le32 offsets[4];
};
```

The rectangle *r* represents the portion of the blob resource being displayed. The rest is the metadata associated with the blob resource. The format MUST be one of *enum virtio\_gpu\_formats*. The format MAY be compressed with header and data planes.

#### 5.7.6.9 Device Operation: controlq (3d)

These commands are supported by the device if the `VIRTIO_GPU_F_VIRGL` feature flag is set.

**VIRTIO\_GPU\_CMD\_CTX\_CREATE** creates a context for submitting an opaque command stream. Request data is *struct virtio\_gpu\_ctx\_create*. Response type is `VIRTIO_GPU_RESP_OK_NODATA`.

```
#define VIRTIO_GPU_CONTEXT_INIT_CAPSET_ID_MASK 0x000000ff;
struct virtio_gpu_ctx_create {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 nlen;
    le32 context_init;
    char debug_name[64];
};
```

The implementation MUST create a context for the given *ctx\_id* in the *hdr*. For debugging purposes, a *debug\_name* and its length *nlen* is provided by the driver. If `VIRTIO_GPU_F_CONTEXT_INIT` is supported, then lower 8 bits of *context\_init* MAY contain the *capset\_id* associated with context. In that case, then the device MUST create a context that can handle the specified command stream.

If the lower 8-bits of the *context\_init* are zero, then the type of the context is determined by the device.

**VIRTIO\_GPU\_CMD\_CTX\_DESTROY**

**VIRTIO\_GPU\_CMD\_CTX\_ATTACH\_RESOURCE**

**VIRTIO\_GPU\_CMD\_CTX\_DETACH\_RESOURCE** Manage virtio-gpu 3d contexts.

**VIRTIO\_GPU\_CMD\_RESOURCE\_CREATE\_3D** Create virtio-gpu 3d resources.

**VIRTIO\_GPU\_CMD\_TRANSFER\_TO\_HOST\_3D**

**VIRTIO\_GPU\_CMD\_TRANSFER\_FROM\_HOST\_3D** Transfer data from and to virtio-gpu 3d resources.

**VIRTIO\_GPU\_CMD\_SUBMIT\_3D** Submit an opaque command stream. The type of the command stream is determined when creating a context.

**VIRTIO\_GPU\_CMD\_RESOURCE\_MAP\_BLOB** maps a host-only blob resource into an offset in the host visible memory region. Request data is *struct virtio\_gpu\_resource\_map\_blob*. The driver MUST not map a blob resource that is already mapped. Response type is VIRTIO\_GPU\_RESP\_OK\_MAP\_INFO. Support is optional and negotiated using the VIRTIO\_GPU\_F\_RESOURCE\_BLOB feature flag and checking for the presence of the host visible memory region.

```
struct virtio_gpu_resource_map_blob {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 padding;
    le64 offset;
};

#define VIRTIO_GPU_MAP_CACHE_MASK      0x0f
#define VIRTIO_GPU_MAP_CACHE_NONE     0x00
#define VIRTIO_GPU_MAP_CACHE_CACHED   0x01
#define VIRTIO_GPU_MAP_CACHE_UNCACHED 0x02
#define VIRTIO_GPU_MAP_CACHE_WC       0x03
struct virtio_gpu_resp_map_info {
    struct virtio_gpu_ctrl_hdr hdr;
    u32 map_info;
    u32 padding;
};
```

**VIRTIO\_GPU\_CMD\_RESOURCE\_UNMAP\_BLOB** unmaps a host-only blob resource from the host visible memory region. Request data is *struct virtio\_gpu\_resource\_unmap\_blob*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA. Support is optional and negotiated using the VIRTIO\_GPU\_F\_RESOURCE\_BLOB feature flag and checking for the presence of the host visible memory region.

```
struct virtio_gpu_resource_unmap_blob {
    struct virtio_gpu_ctrl_hdr hdr;
    le32 resource_id;
    le32 padding;
};
```

#### 5.7.6.10 Device Operation: cursorq

Both cursorq commands use the same command struct.

```
struct virtio_gpu_cursor_pos {
    le32 scanout_id;
    le32 x;
    le32 y;
    le32 padding;
};

struct virtio_gpu_update_cursor {
    struct virtio_gpu_ctrl_hdr hdr;
    struct virtio_gpu_cursor_pos pos;
    le32 resource_id;
    le32 hot_x;
    le32 hot_y;
    le32 padding;
};
```

**VIRTIO\_GPU\_CMD\_UPDATE\_CURSOR** Update cursor. Request data is *struct virtio\_gpu\_update\_cursor*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

Full cursor update. Cursor will be loaded from the specified *resource\_id* and will be moved to *pos*. The driver must transfer the cursor into the resource beforehand (using control queue commands) and make sure the commands to fill the resource are actually processed (using fencing).

**VIRTIO\_GPU\_CMD\_MOVE\_CURSOR** Move cursor. Request data is *struct virtio\_gpu\_update\_cursor*. Response type is VIRTIO\_GPU\_RESP\_OK\_NODATA.

Move cursor to the place specified in *pos*. The other fields are not used and will be ignored by the device.

### 5.7.7 VGA Compatibility

Applies to Virtio Over PCI only. The GPU device can come with and without VGA compatibility. The PCI class should be DISPLAY\_VGA if VGA compatibility is present and DISPLAY\_OTHER otherwise.

VGA compatibility: PCI region 0 has the linear framebuffer, standard vga registers are present. Configuring a scanout (VIRTIO\_GPU\_CMD\_SET\_SCANOUT) switches the device from vga compatibility mode into native virtio mode. A reset switches it back into vga compatibility mode.

Note: qemu implementation also provides bochs dispi interface io ports and mmio bar at pci region 1 and is therefore fully compatible with the qemu stdvga (see [docs/specs/standard-vga.txt](https://docs.specs.standard-vga.txt) in the qemu source tree).

## 5.8 Input Device

The virtio input device can be used to create virtual human interface devices such as keyboards, mice and tablets. An instance of the virtio device represents one such input device. Device behavior mirrors that of the evdev layer in Linux, making pass-through implementations on top of evdev easy.

This specification defines how evdev events are transported over virtio and how the set of supported events is discovered by a driver. It does not, however, define the semantics of input events as this is dependent on the particular evdev implementation. For the list of events used by Linux input devices, see [include/uapi/linux/input-event-codes.h](https://include/uapi/linux/input-event-codes.h) in the Linux source tree.

### 5.8.1 Device ID

18

### 5.8.2 Virtqueues

0 eventq

1 statusq

### 5.8.3 Feature bits

None.

### 5.8.4 Device configuration layout

Device configuration holds all information the guest needs to handle the device, most importantly the events which are supported.

```
enum virtio_input_config_select {
    VIRTIO_INPUT_CFG_UNSET      = 0x00,
    VIRTIO_INPUT_CFG_ID_NAME    = 0x01,
    VIRTIO_INPUT_CFG_ID_SERIAL  = 0x02,
    VIRTIO_INPUT_CFG_ID_DEVIDS  = 0x03,
    VIRTIO_INPUT_CFG_PROP_BITS  = 0x10,
    VIRTIO_INPUT_CFG_EV_BITS    = 0x11,
```



```

    VIRTIO_INPUT_CFG_ABS_INFO    = 0x12,
};

struct virtio_input_absinfo {
    le32  min;
    le32  max;
    le32  fuzz;
    le32  flat;
    le32  res;
};

struct virtio_input_devids {
    le16  bustype;
    le16  vendor;
    le16  product;
    le16  version;
};

struct virtio_input_config {
    u8     select;
    u8     subssel;
    u8     size;
    u8     reserved[5];
    union {
        char string[128];
        u8  bitmap[128];
        struct virtio_input_absinfo abs;
        struct virtio_input_devids ids;
    } u;
};

```

To query a specific piece of information the driver sets *select* and *subssel* accordingly, then checks *size* to see how much information is available. *size* can be zero if no information is available. Strings do not include a NUL terminator. Related evdev ioctl names are provided for reference.

**VIRTIO\_INPUT\_CFG\_ID\_NAME** *subssel* is zero. Returns the name of the device, in *u.string*.

Similar to EVIOCGNAME ioctl for Linux evdev devices.

**VIRTIO\_INPUT\_CFG\_ID\_SERIAL** *subssel* is zero. Returns the serial number of the device, in *u.string*.

**VIRTIO\_INPUT\_CFG\_ID\_DEVIDS** *subssel* is zero. Returns ID information of the device, in *u.ids*.

Similar to EVIOCGID ioctl for Linux evdev devices.

**VIRTIO\_INPUT\_CFG\_PROP\_BITS** *subssel* is zero. Returns input properties of the device, in *u.bitmap*. Individual bits in the bitmap correspond to INPUT\_PROP\_\* constants used by the underlying evdev implementation.

Similar to EVIOCGPROP ioctl for Linux evdev devices.

**VIRTIO\_INPUT\_CFG\_EV\_BITS** *subssel* specifies the event type using EV\_\* constants in the underlying evdev implementation. If *size* is non-zero the event type is supported and a bitmap of supported event codes is returned in *u.bitmap*. Individual bits in the bitmap correspond to implementation-defined input event codes, for example keys or pointing device axes.

Similar to EVIOCGBIT ioctl for Linux evdev devices.

**VIRTIO\_INPUT\_CFG\_ABS\_INFO** *subssel* specifies the absolute axis using ABS\_\* constants in the underlying evdev implementation. Information about the axis will be returned in *u.abs*.

Similar to EVIOCGABS ioctl for Linux evdev devices.

## 5.8.5 Device Initialization

1. The device is queried for supported event types and codes.
2. The eventq is populated with receive buffers.

#### 5.8.5.1 Driver Requirements: Device Initialization

A driver **MUST** set both *select* and *subsel* when querying device configuration, in any order.

A driver **MUST NOT** write to configuration fields other than *select* and *subsel*.

A driver **SHOULD** check the *size* field before accessing the configuration information.

#### 5.8.5.2 Device Requirements: Device Initialization

A device **MUST** set the *size* field to zero if it doesn't support a given *select* and *subsel* combination.

### 5.8.6 Device Operation

1. Input events such as press and release events for keys and buttons, and motion events for pointing devices are sent from the device to the driver using the *eventq*.
2. Status feedback such as keyboard LED updates are sent from the driver to the device using the *statusq*.
3. Both queues use the same *virtio\_input\_event* struct. *type*, *code* and *value* are filled according to the Linux input layer (evdev) interface, except that the fields are in little endian byte order whereas the evdev ioctl interface uses native endian-ness.

```
struct virtio_input_event {  
    le16 type;  
    le16 code;  
    le32 value;  
};
```

#### 5.8.6.1 Driver Requirements: Device Operation

A driver **SHOULD** keep the *eventq* populated with buffers. These buffers **MUST** be device-writable and **MUST** be at least the size of struct *virtio\_input\_event*.

Buffers placed into the *statusq* by a driver **MUST** be at least the size of struct *virtio\_input\_event*.

A driver **SHOULD** ignore *eventq* input events it does not recognize. Note that evdev devices generally maintain backward compatibility by sending redundant events and relying on the consuming side using only the events it understands and ignoring the rest.

#### 5.8.6.2 Device Requirements: Device Operation

A device **MAY** drop input events if the *eventq* does not have enough available buffers. It **SHOULD NOT** drop individual input events if they are part of a sequence forming one input device update. For example, a pointing device update typically consists of several input events, one for each axis, and a terminating *EV\_SYN* event. A device **SHOULD** either buffer or drop the entire sequence.

## 5.9 Crypto Device

The virtio crypto device is a virtual cryptography device as well as a virtual cryptographic accelerator. The virtio crypto device provides the following crypto services: CIPHER, MAC, HASH, and AEAD. Virtio crypto devices have a single control queue and at least one data queue. Crypto operation requests are placed into a data queue, and serviced by the device. Some crypto operation requests are only valid in the context of a session. The role of the control queue is facilitating control operation requests. Sessions management is realized with control operation requests.

### 5.9.1 Device ID

20

## 5.9.2 Virtqueues

0 dataq1

...

N-1 dataqN

N controlq

N is set by *max\_dataqueues*.

## 5.9.3 Feature bits

VIRTIO\_CRYPTO\_F\_REVISION\_1 (0) revision 1. Revision 1 has a specific request format and other enhancements (which result in some additional requirements).

VIRTIO\_CRYPTO\_F\_CIPHER\_STATELESS\_MODE (1) stateless mode requests are supported by the CIPHER service.

VIRTIO\_CRYPTO\_F\_HASH\_STATELESS\_MODE (2) stateless mode requests are supported by the HASH service.

VIRTIO\_CRYPTO\_F\_MAC\_STATELESS\_MODE (3) stateless mode requests are supported by the MAC service.

VIRTIO\_CRYPTO\_F\_AEAD\_STATELESS\_MODE (4) stateless mode requests are supported by the AEAD service.

### 5.9.3.1 Feature bit requirements

Some crypto feature bits require other crypto feature bits (see 2.2.1):

**VIRTIO\_CRYPTO\_F\_CIPHER\_STATELESS\_MODE** Requires VIRTIO\_CRYPTO\_F\_REVISION\_1.

**VIRTIO\_CRYPTO\_F\_HASH\_STATELESS\_MODE** Requires VIRTIO\_CRYPTO\_F\_REVISION\_1.

**VIRTIO\_CRYPTO\_F\_MAC\_STATELESS\_MODE** Requires VIRTIO\_CRYPTO\_F\_REVISION\_1.

**VIRTIO\_CRYPTO\_F\_AEAD\_STATELESS\_MODE** Requires VIRTIO\_CRYPTO\_F\_REVISION\_1.

## 5.9.4 Supported crypto services

The following crypto services are defined:

```
/* CIPHER service */
#define VIRTIO_CRYPTO_SERVICE_CIPHER 0
/* HASH service */
#define VIRTIO_CRYPTO_SERVICE_HASH 1
/* MAC (Message Authentication Codes) service */
#define VIRTIO_CRYPTO_SERVICE_MAC 2
/* AEAD (Authenticated Encryption with Associated Data) service */
#define VIRTIO_CRYPTO_SERVICE_AEAD 3
```

The above constants designate bits used to indicate the which of crypto services are offered by the device as described in, see 5.9.5.

### 5.9.4.1 CIPHER services

The following CIPHER algorithms are defined:

```
#define VIRTIO_CRYPTO_NO_CIPHER 0
#define VIRTIO_CRYPTO_CIPHER_ARC4 1
#define VIRTIO_CRYPTO_CIPHER_AES_ECB 2
#define VIRTIO_CRYPTO_CIPHER_AES_CBC 3
#define VIRTIO_CRYPTO_CIPHER_AES_CTR 4
#define VIRTIO_CRYPTO_CIPHER_DES_ECB 5
#define VIRTIO_CRYPTO_CIPHER_DES_CBC 6
```

#define VIRTIO_CRYPT0_CIPHER_3DES_ECB	7
#define VIRTIO_CRYPT0_CIPHER_3DES_CBC	8
#define VIRTIO_CRYPT0_CIPHER_3DES_CTR	9
#define VIRTIO_CRYPT0_CIPHER_KASUMI_F8	10
#define VIRTIO_CRYPT0_CIPHER_SNOW3G_UA2	11
#define VIRTIO_CRYPT0_CIPHER_AES_F8	12
#define VIRTIO_CRYPT0_CIPHER_AES_XTS	13
#define VIRTIO_CRYPT0_CIPHER_ZUC_EEA3	14

The above constants have two usages:

1. As bit numbers, used to tell the driver which CIPHER algorithms are supported by the device, see [5.9.5](#).
2. As values, used to designate the algorithm in (CIPHER type) crypto operation requests, see [5.9.7.2.1](#).

### 5.9.4.2 HASH services

The following HASH algorithms are defined:

#define VIRTIO_CRYPT0_NO_HASH	0
#define VIRTIO_CRYPT0_HASH_MD5	1
#define VIRTIO_CRYPT0_HASH_SHA1	2
#define VIRTIO_CRYPT0_HASH_SHA_224	3
#define VIRTIO_CRYPT0_HASH_SHA_256	4
#define VIRTIO_CRYPT0_HASH_SHA_384	5
#define VIRTIO_CRYPT0_HASH_SHA_512	6
#define VIRTIO_CRYPT0_HASH_SHA3_224	7
#define VIRTIO_CRYPT0_HASH_SHA3_256	8
#define VIRTIO_CRYPT0_HASH_SHA3_384	9
#define VIRTIO_CRYPT0_HASH_SHA3_512	10
#define VIRTIO_CRYPT0_HASH_SHA3_SHAKE128	11
#define VIRTIO_CRYPT0_HASH_SHA3_SHAKE256	12

The above constants have two usages:

1. As bit numbers, used to tell the driver which HASH algorithms are supported by the device, see [5.9.5](#).
2. As values, used to designate the algorithm in (HASH type) crypto operation requires, see [5.9.7.2.1](#).

### 5.9.4.3 MAC services

The following MAC algorithms are defined:

#define VIRTIO_CRYPT0_NO_MAC	0
#define VIRTIO_CRYPT0_MAC_HMAC_MD5	1
#define VIRTIO_CRYPT0_MAC_HMAC_SHA1	2
#define VIRTIO_CRYPT0_MAC_HMAC_SHA_224	3
#define VIRTIO_CRYPT0_MAC_HMAC_SHA_256	4
#define VIRTIO_CRYPT0_MAC_HMAC_SHA_384	5
#define VIRTIO_CRYPT0_MAC_HMAC_SHA_512	6
#define VIRTIO_CRYPT0_MAC_CMAC_3DES	25
#define VIRTIO_CRYPT0_MAC_CMAC_AES	26
#define VIRTIO_CRYPT0_MAC_KASUMI_F9	27
#define VIRTIO_CRYPT0_MAC_SNOW3G_UIA2	28
#define VIRTIO_CRYPT0_MAC_GMAC_AES	41
#define VIRTIO_CRYPT0_MAC_GMAC_TWOFISH	42
#define VIRTIO_CRYPT0_MAC_CBCMAC_AES	49
#define VIRTIO_CRYPT0_MAC_CBCMAC_KASUMI_F9	50
#define VIRTIO_CRYPT0_MAC_XCBC_AES	53
#define VIRTIO_CRYPT0_MAC_ZUC_EIA3	54

The above constants have two usages:

1. As bit numbers, used to tell the driver which MAC algorithms are supported by the device, see [5.9.5](#).
2. As values, used to designate the algorithm in (MAC type) crypto operation requests, see [5.9.7.2.1](#).

#### 5.9.4.4 AEAD services

The following AEAD algorithms are defined:

```
#define VIRTIO_CRYPT0_NO_AEAD      0
#define VIRTIO_CRYPT0_AEAD_GCM    1
#define VIRTIO_CRYPT0_AEAD_CCM    2
#define VIRTIO_CRYPT0_AEAD_CHACHA20_POLY1305  3
```

The above constants have two usages:

1. As bit numbers, used to tell the driver which AEAD algorithms are supported by the device, see [5.9.5](#).
2. As values, used to designate the algorithm in (DEAD type) crypto operation requests, see [5.9.7.2.1](#).

#### 5.9.5 Device configuration layout

Crypto device configuration uses the following layout structure:

```
struct virtio_crypto_config {
    le32 status;
    le32 max_dataqueues;
    le32 crypto_services;
    /* Detailed algorithms mask */
    le32 cipher_algo_l;
    le32 cipher_algo_h;
    le32 hash_algo;
    le32 mac_algo_l;
    le32 mac_algo_h;
    le32 aead_algo;
    /* Maximum length of cipher key in bytes */
    le32 max_cipher_key_len;
    /* Maximum length of authenticated key in bytes */
    le32 max_auth_key_len;
    le32 reserved;
    /* Maximum size of each crypto request's content in bytes */
    le64 max_size;
};
```

Currently, only one *status* bit is defined: VIRTIO\_CRYPT0\_S\_HW\_READY set indicates that the device is ready to process requests, this bit is read-only for the driver

```
#define VIRTIO_CRYPT0_S_HW_READY  (1 << 0)
```

**max\_dataqueues** is the maximum number of data virtqueues that can be configured by the device. The driver MAY use only one data queue, or it can use more to achieve better performance.

**crypto\_services** crypto service offered, see [5.9.4](#).

**cipher\_algo\_l** CIPHER algorithms bits 0-31, see [5.9.4.1](#).

**cipher\_algo\_h** CIPHER algorithms bits 32-63, see [5.9.4.1](#).

**hash\_algo** HASH algorithms bits, see [5.9.4.2](#).

**mac\_algo\_l** MAC algorithms bits 0-31, see [5.9.4.3](#).

**mac\_algo\_h** MAC algorithms bits 32-63, see [5.9.4.3](#).

**aead\_algo** AEAD algorithms bits, see [5.9.4.4](#).

**max\_cipher\_key\_len** is the maximum length of cipher key supported by the device.

**max\_auth\_key\_len** is the maximum length of authenticated key supported by the device.

**reserved** is reserved for future use.

**max\_size** is the maximum size of the variable-length parameters of data operation of each crypto request's content supported by the device.

**Note:** Unless explicitly stated otherwise all lengths and sizes are in bytes.

#### 5.9.5.1 Device Requirements: Device configuration layout

- The device MUST set *max\_dataqueues* to between 1 and 65535 inclusive.
- The device MUST set the *status* with valid flags, undefined flags MUST NOT be set.
- The device MUST accept and handle requests after *status* is set to `VIRTIO_CRYPTO_S_HW_READY`.
- The device MUST set *crypto\_services* based on the crypto services the device offers.
- The device MUST set detailed algorithms masks for each service advertised by *crypto\_services*. The device MUST NOT set the not defined algorithms bits.
- The device MUST set *max\_size* to show the maximum size of crypto request the device supports.
- The device MUST set *max\_cipher\_key\_len* to show the maximum length of cipher key if the device supports CIPHER service.
- The device MUST set *max\_auth\_key\_len* to show the maximum length of authenticated key if the device supports MAC service.

#### 5.9.5.2 Driver Requirements: Device configuration layout

- The driver MUST read the *status* from the bottom bit of status to check whether the `VIRTIO_CRYPTO_S_HW_READY` is set, and the driver MUST reread it after device reset.
- The driver MUST NOT transmit any requests to the device if the `VIRTIO_CRYPTO_S_HW_READY` is not set.
- The driver MUST read *max\_dataqueues* field to discover the number of data queues the device supports.
- The driver MUST read *crypto\_services* field to discover which services the device is able to offer.
- The driver SHOULD ignore the not defined algorithms bits.
- The driver MUST read the detailed algorithms fields based on *crypto\_services* field.
- The driver SHOULD read *max\_size* to discover the maximum size of the variable-length parameters of data operation of the crypto request's content the device supports and MUST guarantee that the size of each crypto request's content is within the *max\_size*, otherwise the request will fail and the driver MUST reset the device.
- The driver SHOULD read *max\_cipher\_key\_len* to discover the maximum length of cipher key the device supports and MUST guarantee that the *key\_len* (CIPHER service or AEAD service) is within the *max\_cipher\_key\_len* of the device configuration, otherwise the request will fail.
- The driver SHOULD read *max\_auth\_key\_len* to discover the maximum length of authenticated key the device supports and MUST guarantee that the *auth\_key\_len* (MAC service) is within the *max\_auth\_key\_len* of the device configuration, otherwise the request will fail.

### 5.9.6 Device Initialization

#### 5.9.6.1 Driver Requirements: Device Initialization

- The driver MUST configure and initialize all virtqueues.
- The driver MUST read the supported crypto services from bits of *crypto\_services*.
- The driver MUST read the supported algorithms based on *crypto\_services* field.

### 5.9.7 Device Operation

The operation of a virtio crypto device is driven by requests placed on the virtqueues. Requests consist of a queue-type specific header (specifying among others the operation) and an operation specific payload.

If `VIRTIO_CRYPTO_F_REVISION_1` is negotiated the device may support both session mode (See 5.9.7.2.1) and stateless mode operation requests. In stateless mode all operation parameters are supplied as a part of each request, while in session mode, some or all operation parameters are managed within the session. Stateless mode is guarded by feature bits 0-4 on a service level. If stateless mode is negotiated for a service, the service accepts both session mode and stateless requests; otherwise stateless mode requests are rejected (via operation status).

### 5.9.7.1 Operation Status

The device MUST return a status code as part of the operation (both session operation and service operation) result. The valid operation status as follows:

```
enum VIRTIO_CRYPT0_STATUS {
    VIRTIO_CRYPT0_OK = 0,
    VIRTIO_CRYPT0_ERR = 1,
    VIRTIO_CRYPT0_BADMSG = 2,
    VIRTIO_CRYPT0_NOTSUPP = 3,
    VIRTIO_CRYPT0_INVSESS = 4,
    VIRTIO_CRYPT0_NOSPC = 5,
    VIRTIO_CRYPT0_MAX
};
```

- VIRTIO\_CRYPT0\_OK: success.
- VIRTIO\_CRYPT0\_BADMSG: authentication failed (only when AEAD decryption).
- VIRTIO\_CRYPT0\_NOTSUPP: operation or algorithm is unsupported.
- VIRTIO\_CRYPT0\_INVSESS: invalid session ID when executing crypto operations.
- VIRTIO\_CRYPT0\_NOSPC: no free session ID (only when the VIRTIO\_CRYPT0\_F\_REVISION\_1 feature bit is negotiated).
- VIRTIO\_CRYPT0\_ERR: any failure not mentioned above occurs.

### 5.9.7.2 Control Virtqueue

The driver uses the control virtqueue to send control commands to the device, such as session operations (See [5.9.7.2.1](#)).

The header for controlq is of the following form:

```
#define VIRTIO_CRYPT0_OPCODE(service, op) (((service) << 8) | (op))

struct virtio_crypto_ctrl_header {
#define VIRTIO_CRYPT0_CIPHER_CREATE_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_CIPHER, 0x02)
#define VIRTIO_CRYPT0_CIPHER_DESTROY_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_CIPHER, 0x03)
#define VIRTIO_CRYPT0_HASH_CREATE_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_HASH, 0x02)
#define VIRTIO_CRYPT0_HASH_DESTROY_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_HASH, 0x03)
#define VIRTIO_CRYPT0_MAC_CREATE_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_MAC, 0x02)
#define VIRTIO_CRYPT0_MAC_DESTROY_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_MAC, 0x03)
#define VIRTIO_CRYPT0_AEAD_CREATE_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_AEAD, 0x02)
#define VIRTIO_CRYPT0_AEAD_DESTROY_SESSION \
    VIRTIO_CRYPT0_OPCODE(VIRTIO_CRYPT0_SERVICE_AEAD, 0x03)
    le32 opcode;
    /* algo should be service-specific algorithms */
    le32 algo;
    le32 flag;
    le32 reserved;
};
```

The controlq request is composed of four parts:

```
struct virtio_crypto_op_ctrl_req {
    /* Device read only portion */

    struct virtio_crypto_ctrl_header header;

#define VIRTIO_CRYPT0_CTRLQ_OP_SPEC_HDR_LEGACY 56
    /* fixed length fields, opcode specific */
    u8 op_flf[flf_len];

    /* variable length fields, opcode specific */
};
```



```

    u8 op_vlf[vlf_len];

    /* Device write only portion */

    /* op result or completion status */
    u8 op_outcome[outcome_len];
};

```

*header* is a general header (see above).

*op\_flg* is the opcode (in *header*) specific fixed-length parameters.

*flf\_len* depends on the VIRTIO\_CRYPTOF\_REVISION\_1 feature bit (see below).

*op\_vlf* is the opcode (in *header*) specific variable-length parameters.

*vlf\_len* is the size of the specific structure used.

**Note:** The *vlf\_len* of session-destroy operation and the hash-session-create operation is ZERO.

- If the opcode (in *header*) is VIRTIO\_CRYPTOCIPHER\_CREATE\_SESSION then *op\_flg* is struct `virtio_crypto_sym_create_session_flg` if VIRTIO\_CRYPTOF\_REVISION\_1 is negotiated and struct `virtio_crypto_sym_create_session_flg` is padded to 56 bytes if NOT negotiated, and *op\_vlf* is struct `virtio_crypto_sym_create_session_vlf`.
- If the opcode (in *header*) is VIRTIO\_CRYPTOHASH\_CREATE\_SESSION then *op\_flg* is struct `virtio_crypto_hash_create_session_flg` if VIRTIO\_CRYPTOF\_REVISION\_1 is negotiated and struct `virtio_crypto_hash_create_session_flg` is padded to 56 bytes if NOT negotiated.
- If the opcode (in *header*) is VIRTIO\_CRYPTOMAC\_CREATE\_SESSION then *op\_flg* is struct `virtio_crypto_mac_create_session_flg` if VIRTIO\_CRYPTOF\_REVISION\_1 is negotiated and struct `virtio_crypto_mac_create_session_flg` is padded to 56 bytes if NOT negotiated, and *op\_vlf* is struct `virtio_crypto_mac_create_session_vlf`.
- If the opcode (in *header*) is VIRTIO\_CRYPTOAead\_CREATE\_SESSION then *op\_flg* is struct `virtio_crypto_aead_create_session_flg` if VIRTIO\_CRYPTOF\_REVISION\_1 is negotiated and struct `virtio_crypto_aead_create_session_flg` is padded to 56 bytes if NOT negotiated, and *op\_vlf* is struct `virtio_crypto_aead_create_session_vlf`.
- If the opcode (in *header*) is VIRTIO\_CRYPTOCIPHER\_DESTROY\_SESSION or VIRTIO\_CRYPTOHASH\_DESTROY\_SESSION or VIRTIO\_CRYPTOMAC\_DESTROY\_SESSION or VIRTIO\_CRYPTOAead\_DESTROY\_SESSION then *op\_flg* is struct `virtio_crypto_destroy_session_flg` if VIRTIO\_CRYPTOF\_REVISION\_1 is negotiated and struct `virtio_crypto_destroy_session_flg` is padded to 56 bytes if NOT negotiated.

*op\_outcome* stores the result of operation and must be struct `virtio_crypto_destroy_session_input` for destroy session or struct `virtio_crypto_create_session_input` for create session.

*outcome\_len* is the size of the structure used.

### 5.9.7.2.1 Session operation

The session is a handle which describes the cryptographic parameters to be applied to a number of buffers.

The following structure stores the result of session creation set by the device:

```

struct virtio_crypto_create_session_input {
    le64 session_id;
    le32 status;
    le32 padding;
};

```

A request to destroy a session includes the following information:

```

struct virtio_crypto_destroy_session_flg {
    /* Device read only portion */
    le64 session_id;
};

```

```
struct virtio_crypto_destroy_session_input {
    /* Device write only portion */
    u8 status;
};
```

#### 5.9.7.2.1.1 Session operation: HASH session

The fixed-length parameters of HASH session requests is as follows:

```
struct virtio_crypto_hash_create_session_flf {
    /* Device read only portion */

    /* See VIRTIO_CRYPTO_HASH_* above */
    le32 algo;
    /* hash result length */
    le32 hash_result_len;
};
```

#### 5.9.7.2.1.2 Session operation: MAC session

The fixed-length and the variable-length parameters of MAC session requests are as follows:

```
struct virtio_crypto_mac_create_session_flf {
    /* Device read only portion */

    /* See VIRTIO_CRYPTO_MAC_* above */
    le32 algo;
    /* hash result length */
    le32 hash_result_len;
    /* length of authenticated key */
    le32 auth_key_len;
    le32 padding;
};

struct virtio_crypto_mac_create_session_vlf {
    /* Device read only portion */

    /* The authenticated key */
    u8 auth_key[auth_key_len];
};
```

The length of *auth\_key* is specified in *auth\_key\_len* in the struct *virtio\_crypto\_mac\_create\_session\_flf*.

#### 5.9.7.2.1.3 Session operation: Symmetric algorithms session

The request of symmetric session could be the CIPHER algorithms request or the chain algorithms (chaining CIPHER and HASH/MAC) request.

The fixed-length and the variable-length parameters of CIPHER session requests are as follows:

```
struct virtio_crypto_cipher_session_flf {
    /* Device read only portion */

    /* See VIRTIO_CRYPTO_CIPHER* above */
    le32 algo;
    /* length of key */
    le32 key_len;
#define VIRTIO_CRYPTO_OP_ENCRYPT 1
#define VIRTIO_CRYPTO_OP_DECRYPT 2
    /* encryption or decryption */
    le32 op;
    le32 padding;
};

struct virtio_crypto_cipher_session_vlf {
    /* Device read only portion */
```

```

    /* The cipher key */
    u8 cipher_key[key_len];
};

```

The length of *cipher\_key* is specified in *key\_len* in the struct *virtio\_crypto\_cipher\_session\_flf*.

The fixed-length and the variable-length parameters of Chain session requests are as follows:

```

struct virtio_crypto_alg_chain_session_flf {
    /* Device read only portion */

#define VIRTIO_CRYPTO_SYM_ALG_CHAIN_ORDER_HASH_THEN_CIPHER 1
#define VIRTIO_CRYPTO_SYM_ALG_CHAIN_ORDER_CIPHER_THEN_HASH 2
    le32 alg_chain_order;
    /* Plain hash */
#define VIRTIO_CRYPTO_SYM_HASH_MODE_PLAIN 1
    /* Authenticated hash (mac) */
#define VIRTIO_CRYPTO_SYM_HASH_MODE_AUTH 2
    /* Nested hash */
#define VIRTIO_CRYPTO_SYM_HASH_MODE_NESTED 3
    le32 hash_mode;
    struct virtio_crypto_cipher_session_flf cipher_hdr;

#define VIRTIO_CRYPTO_ALG_CHAIN_SESS_OP_SPEC_HDR_SIZE 16
    /* fixed length fields, algo specific */
    u8 algo_flf[VIRTIO_CRYPTO_ALG_CHAIN_SESS_OP_SPEC_HDR_SIZE];

    /* length of the additional authenticated data (AAD) in bytes */
    le32 aad_len;
    le32 padding;
};

struct virtio_crypto_alg_chain_session_vlf {
    /* Device read only portion */

    /* The cipher key */
    u8 cipher_key[key_len];
    /* The authenticated key */
    u8 auth_key[auth_key_len];
};

```

*hash\_mode* decides the type used by *algo\_flf*.

*algo\_flf* is fixed to 16 bytes and MUST contains or be one of the following types:

- struct *virtio\_crypto\_hash\_create\_session\_flf*
- struct *virtio\_crypto\_mac\_create\_session\_flf*

The data of unused part (if has) in *algo\_flf* will be ignored.

The length of *cipher\_key* is specified in *key\_len* in *cipher\_hdr*.

The length of *auth\_key* is specified in *auth\_key\_len* in struct *virtio\_crypto\_mac\_create\_session\_flf*.

The fixed-length parameters of Symmetric session requests are as follows:

```

struct virtio_crypto_sym_create_session_flf {
    /* Device read only portion */

#define VIRTIO_CRYPTO_SYM_SESS_OP_SPEC_HDR_SIZE 48
    /* fixed length fields, opcode specific */
    u8 op_flf[VIRTIO_CRYPTO_SYM_SESS_OP_SPEC_HDR_SIZE];

    /* No operation */
#define VIRTIO_CRYPTO_SYM_OP_NONE 0
    /* Cipher only operation on the data */
#define VIRTIO_CRYPTO_SYM_OP_CIPHER 1
    /* Chain any cipher with any hash or mac operation. The order
       depends on the value of alg_chain_order param */
#define VIRTIO_CRYPTO_SYM_OP_ALGORITHM_CHAINING 2
    le32 op_type;

```

```
    le32 padding;
};
```

*op\_ffl* is fixed to 48 bytes, MUST contains or be one of the following types:

- struct `virtio_crypto_cipher_session_ffl`
- struct `virtio_crypto_alg_chain_session_ffl`

The data of unused part (if has) in *op\_ffl* will be ignored.

*op\_type* decides the type used by *op\_ffl*.

The variable-length parameters of Symmetric session requests are as follows:

```
struct virtio_crypto_sym_create_session_vlf {
    /* Device read only portion */
    /* variable length fields, opcode specific */
    u8 op_vlf[vlf_len];
};
```

*op\_vlf* MUST contains or be one of the following types:

- struct `virtio_crypto_cipher_session_vlf`
- struct `virtio_crypto_alg_chain_session_vlf`

*op\_type* in struct `virtio_crypto_sym_create_session_ffl` decides the type used by *op\_vlf*.

*vlf\_len* is the size of the specific structure used.

#### 5.9.7.2.1.4 Session operation: AEAD session

The fixed-length and the variable-length parameters of AEAD session requests are as follows:

```
struct virtio_crypto_aead_create_session_ffl {
    /* Device read only portion */

    /* See VIRTIO_CRYPT0_AEAD_* above */
    le32 algo;
    /* length of key */
    le32 key_len;
    /* Authentication tag length */
    le32 tag_len;
    /* The length of the additional authenticated data (AAD) in bytes */
    le32 aad_len;
    /* encryption or decryption, See above VIRTIO_CRYPT0_OP_* */
    le32 op;
    le32 padding;
};

struct virtio_crypto_aead_create_session_vlf {
    /* Device read only portion */
    u8 key[key_len];
};
```

The length of *key* is specified in *key\_len* in struct `virtio_crypto_aead_create_session_ffl`.

#### 5.9.7.2.1.5 Driver Requirements: Session operation: create session

- The driver MUST set the *opcode* field based on service type: CIPHER, HASH, MAC, or AEAD.
- The driver MUST set the control general header, the opcode specific header, the opcode specific extra parameters and the opcode specific outcome buffer in turn. See [5.9.7.2](#).
- The driver MUST set the *reversed* field to zero.

#### 5.9.7.2.1.6 Device Requirements: Session operation: create session

- The device MUST use the corresponding opcode specific structure according to the *opcode* in the control general header.
- The device MUST extract extra parameters according to the structures used.
- The device MUST set the *status* field to one of the following values of enum VIRTIO\_CRYPTIO\_STATUS after finish a session creation:
  - VIRTIO\_CRYPTIO\_OK if a session is created successfully.
  - VIRTIO\_CRYPTIO\_NOTSUPP if the requested algorithm or operation is unsupported.
  - VIRTIO\_CRYPTIO\_NOSPC if no free session ID (only when the VIRTIO\_CRYPTIO\_F\_REVISION\_1 feature bit is negotiated).
  - VIRTIO\_CRYPTIO\_ERR if failure not mentioned above occurs.
- The device MUST set the *session\_id* field to a unique session identifier only if the status is set to VIRTIO\_CRYPTIO\_OK.

#### 5.9.7.2.1.7 Driver Requirements: Session operation: destroy session

- The driver MUST set the *opcode* field based on service type: CIPHER, HASH, MAC, or AEAD.
- The driver MUST set the *session\_id* to a valid value assigned by the device when the session was created.

#### 5.9.7.2.1.8 Device Requirements: Session operation: destroy session

- The device MUST set the *status* field to one of the following values of enum VIRTIO\_CRYPTIO\_STATUS.
  - VIRTIO\_CRYPTIO\_OK if a session is created successfully.
  - VIRTIO\_CRYPTIO\_ERR if any failure occurs.

### 5.9.7.3 Data Virtqueue

The driver uses the data virtqueues to transmit crypto operation requests to the device, and completes the crypto operations.

The header for dataq is as follows:

```
struct virtio_crypto_op_header {
#define VIRTIO_CRYPTIO_CIPHER_ENCRYPT \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_CIPHER, 0x00)
#define VIRTIO_CRYPTIO_CIPHER_DECRYPT \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_CIPHER, 0x01)
#define VIRTIO_CRYPTIO_HASH \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_HASH, 0x00)
#define VIRTIO_CRYPTIO_MAC \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_MAC, 0x00)
#define VIRTIO_CRYPTIO_AEAD_ENCRYPT \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_AEAD, 0x00)
#define VIRTIO_CRYPTIO_AEAD_DECRYPT \
    VIRTIO_CRYPTIO_OPCODE(VIRTIO_CRYPTIO_SERVICE_AEAD, 0x01)
    le32 opcode;
    /* algo should be service-specific algorithms */
    le32 algo;
    le64 session_id;
#define VIRTIO_CRYPTIO_FLAG_SESSION_MODE 1
    /* control flag to control the request */
    le32 flag;
    le32 padding;
};
```

**Note:** If VIRTIO\_CRYPTIO\_F\_REVISION\_1 is not negotiated the *flag* is ignored.

If VIRTIO\_CRYPTIO\_F\_REVISION\_1 is negotiated but VIRTIO\_CRYPTIO\_F\_<SERVICE>\_STATELESS\_MODE is not negotiated, then the device SHOULD reject <SERVICE> requests if VIRTIO\_CRYPTIO\_FLAG\_SESSION\_MODE is not set (in *flag*).

The dataq request is composed of four parts:

```
struct virtio_crypto_op_data_req {
    /* Device read only portion */

    struct virtio_crypto_op_header header;

#define VIRTIO_CRYPT0_DATAQ_OP_SPEC_HDR_LEGACY 48
    /* fixed length fields, opcode specific */
    u8 op_flf[flf_len];

    /* Device read && write portion */
    /* variable length fields, opcode specific */
    u8 op_vlf[vlf_len];

    /* Device write only portion */
    struct virtio_crypto_inhdr inhdr;
};
```

*header* is a general header (see above).

*op\_flf* is the opcode (in *header*) specific header.

*flf\_len* depends on the VIRTIO\_CRYPT0\_F\_REVISION\_1 feature bit (see below).

*op\_vlf* is the opcode (in *header*) specific parameters.

*vlf\_len* is the size of the specific structure used.

- If the the opcode (in *header*) is VIRTIO\_CRYPT0\_CIPHER\_ENCRYPT or VIRTIO\_CRYPT0\_CIPHER\_DECRYPT then:
  - If VIRTIO\_CRYPT0\_F\_CIPHER\_STATELESS\_MODE is negotiated, *op\_flf* is struct virtio\_crypto\_sym\_data\_flf\_stateless, and *op\_vlf* is struct virtio\_crypto\_sym\_data\_vlf\_stateless.
  - If VIRTIO\_CRYPT0\_F\_CIPHER\_STATELESS\_MODE is NOT negotiated, *op\_flf* is struct virtio\_crypto\_sym\_data\_flf if VIRTIO\_CRYPT0\_F\_REVISION\_1 is negotiated and struct virtio\_crypto\_sym\_data\_flf is padded to 48 bytes if NOT negotiated, and *op\_vlf* is struct virtio\_crypto\_sym\_data\_vlf.
- If the the opcode (in *header*) is VIRTIO\_CRYPT0\_HASH:
  - If VIRTIO\_CRYPT0\_F\_HASH\_STATELESS\_MODE is negotiated, *op\_flf* is struct virtio\_crypto\_hash\_data\_flf\_stateless, and *op\_vlf* is struct virtio\_crypto\_hash\_data\_vlf\_stateless.
  - If VIRTIO\_CRYPT0\_F\_HASH\_STATELESS\_MODE is NOT negotiated, *op\_flf* is struct virtio\_crypto\_hash\_data\_flf if VIRTIO\_CRYPT0\_F\_REVISION\_1 is negotiated and struct virtio\_crypto\_hash\_data\_flf is padded to 48 bytes if NOT negotiated, and *op\_vlf* is struct virtio\_crypto\_hash\_data\_vlf.
- If the the opcode (in *header*) is VIRTIO\_CRYPT0\_MAC:
  - If VIRTIO\_CRYPT0\_F\_MAC\_STATELESS\_MODE is negotiated, *op\_flf* is struct virtio\_crypto\_mac\_data\_flf\_stateless, and *op\_vlf* is struct virtio\_crypto\_mac\_data\_vlf\_stateless.
  - If VIRTIO\_CRYPT0\_F\_MAC\_STATELESS\_MODE is NOT negotiated, *op\_flf* is struct virtio\_crypto\_mac\_data\_flf if VIRTIO\_CRYPT0\_F\_REVISION\_1 is negotiated and struct virtio\_crypto\_mac\_data\_flf is padded to 48 bytes if NOT negotiated, and *op\_vlf* is struct virtio\_crypto\_mac\_data\_vlf.
- If the the opcode (in *header*) is VIRTIO\_CRYPT0\_AEAD\_ENCRYPT or VIRTIO\_CRYPT0\_AEAD\_DECRYPT then:
  - If VIRTIO\_CRYPT0\_F\_AEAD\_STATELESS\_MODE is negotiated, *op\_flf* is struct virtio\_crypto\_aead\_data\_flf\_stateless, and *op\_vlf* is struct virtio\_crypto\_aead\_data\_vlf\_stateless.
  - If VIRTIO\_CRYPT0\_F\_AEAD\_STATELESS\_MODE is NOT negotiated, *op\_flf* is struct virtio\_crypto\_aead\_data\_flf if VIRTIO\_CRYPT0\_F\_REVISION\_1 is negotiated and struct virtio\_crypto\_aead\_data\_flf is padded to 48 bytes if NOT negotiated, and *op\_vlf* is struct virtio\_crypto\_aead\_data\_vlf.

*inhdr* is a unified input header that used to return the status of the operations, is defined as follows:

```
struct virtio_crypto_inhdr {
    u8 status;
};
```

#### 5.9.7.4 HASH Service Operation

Session mode HASH service requests are as follows:

```
struct virtio_crypto_hash_data_flf {
    /* length of source data */
    le32 src_data_len;
    /* hash result length */
    le32 hash_result_len;
};

struct virtio_crypto_hash_data_vlf {
    /* Device read only portion */
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Hash result data */
    u8 hash_result[hash_result_len];
};
```

Each data request uses the `virtio_crypto_hash_data_flf` structure and the `virtio_crypto_hash_data_vlf` structure to store information used to run the HASH operations.

`src_data` is the source data that will be processed. `src_data_len` is the length of source data. `hash_result` is the result data and `hash_result_len` is the length of it.

Stateless mode HASH service requests are as follows:

```
struct virtio_crypto_hash_data_flf_stateless {
    struct {
        /* See VIRTIO_CRYPT_HASH_* above */
        le32 algo;
    } sess_para;

    /* length of source data */
    le32 src_data_len;
    /* hash result length */
    le32 hash_result_len;
    le32 reserved;
};

struct virtio_crypto_hash_data_vlf_stateless {
    /* Device read only portion */
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Hash result data */
    u8 hash_result[hash_result_len];
};
```

##### 5.9.7.4.1 Driver Requirements: HASH Service Operation

- If the driver uses the session mode, then the driver MUST set `session_id` in struct `virtio_crypto_op_header` to a valid value assigned by the device when the session was created.
- If the `VIRTIO_CRYPT_F_HASH_STATELESS_MODE` feature bit is negotiated, 1) if the driver uses the stateless mode, then the driver MUST set the `flag` field in struct `virtio_crypto_op_header` to ZERO and MUST set the fields in struct `virtio_crypto_hash_data_flf_stateless.sess_para`, 2) if the driver uses the session mode, then the driver MUST set the `flag` field in struct `virtio_crypto_op_header` to `VIRTIO_CRYPT_FLAG_SESSION_MODE`.
- The driver MUST set `opcode` in struct `virtio_crypto_op_header` to `VIRTIO_CRYPT_HASH`.

##### 5.9.7.4.2 Device Requirements: HASH Service Operation

- The device MUST use the corresponding structure according to the `opcode` in the data general header.



- If the VIRTIO\_CRYPTOF\_HASH\_STATELESS\_MODE feature bit is negotiated, the device MUST parse *flag* field in struct `virtio_crypto_op_header` in order to decide which mode the driver uses.
- The device MUST copy the results of HASH operations in the `hash_result[]` if HASH operations success.
- The device MUST set *status* in struct `virtio_crypto_inhdr` to one of the following values of enum `VIRTIO_CRYPTOSTATUS`:
  - `VIRTIO_CRYPTO_OK` if the operation success.
  - `VIRTIO_CRYPTO_NOTSUPP` if the requested algorithm or operation is unsupported.
  - `VIRTIO_CRYPTO_INVSESS` if the session ID invalid when in session mode.
  - `VIRTIO_CRYPTO_ERR` if any failure not mentioned above occurs.

#### 5.9.7.5 MAC Service Operation

Session mode MAC service requests are as follows:

```
struct virtio_crypto_mac_data_flf {
    struct virtio_crypto_hash_data_flf hdr;
};

struct virtio_crypto_mac_data_vlf {
    /* Device read only portion */
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Hash result data */
    u8 hash_result[hash_result_len];
};
```

Each request uses the `virtio_crypto_mac_data_flf` structure and the `virtio_crypto_mac_data_vlf` structure to store information used to run the MAC operations.

*src\_data* is the source data that will be processed. *src\_data\_len* is the length of source data. *hash\_result* is the result data and *hash\_result\_len* is the length of it.

Stateless mode MAC service requests are as follows:

```
struct virtio_crypto_mac_data_flf_stateless {
    struct {
        /* See VIRTIO_CRYPTOMAC_* above */
        le32 algo;
        /* length of authenticated key */
        le32 auth_key_len;
    } sess_para;

    /* length of source data */
    le32 src_data_len;
    /* hash result length */
    le32 hash_result_len;
};

struct virtio_crypto_mac_data_vlf_stateless {
    /* Device read only portion */
    /* The authenticated key */
    u8 auth_key[auth_key_len];
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Hash result data */
    u8 hash_result[hash_result_len];
};
```

*auth\_key* is the authenticated key that will be used during the process. *auth\_key\_len* is the length of the key.

#### 5.9.7.5.1 Driver Requirements: MAC Service Operation

- If the driver uses the session mode, then the driver MUST set *session\_id* in struct *virtio\_crypto\_op\_* header to a valid value assigned by the device when the session was created.
- If the *VIRTIO\_CRYPTOF\_MAC\_STATELESS\_MODE* feature bit is negotiated, 1) if the driver uses the stateless mode, then the driver MUST set the *flag* field in struct *virtio\_crypto\_op\_header* to ZERO and MUST set the fields in struct *virtio\_crypto\_mac\_data\_flf\_stateless.sess\_para*, 2) if the driver uses the session mode, then the driver MUST set the *flag* field in struct *virtio\_crypto\_op\_header* to *VIRTIO\_CRYPTOF\_FLAG\_SESSION\_MODE*.
- The driver MUST set *opcode* in struct *virtio\_crypto\_op\_header* to *VIRTIO\_CRYPTOF\_MAC*.

#### 5.9.7.5.2 Device Requirements: MAC Service Operation

- If the *VIRTIO\_CRYPTOF\_MAC\_STATELESS\_MODE* feature bit is negotiated, the device MUST parse *flag* field in struct *virtio\_crypto\_op\_header* in order to decide which mode the driver uses.
- The device MUST copy the results of MAC operations in the *hash\_result[]* if HASH operations success.
- The device MUST set *status* in struct *virtio\_crypto\_inhdr* to one of the following values of enum *VIRTIO\_CRYPTOF\_STATUS*:
  - *VIRTIO\_CRYPTOF\_OK* if the operation success.
  - *VIRTIO\_CRYPTOF\_NOTSUPP* if the requested algorithm or operation is unsupported.
  - *VIRTIO\_CRYPTOF\_INVSESS* if the session ID invalid when in session mode.
  - *VIRTIO\_CRYPTOF\_ERR* if any failure not mentioned above occurs.

#### 5.9.7.6 Symmetric algorithms Operation

Session mode CIPHER service requests are as follows:

```
struct virtio_crypto_cipher_data_flf {
    /*
     * Byte Length of valid IV/Counter data pointed to by the below iv data.
     *
     * For block ciphers in CBC or F8 mode, or for Kasumi in F8 mode, or for
     * SNOW3G in UEA2 mode, this is the length of the IV (which
     * must be the same as the block length of the cipher).
     * For block ciphers in CTR mode, this is the length of the counter
     * (which must be the same as the block length of the cipher).
     */
    le32 iv_len;
    /* length of source data */
    le32 src_data_len;
    /* length of destination data */
    le32 dst_data_len;
    le32 padding;
};

struct virtio_crypto_cipher_data_vlf {
    /* Device read only portion */

    /*
     * Initialization Vector or Counter data.
     *
     * For block ciphers in CBC or F8 mode, or for Kasumi in F8 mode, or for
     * SNOW3G in UEA2 mode, this is the Initialization Vector (IV)
     * value.
     * For block ciphers in CTR mode, this is the counter.
     * For AES-XTS, this is the 128bit tweak, i, from IEEE Std 1619-2007.
     *
     * The IV/Counter will be updated after every partial cryptographic
     * operation.
     */
    u8 iv[iv_len];
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Destination data */
};
```

```

    u8 dst_data[dst_data_len];
};

```

Session mode requests of algorithm chaining are as follows:

```

struct virtio_crypto_alg_chain_data_flf {
    le32 iv_len;
    /* Length of source data */
    le32 src_data_len;
    /* Length of destination data */
    le32 dst_data_len;
    /* Starting point for cipher processing in source data */
    le32 cipher_start_src_offset;
    /* Length of the source data that the cipher will be computed on */
    le32 len_to_cipher;
    /* Starting point for hash processing in source data */
    le32 hash_start_src_offset;
    /* Length of the source data that the hash will be computed on */
    le32 len_to_hash;
    /* Length of the additional auth data */
    le32 aad_len;
    /* Length of the hash result */
    le32 hash_result_len;
    le32 reserved;
};

struct virtio_crypto_alg_chain_data_vlf {
    /* Device read only portion */

    /* Initialization Vector or Counter data */
    u8 iv[iv_len];
    /* Source data */
    u8 src_data[src_data_len];
    /* Additional authenticated data if exists */
    u8 aad[aad_len];

    /* Device write only portion */

    /* Destination data */
    u8 dst_data[dst_data_len];
    /* Hash result data */
    u8 hash_result[hash_result_len];
};

```

Session mode requests of symmetric algorithm are as follows:

```

struct virtio_crypto_sym_data_flf {
    /* Device read only portion */

#define VIRTIO_CRYPTO_SYM_DATA_REQ_HDR_SIZE    40
    u8 op_type_flf[VIRTIO_CRYPTO_SYM_DATA_REQ_HDR_SIZE];

    /* See above VIRTIO_CRYPTO_SYM_OP_* */
    le32 op_type;
    le32 padding;
};

struct virtio_crypto_sym_data_vlf {
    u8 op_type_vlf[sym_para_len];
};

```

Each request uses the `virtio_crypto_sym_data_flf` structure and the `virtio_crypto_sym_data_vlf` structure to store information used to run the CIPHER operations.

`op_type_flf` is the `op_type` specific header, it MUST starts with or be one of the following structures:

- `struct virtio_crypto_cipher_data_flf`
- `struct virtio_crypto_alg_chain_data_flf`

The length of `op_type_flf` is fixed to 40 bytes, the data of unused part (if has) will be ingored.

*op\_type\_vlf* is the *op\_type* specific parameters, it MUST starts with or be one of the following structures:

- struct *virtio\_crypto\_cipher\_data\_vlf*
- struct *virtio\_crypto\_alg\_chain\_data\_vlf*

*sym\_para\_len* is the size of the specific structure used.

Stateless mode CIPHER service requests are as follows:

```
struct virtio_crypto_cipher_data_flf_stateless {
    struct {
        /* See VIRTIO_CRYPT0_CIPHER* above */
        le32 algo;
        /* length of key */
        le32 key_len;

        /* See VIRTIO_CRYPT0_OP_* above */
        le32 op;
    } sess_para;

    /*
     * Byte Length of valid IV/Counter data pointed to by the below iv data.
     */
    le32 iv_len;
    /* length of source data */
    le32 src_data_len;
    /* length of destination data */
    le32 dst_data_len;
};

struct virtio_crypto_cipher_data_vlf_stateless {
    /* Device read only portion */

    /* The cipher key */
    u8 cipher_key[key_len];

    /* Initialization Vector or Counter data. */
    u8 iv[iv_len];
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */
    /* Destination data */
    u8 dst_data[dst_data_len];
};
```

Stateless mode requests of algorithm chaining are as follows:

```
struct virtio_crypto_alg_chain_data_flf_stateless {
    struct {
        /* See VIRTIO_CRYPT0_SYM_ALG_CHAIN_ORDER_* above */
        le32 alg_chain_order;
        /* length of the additional authenticated data in bytes */
        le32 aad_len;

        struct {
            /* See VIRTIO_CRYPT0_CIPHER* above */
            le32 algo;
            /* length of key */
            le32 key_len;
            /* See VIRTIO_CRYPT0_OP_* above */
            le32 op;
        } cipher;

        struct {
            /* See VIRTIO_CRYPT0_HASH_* or VIRTIO_CRYPT0_MAC_* above */
            le32 algo;
            /* length of authenticated key */
            le32 auth_key_len;
            /* See VIRTIO_CRYPT0_SYM_HASH_MODE_* above */
            le32 hash_mode;
        } hash;
    }
};
```

```

    } sess_para;

    le32 iv_len;
    /* Length of source data */
    le32 src_data_len;
    /* Length of destination data */
    le32 dst_data_len;
    /* Starting point for cipher processing in source data */
    le32 cipher_start_src_offset;
    /* Length of the source data that the cipher will be computed on */
    le32 len_to_cipher;
    /* Starting point for hash processing in source data */
    le32 hash_start_src_offset;
    /* Length of the source data that the hash will be computed on */
    le32 len_to_hash;
    /* Length of the additional auth data */
    le32 aad_len;
    /* Length of the hash result */
    le32 hash_result_len;
    le32 reserved;
};

struct virtio_crypto_alg_chain_data_vlf_stateless {
    /* Device read only portion */

    /* The cipher key */
    u8 cipher_key[key_len];
    /* The auth key */
    u8 auth_key[auth_key_len];
    /* Initialization Vector or Counter data */
    u8 iv[iv_len];
    /* Additional authenticated data if exists */
    u8 aad[aad_len];
    /* Source data */
    u8 src_data[src_data_len];

    /* Device write only portion */

    /* Destination data */
    u8 dst_data[dst_data_len];
    /* Hash result data */
    u8 hash_result[hash_result_len];
};

```

Stateless mode requests of symmetric algorithm are as follows:

```

struct virtio_crypto_sym_data_flf_stateless {
    /* Device read only portion */
#define VIRTIO_CRYPTO_SYM_DATA_REQ_HDR_STATELESS_SIZE    72
    u8 op_type_flf[VIRTIO_CRYPTO_SYM_DATA_REQ_HDR_STATELESS_SIZE];

    /* Device write only portion */
    /* See above VIRTIO_CRYPTO_SYM_OP_* */
    le32 op_type;
};

struct virtio_crypto_sym_data_vlf_stateless {
    u8 op_type_vlf[sym_para_len];
};

```

*op\_type\_flf* is the *op\_type* specific header, it MUST starts with or be one of the following structures:

- struct virtio\_crypto\_cipher\_data\_flf\_stateless
- struct virtio\_crypto\_alg\_chain\_data\_flf\_stateless

The length of *op\_type\_flf* is fixed to 72 bytes, the data of unused part (if has) will be ingored.

*op\_type\_vlf* is the *op\_type* specific parameters, it MUST starts with or be one of the following structures:

- struct virtio\_crypto\_cipher\_data\_vlf\_stateless

- struct virtio\_crypto\_alg\_chain\_data\_vlf\_stateless

*sym\_para\_len* is the size of the specific structure used.

#### 5.9.7.6.1 Driver Requirements: Symmetric algorithms Operation

- If the driver uses the session mode, then the driver MUST set *session\_id* in struct virtio\_crypto\_op\_header to a valid value assigned by the device when the session was created.
- If the VIRTIO\_CRYPTO\_F\_CIPHER\_STATELESS\_MODE feature bit is negotiated, 1) if the driver uses the stateless mode, then the driver MUST set the *flag* field in struct virtio\_crypto\_op\_header to ZERO and MUST set the fields in struct virtio\_crypto\_cipher\_data\_flf\_stateless.sess\_para or struct virtio\_crypto\_alg\_chain\_data\_flf\_stateless.sess\_para, 2) if the driver uses the session mode, then the driver MUST set the *flag* field in struct virtio\_crypto\_op\_header to VIRTIO\_CRYPTO\_FLAG\_SESSION\_MODE.
- The driver MUST set the *opcode* field in struct virtio\_crypto\_op\_header to VIRTIO\_CRYPTO\_CIPHER\_ENCRYPT or VIRTIO\_CRYPTO\_CIPHER\_DECRYPT.
- The driver MUST specify the fields of struct virtio\_crypto\_cipher\_data\_flf in struct virtio\_crypto\_sym\_data\_flf and struct virtio\_crypto\_cipher\_data\_vlf in struct virtio\_crypto\_sym\_data\_vlf if the request is based on VIRTIO\_CRYPTO\_SYM\_OP\_CIPHER.
- The driver MUST specify the fields of struct virtio\_crypto\_alg\_chain\_data\_flf in struct virtio\_crypto\_sym\_data\_flf and struct virtio\_crypto\_alg\_chain\_data\_vlf in struct virtio\_crypto\_sym\_data\_vlf if the request is of the VIRTIO\_CRYPTO\_SYM\_OP\_ALGORITHM\_CHAINING type.

#### 5.9.7.6.2 Device Requirements: Symmetric algorithms Operation

- If the VIRTIO\_CRYPTO\_F\_CIPHER\_STATELESS\_MODE feature bit is negotiated, the device MUST parse *flag* field in struct virtio\_crypto\_op\_header in order to decide which mode the driver uses.
- The device MUST parse the virtio\_crypto\_sym\_data\_req based on the *opcode* field in general header.
- The device MUST parse the fields of struct virtio\_crypto\_cipher\_data\_flf in struct virtio\_crypto\_sym\_data\_flf and struct virtio\_crypto\_cipher\_data\_vlf in struct virtio\_crypto\_sym\_data\_vlf if the request is based on VIRTIO\_CRYPTO\_SYM\_OP\_CIPHER.
- The device MUST parse the fields of struct virtio\_crypto\_alg\_chain\_data\_flf in struct virtio\_crypto\_sym\_data\_flf and struct virtio\_crypto\_alg\_chain\_data\_vlf in struct virtio\_crypto\_sym\_data\_vlf if the request is of the VIRTIO\_CRYPTO\_SYM\_OP\_ALGORITHM\_CHAINING type.
- The device MUST copy the result of cryptographic operation in the *dst\_data[]* in both plain CIPHER mode and algorithms chain mode.
- The device MUST check the *para.add\_len* is bigger than 0 before parse the additional authenticated data in plain algorithms chain mode.
- The device MUST copy the result of HASH/MAC operation in the *hash\_result[]* is of the VIRTIO\_CRYPTO\_SYM\_OP\_ALGORITHM\_CHAINING type.
- The device MUST set the *status* field in struct virtio\_crypto\_inhdr to one of the following values of enum VIRTIO\_CRYPTO\_STATUS:
  - VIRTIO\_CRYPTO\_OK if the operation success.
  - VIRTIO\_CRYPTO\_NOTSUPP if the requested algorithm or operation is unsupported.
  - VIRTIO\_CRYPTO\_INVSESS if the session ID is invalid in session mode.
  - VIRTIO\_CRYPTO\_ERR if failure not mentioned above occurs.

#### 5.9.7.7 AEAD Service Operation

Session mode requests of symmetric algorithm are as follows:

```
struct virtio_crypto_aead_data_flf {
    /*
     * Byte Length of valid IV data.
     *
     * For GCM mode, this is either 12 (for 96-bit IVs) or 16, in which
     * case iv points to J0.
     * For CCM mode, this is the length of the nonce, which can be in the
     * range 7 to 13 inclusive.
     */
}
```

```

le32 iv_len;
/* length of additional auth data */
le32 aad_len;
/* length of source data */
le32 src_data_len;
/* length of dst data, this should be at least src_data_len + tag_len */
le32 dst_data_len;
/* Authentication tag length */
le32 tag_len;
le32 reserved;
};

struct virtio_crypto_aead_data_vlf {
/* Device read only portion */

/*
 * Initialization Vector data.
 *
 * For GCM mode, this is either the IV (if the length is 96 bits) or J0
 * (for other sizes), where J0 is as defined by NIST SP800-38D.
 * Regardless of the IV length, a full 16 bytes needs to be allocated.
 * For CCM mode, the first byte is reserved, and the nonce should be
 * written starting at &iv[1] (to allow space for the implementation
 * to write in the flags in the first byte). Note that a full 16 bytes
 * should be allocated, even though the iv_len field will have
 * a value less than this.
 *
 * The IV will be updated after every partial cryptographic operation.
 */
u8 iv[iv_len];
/* Source data */
u8 src_data[src_data_len];
/* Additional authenticated data if exists */
u8 aad[aad_len];

/* Device write only portion */
/* Pointer to output data */
u8 dst_data[dst_data_len];
};

```

Each request uses the `virtio_crypto_aead_data_flf` structure and the `virtio_crypto_aead_data_vlf` structure to store information used to run the AEAD operations.

Stateless mode AEAD service requests are as follows:

```

struct virtio_crypto_aead_data_flf_stateless {
struct {
/* See VIRTIO_CRYPT0_AEAD_* above */
le32 algo;
/* length of key */
le32 key_len;
/* encrypt or decrypt, See above VIRTIO_CRYPT0_OP_* */
le32 op;
} sess_para;

/* Byte Length of valid IV data. */
le32 iv_len;
/* Authentication tag length */
le32 tag_len;
/* length of additional auth data */
le32 aad_len;
/* length of source data */
le32 src_data_len;
/* length of dst data, this should be at least src_data_len + tag_len */
le32 dst_data_len;
};

struct virtio_crypto_aead_data_vlf_stateless {
/* Device read only portion */

/* The cipher key */

```



```

u8 key[key_len];
/* Initialization Vector data. */
u8 iv[iv_len];
/* Source data */
u8 src_data[src_data_len];
/* Additional authenticated data if exists */
u8 aad[aad_len];

/* Device write only portion */
/* Pointer to output data */
u8 dst_data[dst_data_len];
};

```

#### 5.9.7.7.1 Driver Requirements: AEAD Service Operation

- If the driver uses the session mode, then the driver MUST set *session\_id* in struct *virtio\_crypto\_op\_* header to a valid value assigned by the device when the session was created.
- If the *VIRTIO\_CRYPTOF\_AEAD\_STATELESS\_MODE* feature bit is negotiated, 1) if the driver uses the stateless mode, then the driver MUST set the *flag* field in struct *virtio\_crypto\_op\_header* to ZERO and MUST set the fields in struct *virtio\_crypto\_aead\_data\_vlf\_stateless.ssess\_para*, 2) if the driver uses the session mode, then the driver MUST set the *flag* field in struct *virtio\_crypto\_op\_header* to *VIRTIO\_CRYPTOF\_FLAG\_SESSION\_MODE*.
- The driver MUST set the *opcode* field in struct *virtio\_crypto\_op\_header* to *VIRTIO\_CRYPTOF\_AEAD\_ENCRYPT* or *VIRTIO\_CRYPTOF\_AEAD\_DECRYPT*.

#### 5.9.7.7.2 Device Requirements: AEAD Service Operation

- If the *VIRTIO\_CRYPTOF\_AEAD\_STATELESS\_MODE* feature bit is negotiated, the device MUST parse the *virtio\_crypto\_aead\_data\_vlf\_stateless* based on the *opcode* field in general header.
- The device MUST copy the result of cryptographic operation in the *dst\_data[]*.
- The device MUST copy the authentication tag in the *dst\_data[]* offset the cipher result.
- The device MUST set the *status* field in struct *virtio\_crypto\_inhdr* to one of the following values of enum *VIRTIO\_CRYPTOF\_STATUS*:
- When the *opcode* field is *VIRTIO\_CRYPTOF\_AEAD\_DECRYPT*, the device MUST verify and return the verification result to the driver.
  - *VIRTIO\_CRYPTOF\_OK* if the operation success.
  - *VIRTIO\_CRYPTOF\_NOTSUPP* if the requested algorithm or operation is unsupported.
  - *VIRTIO\_CRYPTOF\_BADMSG* if the verification result is incorrect.
  - *VIRTIO\_CRYPTOF\_INVSESS* if the session ID invalid when in session mode.
  - *VIRTIO\_CRYPTOF\_ERR* if any failure not mentioned above occurs.

## 5.10 Socket Device

The virtio socket device is a zero-configuration socket communications device. It facilitates data transfer between the guest and device without using the Ethernet or IP protocols.

### 5.10.1 Device ID

19

### 5.10.2 Virtqueues

0 rx

1 tx

2 event

### 5.10.3 Feature bits

If no feature bit is set, only stream socket type is supported. If `VIRTIO_VSOCK_F_SEQPACKET` has been negotiated, the device MAY act as if `VIRTIO_VSOCK_F_STREAM` has also been negotiated.

**VIRTIO\_VSOCK\_F\_STREAM (0)** stream socket type is supported.

**VIRTIO\_VSOCK\_F\_SEQPACKET (1)** seqpacket socket type is supported.

### 5.10.4 Device configuration layout

Socket device configuration uses the following layout structure:

```
struct virtio_vsock_config {
    le64 guest_cid;
};
```

The *guest\_cid* field contains the guest's context ID, which uniquely identifies the device for its lifetime. The upper 32 bits of the CID are reserved and zeroed.

The following CIDs are reserved and cannot be used as the guest's context ID:

CID	Notes
0	Reserved
1	Reserved
2	Well-known CID for the host
0xffffffff	Reserved
0xffffffffffffffff	Reserved

### 5.10.5 Device Initialization

1. The guest's cid is read from *guest\_cid*.
2. Buffers are added to the event virtqueue to receive events from the device.
3. Buffers are added to the rx virtqueue to start receiving packets.

### 5.10.6 Device Operation

Packets transmitted or received contain a header before the payload:

```
struct virtio_vsock_hdr {
    le64 src_cid;
    le64 dst_cid;
    le32 src_port;
    le32 dst_port;
    le32 len;
    le16 type;
    le16 op;
    le32 flags;
    le32 buf_alloc;
    le32 fwd_cnt;
};
```

The upper 32 bits of *src\_cid* and *dst\_cid* are reserved and zeroed.

Most packets simply transfer data but control packets are also used for connection and buffer space management. *op* is one of the following operation constants:

```
#define VIRTIO_VSOCK_OP_INVALID        0
/* Connect operations */
#define VIRTIO_VSOCK_OP_REQUEST        1
#define VIRTIO_VSOCK_OP_RESPONSE      2
#define VIRTIO_VSOCK_OP_RST           3
```

```
#define VIRTIO_VSOCK_OP_SHUTDOWN      4
/* To send payload */
#define VIRTIO_VSOCK_OP_RW            5
/* Tell the peer our credit info */
#define VIRTIO_VSOCK_OP_CREDIT_UPDATE 6
/* Request the peer to send the credit info to us */
#define VIRTIO_VSOCK_OP_CREDIT_REQUEST 7
```

### 5.10.6.1 Virtqueue Flow Control

The tx virtqueue carries packets initiated by applications and replies to received packets. The rx virtqueue carries packets initiated by the device and replies to previously transmitted packets.

If both rx and tx virtqueues are filled by the driver and device at the same time then it appears that a deadlock is reached. The driver has no free tx descriptors to send replies. The device has no free rx descriptors to send replies either. Therefore neither device nor driver can process virtqueues since that may involve sending new replies.

This is solved using additional resources outside the virtqueue to hold packets. With additional resources, it becomes possible to process incoming packets even when outgoing packets cannot be sent.

Eventually even the additional resources will be exhausted and further processing is not possible until the other side processes the virtqueue that it has neglected. This stop to processing prevents one side from causing unbounded resource consumption in the other side.

#### 5.10.6.1.1 Driver Requirements: Device Operation: Virtqueue Flow Control

The rx virtqueue MUST be processed even when the tx virtqueue is full so long as there are additional resources available to hold packets outside the tx virtqueue.

#### 5.10.6.1.2 Device Requirements: Device Operation: Virtqueue Flow Control

The tx virtqueue MUST be processed even when the rx virtqueue is full so long as there are additional resources available to hold packets outside the rx virtqueue.

### 5.10.6.2 Addressing

Flows are identified by a (source, destination) address tuple. An address consists of a (cid, port number) tuple. The header fields used for this are *src\_cid*, *src\_port*, *dst\_cid*, and *dst\_port*.

Currently stream and seqpacket sockets are supported. *type* is 1 (VIRTIO\_VSOCK\_TYPE\_STREAM) for stream socket types, and 2 (VIRTIO\_VSOCK\_TYPE\_SEQPACKET) for seqpacket socket types.

```
#define VIRTIO_VSOCK_TYPE_STREAM      1
#define VIRTIO_VSOCK_TYPE_SEQPACKET  2
```

Stream sockets provide in-order, guaranteed, connection-oriented delivery without message boundaries. Seqpacket sockets provide in-order, guaranteed, connection-oriented delivery with message and record boundaries.

### 5.10.6.3 Buffer Space Management

*buf\_alloc* and *fwd\_cnt* are used for buffer space management of stream sockets. The guest and the device publish how much buffer space is available per socket. Only payload bytes are counted and header bytes are not included. This facilitates flow control so data is never dropped.

*buf\_alloc* is the total receive buffer space, in bytes, for this socket. This includes both free and in-use buffers. *fwd\_cnt* is the free-running bytes received counter. The sender calculates the amount of free receive buffer space as follows:

```
/* tx_cnt is the sender's free-running bytes transmitted counter */
u32 peer_free = peer_buf_alloc - (tx_cnt - peer_fwd_cnt);
```

If there is insufficient buffer space, the sender waits until virtqueue buffers are returned and checks *buf\_alloc* and *fwd\_cnt* again. Sending the `VIRTIO_VSOCK_OP_CREDIT_REQUEST` packet queries how much buffer space is available. The reply to this query is a `VIRTIO_VSOCK_OP_CREDIT_UPDATE` packet. It is also valid to send a `VIRTIO_VSOCK_OP_CREDIT_UPDATE` packet without previously receiving a `VIRTIO_VSOCK_OP_CREDIT_REQUEST` packet. This allows communicating updates any time a change in buffer space occurs.

#### 5.10.6.3.1 Driver Requirements: Device Operation: Buffer Space Management

`VIRTIO_VSOCK_OP_RW` data packets MUST only be transmitted when the peer has sufficient free buffer space for the payload.

All packets associated with a stream flow MUST contain valid information in *buf\_alloc* and *fwd\_cnt* fields.

#### 5.10.6.3.2 Device Requirements: Device Operation: Buffer Space Management

`VIRTIO_VSOCK_OP_RW` data packets MUST only be transmitted when the peer has sufficient free buffer space for the payload.

All packets associated with a stream flow MUST contain valid information in *buf\_alloc* and *fwd\_cnt* fields.

#### 5.10.6.4 Receive and Transmit

The driver queues outgoing packets on the tx virtqueue and incoming packet receive buffers on the rx virtqueue. Packets are of the following form:

```
struct virtio_vsock_packet {
    struct virtio_vsock_hdr hdr;
    u8 data[];
};
```

Virtqueue buffers for outgoing packets are read-only. Virtqueue buffers for incoming packets are write-only.

##### 5.10.6.4.1 Driver Requirements: Device Operation: Receive and Transmit

The *guest\_cid* configuration field MUST be used as the source CID when sending outgoing packets.

A `VIRTIO_VSOCK_OP_RST` reply MUST be sent if a packet is received with an unknown *type* value.

##### 5.10.6.4.2 Device Requirements: Device Operation: Receive and Transmit

The *guest\_cid* configuration field MUST NOT contain a reserved CID as listed in 5.10.4.

A `VIRTIO_VSOCK_OP_RST` reply MUST be sent if a packet is received with an unknown *type* value.

#### 5.10.6.5 Stream Sockets

Connections are established by sending a `VIRTIO_VSOCK_OP_REQUEST` packet. If a listening socket exists on the destination a `VIRTIO_VSOCK_OP_RESPONSE` reply is sent and the connection is established. A `VIRTIO_VSOCK_OP_RST` reply is sent if a listening socket does not exist on the destination or the destination has insufficient resources to establish the connection.

When a connected socket receives `VIRTIO_VSOCK_OP_SHUTDOWN` the header *flags* field bit `VIRTIO_VSOCK_SHUTDOWN_F_RECEIVE` (bit 0) set indicates that the peer will not receive any more data and bit `VIRTIO_VSOCK_SHUTDOWN_F_SEND` (bit 1) set indicates that the peer will not send any more data. These hints are permanent once sent and successive packets with bits clear do not reset them.

```
#define VIRTIO_VSOCK_SHUTDOWN_F_RECEIVE 0
#define VIRTIO_VSOCK_SHUTDOWN_F_SEND 1
```

The VIRTIO\_VSOCK\_OP\_RST packet aborts the connection process or forcibly disconnects a connected socket.

Clean disconnect is achieved by one or more VIRTIO\_VSOCK\_OP\_SHUTDOWN packets that indicate no more data will be sent and received, followed by a VIRTIO\_VSOCK\_OP\_RST response from the peer. If no VIRTIO\_VSOCK\_OP\_RST response is received within an implementation-specific amount of time, a VIRTIO\_VSOCK\_OP\_RST packet is sent to forcibly disconnect the socket.

The clean disconnect process ensures that neither peer reuses the (source, destination) address tuple for a new connection while the other peer is still processing the old connection.

### 5.10.6.6 Seqpacket Sockets

#### 5.10.6.6.1 Message and record boundaries

Two types of boundaries are supported: message and record boundaries.

A message contains data sent in a single operation. A single message can be split into multiple RW packets. To provide message boundaries, last RW packet of each message has VIRTIO\_VSOCK\_SEQ\_EOM bit (bit 0) set in the *flags* of packet's header.

Record is any number of subsequent messages, where last message is sent with POSIX MSG\_EOR flag set. Record boundary means that receiver gets MSG\_EOR flag set in the corresponding message where sender set it. To provide record boundaries, last RW packet of each record has VIRTIO\_VSOCK\_SEQ\_EOR bit (bit 1) set in the *flags* of packet's header.

```
#define VIRTIO_VSOCK_SEQ_EOM (1 << 0)
#define VIRTIO_VSOCK_SEQ_EOR (1 << 1)
```

### 5.10.6.7 Device Events

Certain events are communicated by the device to the driver using the event virtqueue.

The event buffer is as follows:

```
#define VIRTIO_VSOCK_EVENT_TRANSPORT_RESET 0

struct virtio_vsock_event {
    le32 id;
};
```

The VIRTIO\_VSOCK\_EVENT\_TRANSPORT\_RESET event indicates that communication has been interrupted. This usually occurs if the guest has been physically migrated. The driver shuts down established connections and the *guest\_cid* configuration field is fetched again. Existing listen sockets remain but their CID is updated to reflect the current *guest\_cid*.

#### 5.10.6.7.1 Driver Requirements: Device Operation: Device Events

Event virtqueue buffers SHOULD be replenished quickly so that no events are missed.

The *guest\_cid* configuration field MUST be fetched to determine the current CID when a VIRTIO\_VSOCK\_EVENT\_TRANSPORT\_RESET event is received.

Existing connections MUST be shut down when a VIRTIO\_VSOCK\_EVENT\_TRANSPORT\_RESET event is received.

Listen connections MUST remain operational with the current CID when a VIRTIO\_VSOCK\_EVENT\_TRANSPORT\_RESET event is received.

## 5.11 File System Device

The virtio file system device provides file system access. The device either directly manages a file system or it acts as a gateway to a remote file system. The details of how the device implementation accesses files are hidden by the device interface, allowing for a range of use cases.

Unlike block-level storage devices such as virtio block and SCSI, the virtio file system device provides file-level access to data. The device interface is based on the Linux Filesystem in Userspace (FUSE) protocol. This consists of requests for file system traversal and access the files and directories within it. The protocol details are defined by [FUSE](#).

The device acts as the FUSE file system daemon and the driver acts as the FUSE client mounting the file system. The virtio file system device provides the mechanism for transporting FUSE requests, much like `/dev/fuse` in a traditional FUSE application.

This section relies on definitions from [FUSE](#).

### 5.11.1 Device ID

26

### 5.11.2 Virtqueues

**0** hiprio

**1** notification queue

**2...n** request queues

The notification queue only exists if `VIRTIO_FS_F_NOTIFICATION` is set.

### 5.11.3 Feature bits

**VIRTIO\_FS\_F\_NOTIFICATION (0)** Device has support for FUSE notify messages. The notification queue is virtqueue 1.

### 5.11.4 Device configuration layout

```
struct virtio_fs_config {
    char tag[36];
    le32 num_request_queues;
    le32 notify_buf_size;
};
```

The *tag* and *num\_request\_queues* fields are always available. The *notify\_buf\_size* field is only available when `VIRTIO_FS_F_NOTIFICATION` is set.

**tag** is the name associated with this file system. The tag is encoded in UTF-8 and padded with NUL bytes if shorter than the available space. This field is not NUL-terminated if the encoded bytes take up the entire field.

**num\_request\_queues** is the total number of request virtqueues exposed by the device. Each virtqueue offers identical functionality and there are no ordering guarantees between requests made available on different queues. Use of multiple queues is intended to increase performance.

**notify\_buf\_size** is the minimum number of bytes required for each buffer in the notification queue.

#### 5.11.4.1 Driver Requirements: Device configuration layout

The driver **MUST NOT** write to device configuration fields.

The driver **MAY** use from one up to *num\_request\_queues* request virtqueues.

#### 5.11.4.2 Device Requirements: Device configuration layout

The device MUST set *num\_request\_queues* to 1 or greater.

The device MUST set *notify\_buf\_size* to be large enough to hold any of the FUSE notify messages that this device emits.

#### 5.11.5 Device Initialization

On initialization the driver first discovers the device's virtqueues.

The driver populates the notification queue with buffers for receiving FUSE notify messages if *VIRTIO\_FS\_F\_NOTIFICATION* is set.

The FUSE session is started by sending a *FUSE\_INIT* request as defined by the FUSE protocol on one request virtqueue. All virtqueues provide access to the same FUSE session and therefore only one *FUSE\_INIT* request is required regardless of the number of available virtqueues.

#### 5.11.6 Device Operation

Device operation consists of operating the virtqueues to facilitate file system access.

The FUSE request types are as follows:

- Normal requests are made available by the driver on request queues and are used by the device.
- High priority requests (*FUSE\_INTERRUPT*, *FUSE\_FORGET*, and *FUSE\_BATCH\_FORGET*) are made available by the driver on the *hiprio* queue so the device is able to process them even if the request queues are full.

FUSE notify messages are received on the notification queue if *VIRTIO\_FS\_F\_NOTIFICATION* is set.

##### 5.11.6.1 Device Operation: Request Queues

The driver enqueues normal requests on an arbitrary request queue. High priority requests are not placed on request queues. The device processes requests in any order. The driver is responsible for ensuring that ordering constraints are met by making available a dependent request only after its prerequisite request has been used.

Requests have the following format with endianness chosen by the driver in the *FUSE\_INIT* request used to initiate the session as detailed below:

```
struct virtio_fs_req {
    // Device-readable part
    struct fuse_in_header in;
    u8 datain[];

    // Device-writable part
    struct fuse_out_header out;
    u8 dataout[];
};
```

Note that the words "in" and "out" follow the FUSE meaning and do not indicate the direction of data transfer under *VIRTIO*. "In" means input to a request and "out" means output from processing a request.

*in* is the common header for all types of FUSE requests.

*datain* consists of request-specific data, if any. This is identical to the data read from the */dev/fuse* device by a FUSE daemon.

*out* is the completion header common to all types of FUSE requests.

*dataout* consists of request-specific data, if any. This is identical to the data written to the */dev/fuse* device by a FUSE daemon.

For example, the full layout of a *FUSE\_READ* request is as follows:



```

struct virtio_fs_read_req {
    // Device-readable part
    struct fuse_in_header in;
    union {
        struct fuse_read_in readin;
        u8 datain[sizeof(struct fuse_read_in)];
    };

    // Device-writable part
    struct fuse_out_header out;
    u8 dataout[out.len - sizeof(struct fuse_out_header)];
};

```

The FUSE protocol documented in [FUSE](#) specifies the set of request types and their contents.

The endianness of the FUSE protocol session is detectable by inspecting the `uint32_t in.opcode` field of the FUSE\_INIT request sent by the driver to the device. This allows the device to determine whether the session is little-endian or big-endian. The next FUSE\_INIT message terminates the current session and starts a new session with the possibility of changing endianness.

#### 5.11.6.2 Device Operation: High Priority Queue

The hiprio queue follows the same request format as the request queues. This queue only contains FUSE\_INTERRUPT, FUSE\_FORGET, and FUSE\_BATCH\_FORGET requests.

Interrupt and forget requests have a higher priority than normal requests. The separate hiprio queue is used for these requests to ensure they can be delivered even when all request queues are full.

##### 5.11.6.2.1 Device Requirements: Device Operation: High Priority Queue

The device MUST NOT pause processing of the hiprio queue due to activity on a normal request queue.

The device MAY process request queues concurrently with the hiprio queue.

##### 5.11.6.2.2 Driver Requirements: Device Operation: High Priority Queue

The driver MUST submit FUSE\_INTERRUPT, FUSE\_FORGET, and FUSE\_BATCH\_FORGET requests solely on the hiprio queue.

The driver MUST not submit normal requests on the hiprio queue.

The driver MUST anticipate that request queues are processed concurrently with the hiprio queue.

#### 5.11.6.3 Device Operation: Notification Queue

The notification queue is populated with buffers by the driver and these buffers are used by the device to emit FUSE notify messages. Notification queue buffer layout is as follows:

```

struct virtio_fs_notify {
    // Device-writable part
    struct fuse_out_header out_hdr;
    char outarg[notify_buf_size - sizeof(struct fuse_out_header)];
};

```

*outarg* contains the FUSE notify message payload that depends on the type of notification being emitted.

If the driver provides notification queue buffers at a slower rate than the device emits FUSE notify messages then the virtqueue will eventually become empty. The behavior in response to an empty virtqueue depends on the FUSE notify message type. The following FUSE notify message types are supported:

- FUSE\_NOTIFY\_LOCK messages are delivered when buffers become available again. The device has resources for a certain number of lock requests. If the device runs out of resources new lock requests fail with ENOLCK.

#### 5.11.6.3.1 Driver Requirements: Device Operation: Notification Queue

The driver MUST provide buffers of at least *notify\_buf\_size* bytes.

The driver SHOULD replenish notification queue buffers sufficiently quickly so that there is always at least one available buffer.

#### 5.11.6.4 Device Operation: DAX Window

FUSE\_READ and FUSE\_WRITE requests transfer file contents between the driver-provided buffer and the device. In cases where data transfer is undesirable, the device can map file contents into the DAX window shared memory region. The driver then accesses file contents directly in device-owned memory without a data transfer.

The DAX Window is an alternative mechanism for accessing file contents. FUSE\_READ/FUSE\_WRITE requests and DAX Window accesses are possible at the same time. Providing the DAX Window is optional for devices. Using the DAX Window is optional for drivers.

Shared memory region ID 0 is called the DAX window. Drivers map this shared memory region with writeback caching as if it were regular RAM. The contents of the DAX window are undefined unless a mapping exists for that range.

The driver maps a file range into the DAX window using the FUSE\_SETUPMAPPING request. Alignment constraints for FUSE\_SETUPMAPPING and FUSE\_REMOVEMAPPING requests are communicated during FUSE\_INIT negotiation.

When a FUSE\_SETUPMAPPING request perfectly overlaps a previous mapping, the previous mapping is replaced. When a mapping partially overlaps a previous mapping, the previous mapping is split into one or two smaller mappings. When a mapping is partially unmapped it is also split into one or two smaller mappings.

Establishing new mappings or splitting existing mappings consumes resources. If the device runs out of resources the FUSE\_SETUPMAPPING request fails until resources are available again following FUSE\_REMOVEMAPPING.

After FUSE\_SETUPMAPPING has completed successfully the file range is accessible from the DAX window at the offset provided by the driver in the request. A mapping is removed using the FUSE\_REMOVEMAPPING request.

Data is only guaranteed to be persistent when a FUSE\_FSYNC request is used by the device after having been made available by the driver following the write.

##### 5.11.6.4.1 Device Requirements: Device Operation: DAX Window

The device MAY provide the DAX Window to memory-mapped access to file contents. If present, the DAX Window MUST be shared memory region ID 0.

The device MUST support FUSE\_READ and FUSE\_WRITE requests regardless of whether the DAX Window is being used or not.

The device MUST allow mappings that completely or partially overlap existing mappings within the DAX window.

The device MUST reject mappings that would go beyond the end of the DAX window.

##### 5.11.6.4.2 Driver Requirements: Device Operation: DAX Window

The driver SHOULD be prepared to find shared memory region ID 0 absent and fall back to FUSE\_READ and FUSE\_WRITE requests.

The driver MAY use both FUSE\_READ/FUSE\_WRITE requests and the DAX Window to access file contents.

The driver MUST NOT access DAX window areas that have not been mapped.

#### 5.11.6.5 Security Considerations

The device provides access to a file system containing files owned by one or more POSIX user ids and group ids. The device has no secure way of differentiating between users originating requests via the driver. Therefore the device accepts the POSIX user ids and group ids provided by the driver and security is enforced by the driver rather than the device. It is nevertheless possible for devices to implement POSIX user id and group id mapping or whitelisting to control the ownership and access available to the driver.

File systems containing special files including device nodes and setuid executable files pose a security concern. These properties are defined by the file type and mode, which are set by the driver when creating new files or by changes at a later time. These special files present a security risk when the file system is shared with another machine. A setuid executable or a device node placed by a malicious machine make it possible for unprivileged users on other machines to elevate their privileges through the shared file system. This issue can be solved on some operating systems using mount options that ignore special files. It is also possible for devices to implement restrictions on special files by refusing their creation.

When the device provides shared access to a file system between multiple machines, symlink race conditions, exhausting file system capacity, and overwriting or deleting files used by others are factors to consider. These issues have a long history in multi-user operating systems and also apply to virtio-fs. They are typically managed at the file system administration level by providing shared access only to mutually trusted users.

Multiple machines sharing access to a file system are susceptible to timing side-channel attacks. By measuring the latency of accesses to file contents or file system metadata it is possible to infer whether other machines also accessed the same information. Short latencies indicate that the information was cached due to a previous access. This can reveal sensitive information, such as whether certain code paths were taken. The DAX Window provides direct access to file contents and is therefore a likely target of such attacks. These attacks are also possible with traditional FUSE requests. The safest approach is to avoid sharing file systems between untrusted machines.

#### 5.11.6.6 Live migration considerations

When a driver is migrated to a new device it is necessary to consider the FUSE session and its state. The continuity of FUSE inode numbers (also known as nodeids) and fh values is necessary so the driver can continue operation without disruption.

It is possible to maintain the FUSE session across live migration either by transferring the state or by redirecting requests from the new device to the old device where the state resides. The details of how to achieve this are implementation-dependent and are not visible at the device interface level.

Maintaining version and feature information negotiated by FUSE\_INIT is necessary so that no FUSE protocol feature changes are visible to the driver across live migration. The FUSE\_INIT information forms part of the FUSE session state that needs to be transferred during live migration.

## 5.12 RPMB Device

virtio-rpmb is a virtio based RPMB (Replay Protected Memory Block) device. It is used as a tamper-resistant and anti-replay storage. The device is driven via requests including read, write, get write counter and program key, which are submitted via a request queue. This section relies on definitions from paragraph 6.6.22 of [eMMC](#).

### 5.12.1 Device ID

28

### 5.12.2 Virtqueues

0 requestq

### 5.12.3 Feature bits

None.

### 5.12.4 Device configuration layout

All fields of this configuration are always available and read-only for the driver.

```
struct virtio_rpmb_config {
    u8 capacity;
    u8 max_wr_cnt;
    u8 max_rd_cnt;
}
```

**capacity** is the capacity of the device (expressed in 128KB units). The values MUST range between 0x00 and 0x80 inclusive.

**max\_wr\_cnt and max\_rd\_cnt** are the maximum numbers of RPMB block count (256B) that can be performed to device in one request. 0 implies no limitation.

### 5.12.5 Device Requirements: Device Initialization

1. The virtqueue is initialized.
2. The device capacity MUST be initialized to a multiple of 128Kbytes and up to 16Mbytes.

### 5.12.6 Device Operation

The operation of a virtio RPMB device is driven by the requests placed on the virtqueue. The type of request can be program key (VIRTIO\_RPMB\_REQ\_PROGRAM\_KEY), get write counter (VIRTIO\_RPMB\_REQ\_GET\_WRITE\_COUNTER), write (VIRTIO\_RPMB\_REQ\_DATA\_WRITE), and read (VIRTIO\_RPMB\_REQ\_DATA\_READ). A program key or write request can also combine with a result read (VIRTIO\_RPMB\_REQ\_RESULT\_READ) for a returned result.

```
/* RPMB Request Types */
#define VIRTIO_RPMB_REQ_PROGRAM_KEY      0x0001
#define VIRTIO_RPMB_REQ_GET_WRITE_COUNTER 0x0002
#define VIRTIO_RPMB_REQ_DATA_WRITE      0x0003
#define VIRTIO_RPMB_REQ_DATA_READ       0x0004
#define VIRTIO_RPMB_REQ_RESULT_READ     0x0005

/* RPMB Response Types */
#define VIRTIO_RPMB_RESP_PROGRAM_KEY     0x0100
#define VIRTIO_RPMB_RESP_GET_COUNTER     0x0200
#define VIRTIO_RPMB_RESP_DATA_WRITE      0x0300
#define VIRTIO_RPMB_RESP_DATA_READ       0x0400
```

**VIRTIO\_RPMB\_REQ\_PROGRAM\_KEY** requests for authentication key programming. If VIRTIO\_RPMB\_REQ\_RESULT\_READ is requested, the device returns the RPMB frame with the response (VIRTIO\_RPMB\_RESP\_PROGRAM\_KEY), the calculated MAC and the result.

**VIRTIO\_RPMB\_REQ\_GET\_WRITE\_COUNTER** requests for reading the write counter. The device returns the RPMB frame with the response (VIRTIO\_RPMB\_RESP\_GET\_COUNTER), the writer counter, a copy of the nonce received in the request, the calculated MAC and the result.

**VIRTIO\_RPMB\_REQ\_DATA\_WRITE** requests for authenticated data write. If VIRTIO\_RPMB\_REQ\_RESULT\_READ is requested, the device returns the RPMB data frame with the response (VIRTIO\_RPMB\_RESP\_DATA\_WRITE), the incremented counter value, the data address, the calculated MAC and the result.

**VIRTIO\_RPMB\_REQ\_DATA\_READ** requests for authenticated data read. The device returns the RPMB frame with the response (VIRTIO\_RPMB\_RESP\_DATA\_READ), the block count, a copy of the nonce received in the request, the address, the data, the calculated MAC and the result.

**VIRTIO\_RPMB\_REQ\_RESULT\_READ** requests for a returned result. It is used following with **VIRTIO\_RPMB\_REQ\_PROGRAM\_KEY** or **VIRTIO\_RPMB\_REQ\_DATA\_WRITE** request types for a returned result in one or multiple RPMB frames. If it's not requested, the device will not return result frame to the driver.

#### 5.12.6.1 Device Operation: Request Queue

The request information is delivered in RPMB frame. The frame is in size of 512B.

```
struct virtio_rpmb_frame {
    u8 stuff[196];
    u8 key_mac[32];
    u8 data[256];
    u8 nonce[16];
    be32 write_counter;
    be16 address;
    be16 block_count;
    be16 result;
    be16 req_resp;
};

/* RPMB Operation Results */
#define VIRTIO_RPMB_RES_OK 0x0000
#define VIRTIO_RPMB_RES_GENERAL_FAILURE 0x0001
#define VIRTIO_RPMB_RES_AUTH_FAILURE 0x0002
#define VIRTIO_RPMB_RES_COUNT_FAILURE 0x0003
#define VIRTIO_RPMB_RES_ADDR_FAILURE 0x0004
#define VIRTIO_RPMB_RES_WRITE_FAILURE 0x0005
#define VIRTIO_RPMB_RES_READ_FAILURE 0x0006
#define VIRTIO_RPMB_RES_NO_AUTH_KEY 0x0007
#define VIRTIO_RPMB_RES_WRITE_COUNTER_EXPIRED 0x0080
```

**stuff** Padding for the frame.

**key\_mac** is the authentication key or the message authentication code (MAC) depending on the request/response type. If the request is **VIRTIO\_RPMB\_REQ\_PROGRAM\_KEY**, it's used as an authentication key. Otherwise, it's used as MAC. The MAC is calculated using HMAC SHA-256. It takes as input a key and a message. The key used for the MAC calculation is always the 256-bit RPMB authentication key. The message used as input to the MAC calculation is the concatenation of the fields in the RPMB frames excluding stuff bytes and the MAC itself.

**data** is used to be written or read via authenticated read/write access. It's fixed 256B.

**nonce** is a random number generated by the user for the read or get write counter requests and copied to the response by the device. It's used for anti-replay protection.

**writer\_counter** is the counter value for the total amount of the successful authenticated data write requests.

**address** is the address of the data to be written to or read from the RPMB virtio device. It is the number of the accessed half sector (256B).

**block\_count** is the number of blocks (256B) requested to be read/written. It's limited by *max\_wr\_cnt* or *max\_rd\_cnt*. For RPMB read request, one virtio buffer including request command and the subsequent [*block\_count*] virtio buffers for response data are placed in the queue. For RPMB write request, [*block\_count*] virtio buffers including request command and data are placed in the queue.

**result** includes information about the status of access made to the device. It is written by the device.

**req\_resp** is the type of request or response, to/from the device.

##### 5.12.6.1.1 Device Requirements: Device Operation: Program Key

If **VIRTIO\_RPMB\_REQ\_RESULT\_READ** is requested, the device SHOULD return the RPMB frame with the response, the calculated MAC and the result:

1. If the *block\_count* is not set to 1 then **VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE** SHOULD be responded as *result*.

2. If the programming of authentication key fails, then VIRTIO\_RPMB\_RES\_WRITE\_FAILURE SHOULD be responded as *result*.
3. If some other error occurs then returned result VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE SHOULD be responded as *result*.
4. The *req\_resp* value VIRTIO\_RPMB\_RESP\_PROGRAM\_KEY SHOULD be responded.

#### 5.12.6.1.2 Device Requirements: Device Operation: Get Write Counter

If the authentication key is not yet programmed then VIRTIO\_RPMB\_RES\_NO\_AUTH\_KEY SHOULD be returned in *result*.

If block count has not been set to 1 then VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE SHOULD be responded as *result*.

The *req\_resp* value VIRTIO\_RPMB\_RESP\_GET\_COUNTER SHOULD be responded.

#### 5.12.6.1.3 Device Requirements: Device Operation: Data Write

If VIRTIO\_RPMB\_REQ\_RESULT\_READ is requested, the device SHOULD return the RPMB data frame with the response VIRTIO\_RPMB\_RESP\_DATA\_WRITE, the incremented counter value, the data address, the calculated MAC and the result:

1. If the authentication key is not yet programmed, then VIRTIO\_RPMB\_RES\_NO\_AUTH\_KEY SHOULD be returned in *result*.
2. If block count is zero or greater than *max\_wr\_cnt* then VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE SHOULD be responded.
3. The device MUST check whether the write counter has expired. If the write counter is expired then the *result* SHOULD be set to VIRTIO\_RPMB\_RES\_WRITE\_COUNTER\_EXPIRED.
4. If there is an error in the address (out of range) then the *result* SHOULD be set to VIRTIO\_RPMB\_RES\_ADDR\_FAILURE.
5. The device MUST calculate the MAC taking authentication key and frame as input, and compare this with the MAC in the request. If the two MAC's are different then VIRTIO\_RPMB\_RES\_AUTH\_FAILURE SHOULD be returned in *result*.
6. If the writer counter in the request is different from the one maintained by the device then VIRTIO\_RPMB\_RES\_COUNT\_FAILURE SHOULD be returned in *result*.
7. If the MAC and write counter comparisons are matched then the write request is considered to be authenticated. The data from the request SHOULD be written to the address indicated in the request and the write counter SHOULD be incremented by 1.
8. If the write fails then the *result* SHOULD be VIRTIO\_RPMB\_RES\_WRITE\_FAILURE.
9. If some other error occurs during the writing procedure then the *result* SHOULD be VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE.
10. The *req\_resp* value VIRTIO\_RPMB\_RESP\_DATA\_WRITE SHOULD be responded.

#### 5.12.6.1.4 Device Requirements: Device Operation: Data Read

1. If the authentication key is not yet programmed then VIRTIO\_RPMB\_RES\_NO\_AUTH\_KEY SHOULD be returned in *result*.
2. If block count is zero or greater than *max\_rd\_cnt* then VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE SHOULD be responded as *result*.
3. If there is an error in the address (out of range) then the *result* SHOULD be set to VIRTIO\_RPMB\_RES\_ADDR\_FAILURE.



4. If data fetch from addressed location inside the device fails then the *result* SHOULD be VIRTIO\_RPMB\_RES\_READ\_FAILURE.
5. If some other error occurs during the read procedure then the *result* SHOULD be VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE.
6. The device SHOULD respond with *block\_count* frames containing the data and *req\_resp* value set to VIRTIO\_RPMB\_RESP\_DATA\_READ.

#### 5.12.6.1.5 Device Requirements: Device Operation: Result Read

If the *block\_count* has not been set to 1 of VIRTIO\_RPMB\_REQ\_RESULT\_READ request then VIRTIO\_RPMB\_RES\_GENERAL\_FAILURE SHOULD be responded as *result*.

#### 5.12.6.2 Driver Requirements: Device Operation

The RPMB frames MUST not be packed by the driver. The driver MUST configure, initialize and format virtqueue for the RPMB requests received from its caller then send it to the device.

#### 5.12.6.3 Device Requirements: Device Operation

The virtio-rpmb device could be backed in a number of ways. It SHOULD keep consistent behaviors with hardware as described in paragraph 6.6.22 of eMMC. Some elements are maintained by the device:

1. The device MUST maintain a one-time programmable authentication key. It cannot be overwritten, erased or read. The key is used to authenticate the accesses when MAC is calculated. This key MUST be kept regardless of device reset or reboot.
2. The device MUST maintain a read-only monotonic write counter. It MUST be initialized to zero and added by one automatically along with successful write operation. The value cannot be reset. After the counter has reached its maximum value 0xFFFF FFFF, it will not be incremented anymore. This counter MUST be kept regardless of device reset or reboot.
3. The device MUST maintain the data for read/write via authenticated access.

## 5.13 IOMMU device

The virtio-iommu device manages Direct Memory Access (DMA) from one or more endpoints. It may act both as a proxy for physical IOMMUs managing devices assigned to the guest, and as virtual IOMMU managing emulated and paravirtualized devices.

The driver first discovers endpoints managed by the virtio-iommu device using platform specific mechanisms. It then sends requests to create virtual address spaces and virtual-to-physical mappings for these endpoints. In its simplest form, the virtio-iommu supports four request types:

1. Create a domain and attach an endpoint to it.  
attach(endpoint = 0x8, domain = 1)
2. Create a mapping between a range of guest-virtual and guest-physical address.  
map(domain = 1, virt\_start = 0x1000, virt\_end = 0x1fff, phys = 0xa000, flags = READ)

Endpoint 0x8, for example a hardware PCI endpoint with BDF 00:01.0, can now read at addresses 0x1000-0x1fff. These accesses are translated into system-physical addresses by the IOMMU.

3. Remove the mapping.  
unmap(domain = 1, virt\_start = 0x1000, virt\_end = 0x1fff)  
Any access to addresses 0x1000-0x1fff by endpoint 0x8 would now be rejected.
4. Detach the device and remove the domain.  
detach(endpoint = 0x8, domain = 1)



### 5.13.1 Device ID

23

### 5.13.2 Virtqueues

0 requestq

1 eventq

### 5.13.3 Feature bits

**VIRTIO\_IOMMU\_F\_INPUT\_RANGE (0)** Available range of virtual addresses is described in *input\_range*.

**VIRTIO\_IOMMU\_F\_DOMAIN\_RANGE (1)** The number of domains supported is described in *domain\_range*.

**VIRTIO\_IOMMU\_F\_MAP\_UNMAP (2)** Map and unmap requests are available.<sup>12</sup>

**VIRTIO\_IOMMU\_F\_BYPASS (3)** Endpoints that are not attached to a domain are in bypass mode.

**VIRTIO\_IOMMU\_F\_PROBE (4)** The PROBE request is available.

**VIRTIO\_IOMMU\_F\_MMIO (5)** The VIRTIO\_IOMMU\_MAP\_F\_MMIO flag is available.

**VIRTIO\_IOMMU\_F\_BYPASS\_CONFIG (6)** Field *bypass* of struct *virtio\_iommu\_config* determines whether endpoints that are not attached to a domain are in bypass mode. Flag *VIRTIO\_IOMMU\_ATTACH\_F\_BYPASS* determines whether endpoints that are attached to a domain are in bypass mode.

#### 5.13.3.1 Driver Requirements: Feature bits

The driver SHOULD accept any of the *VIRTIO\_IOMMU\_F\_INPUT\_RANGE*, *VIRTIO\_IOMMU\_F\_DOMAIN\_RANGE* and *VIRTIO\_IOMMU\_F\_PROBE* feature bits if offered by the device.

#### 5.13.3.2 Device Requirements: Feature bits

The device SHOULD offer feature bit *VIRTIO\_IOMMU\_F\_MAP\_UNMAP*.

The *VIRTIO\_IOMMU\_F\_BYPASS\_CONFIG* feature supersedes *VIRTIO\_IOMMU\_F\_BYPASS*. New devices SHOULD NOT offer *VIRTIO\_IOMMU\_F\_BYPASS*. Devices SHOULD NOT offer both *VIRTIO\_IOMMU\_F\_BYPASS* and *VIRTIO\_IOMMU\_F\_BYPASS\_CONFIG*.

### 5.13.4 Device configuration layout

The *page\_size\_mask* field is always present. Availability of the others all depend on feature bits described in 5.13.3.

```
struct virtio_iommu_config {
    le64 page_size_mask;
    struct virtio_iommu_range_64 {
        le64 start;
        le64 end;
    } input_range;
    struct virtio_iommu_range_32 {
        le32 start;
        le32 end;
    } domain_range;
    le32 probe_size;
    u8 bypass;
    u8 reserved[3];
};
```

<sup>12</sup>Future extensions may add different modes of operations. At the moment, only *VIRTIO\_IOMMU\_F\_MAP\_UNMAP* is supported.

#### 5.13.4.1 Driver Requirements: Device configuration layout

When the `VIRTIO_IOMMU_F_BYPASS_CONFIG` feature is negotiated, the driver MAY write to *bypass*. The driver MUST NOT write to any other device configuration field.

The driver MUST NOT write a value different than 0 or 1 to *bypass*. The driver SHOULD ignore bits 1-7 of *bypass*.

#### 5.13.4.2 Device Requirements: Device configuration layout

The device MUST set at least one bit in *page\_size\_mask*, describing the page granularity. The device MAY set more than one bit in *page\_size\_mask*.

If the device offers the `VIRTIO_IOMMU_F_BYPASS_CONFIG` feature, it MAY initialize the *bypass* field to 1. Field *bypass* SHOULD NOT change on device reset, but SHOULD be restored to its initial value on system reset.

The device MUST NOT present a value different than 0 or 1 in *bypass*.

### 5.13.5 Device initialization

When the device is reset, endpoints are not attached to any domain.

Future devices might support more modes of operation besides MAP/UNMAP. Drivers verify that devices set `VIRTIO_IOMMU_F_MAP_UNMAP` and fail gracefully if they don't.

#### 5.13.5.1 Driver Requirements: Device Initialization

The driver MUST NOT negotiate `VIRTIO_IOMMU_F_MAP_UNMAP` if it is incapable of sending `VIRTIO_IOMMU_T_MAP` and `VIRTIO_IOMMU_T_UNMAP` requests.

If the `VIRTIO_IOMMU_F_PROBE` feature is negotiated, the driver SHOULD send a `VIRTIO_IOMMU_T_PROBE` request for each endpoint before attaching the endpoint to a domain.

### 5.13.6 Device operations

Driver send requests on the request virtqueue, notifies the device and waits for the device to return the request with a status in the used ring. All requests are split in two parts: one device-readable, one device-writable.

```
struct virtio_iommu_req_head {
    u8    type;
    u8    reserved[3];
};

struct virtio_iommu_req_tail {
    u8    status;
    u8    reserved[3];
};
```

Type may be one of:

```
#define VIRTIO_IOMMU_T_ATTACH    1
#define VIRTIO_IOMMU_T_DETACH    2
#define VIRTIO_IOMMU_T_MAP      3
#define VIRTIO_IOMMU_T_UNMAP    4
#define VIRTIO_IOMMU_T_PROBE    5
```

A few general-purpose status codes are defined here.

```
/* All good! Carry on. */
#define VIRTIO_IOMMU_S_OK        0
/* Virtio communication error */
#define VIRTIO_IOMMU_S_IOERR     1
/* Unsupported request */
#define VIRTIO_IOMMU_S_UNSUPP    2
```

```

/* Internal device error */
#define VIRTIO_IOMMU_S_DEVERR      3
/* Invalid parameters */
#define VIRTIO_IOMMU_S_INVAL      4
/* Out-of-range parameters */
#define VIRTIO_IOMMU_S_RANGE      5
/* Entry not found */
#define VIRTIO_IOMMU_S_NOENT      6
/* Bad address */
#define VIRTIO_IOMMU_S_FAULT      7
/* Insufficient resources */
#define VIRTIO_IOMMU_S_NOMEM      8

```

When the device fails to parse a request, for instance if a request is too small for its type and the device cannot find the tail, then it is unable to set *status*. In that case, it returns the buffers without writing to them.

Range limits of some request fields are described in the device configuration:

- *page\_size\_mask* contains the bitmask of all page sizes that can be mapped. The least significant bit set defines the page granularity of IOMMU mappings.

The smallest page granularity supported by the IOMMU is one byte. It is legal for the driver to map one byte at a time if bit 0 of *page\_size\_mask* is set.

Other bits in *page\_size\_mask* are hints and describe larger page sizes that the IOMMU device handles efficiently. For example, when the device stores mappings using a page table tree, it may be able to describe large mappings using a few leaf entries in intermediate tables, rather than using lots of entries in the last level of the tree. Creating mappings aligned on large page sizes can improve performance since they require fewer page table and TLB entries.

- If the VIRTIO\_IOMMU\_F\_DOMAIN\_RANGE feature is offered, *domain\_range* describes the values supported in a *domain* field. If the feature is not offered, any *domain* value is valid.
- If the VIRTIO\_IOMMU\_F\_INPUT\_RANGE feature is offered, *input\_range* contains the virtual address range that the IOMMU is able to translate. Any mapping request to virtual addresses outside of this range fails.

If the feature is not offered, virtual mappings span over the whole 64-bit address space (start = 0, end = 0xffffffff ffffffff)

An endpoint is in bypass mode if:

- the VIRTIO\_IOMMU\_F\_BYPASS\_CONFIG feature is offered and:
  - config field *bypass* is 1 and the endpoint is not attached to a domain. This applies even if the driver does not accept the VIRTIO\_IOMMU\_F\_BYPASS\_CONFIG feature and the device initializes field *bypass* to 1.
  - or
  - the endpoint is attached to a domain with VIRTIO\_IOMMU\_ATTACH\_F\_BYPASS.
  - or
- the VIRTIO\_IOMMU\_F\_BYPASS feature is negotiated and the endpoint is not attached to a domain.

All accesses from an endpoint in bypass mode are allowed and translated by the IOMMU using the identity function.

### 5.13.6.1 Driver Requirements: Device operations

The driver SHOULD set field *reserved* of struct *virtio\_iommu\_req\_head* to zero and MUST ignore field *reserved* of struct *virtio\_iommu\_req\_tail*.

When a device uses a buffer without having written to it (i.e. used length is zero), the driver SHOULD interpret it as a request failure.

If the `VIRTIO_IOMMU_F_INPUT_RANGE` feature is negotiated, the driver MUST NOT send requests with *virt\_start* less than *input\_range.start* or *virt\_end* greater than *input\_range.end*.

If the `VIRTIO_IOMMU_F_DOMAIN_RANGE` feature is negotiated, the driver MUST NOT send requests with *domain* less than *domain\_range.start* or greater than *domain\_range.end*.

### 5.13.6.2 Device Requirements: Device operations

The device SHOULD set *status* to `VIRTIO_IOMMU_S_OK` if a request succeeds.

If a request *type* is not recognized, the device SHOULD NOT write the buffer and SHOULD set the used length to zero.

The device MUST ignore field *reserved* of struct `virtio_iommu_req_head` and SHOULD set field *reserved* of struct `virtio_iommu_req_tail` to zero.

The device SHOULD NOT let unattached endpoints that are not in bypass mode access the guest-physical address space.

### 5.13.6.3 ATTACH request

```
struct virtio_iommu_req_attach {
    struct virtio_iommu_req_head head;
    le32 domain;
    le32 endpoint;
    le32 flags;
    u8 reserved[4];
    struct virtio_iommu_req_tail tail;
};

#define VIRTIO_IOMMU_ATTACH_F_BYPASS (1 << 0)
```

Attach an endpoint to a domain. *domain* uniquely identifies a domain within the virtio-iommu device. If the domain doesn't exist in the device, it is created. Semantics of the *endpoint* identifier are platform specific, but the following rules apply:

- The endpoint ID uniquely identifies an endpoint from the virtio-iommu point of view. Multiple endpoints whose DMA transactions are not translated by the same virtio-iommu device can have the same endpoint ID. Endpoints whose DMA transactions may be translated by the same virtio-iommu device have different endpoint IDs.
- On some platforms, it might not be possible to completely isolate two endpoints from each other. For example on a conventional PCI bus, endpoints can snoop DMA transactions from other endpoints on the same bus. Such limitations need to be communicated in a platform specific way.

Multiple endpoints can be attached to the same domain. An endpoint can be attached to a single domain at a time. Endpoints attached to different domains are isolated from each other.

When the `VIRTIO_IOMMU_F_BYPASS_CONFIG` is negotiated, the driver can set the `VIRTIO_IOMMU_ATTACH_F_BYPASS` flag to create a bypass domain. Endpoints attached to this domain are in bypass mode.

#### 5.13.6.3.1 Driver Requirements: ATTACH request

The driver SHOULD set *reserved* to zero.

The driver SHOULD ensure that endpoints that cannot be isolated from each other are attached to the same domain.

If the domain already exists and is a bypass domain, the driver SHOULD set the `VIRTIO_IOMMU_ATTACH_F_BYPASS` flag. If the domain exists and is not a bypass domain, the driver SHOULD NOT set the `VIRTIO_IOMMU_ATTACH_F_BYPASS` flag.

#### 5.13.6.3.2 Device Requirements: ATTACH request

If the *reserved* field of an ATTACH request is not zero, the device MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If the device does not recognize a *flags* bit, it MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If the endpoint identified by *endpoint* doesn't exist, the device MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_NOENT.

If another endpoint is already attached to the domain identified by *domain*, then the device MAY attach the endpoint identified by *endpoint* to the domain. If it cannot do so, the device MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_UNSUPP.

If the domain already exists and the VIRTIO\_IOMMU\_ATTACH\_F\_BYPASS flag is not consistent with that domain, the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If the endpoint identified by *endpoint* is already attached to another domain, then the device SHOULD first detach it from that domain and attach it to the one identified by *domain*. In that case the device SHOULD behave as if the driver issued a DETACH request with this *endpoint*, followed by the ATTACH request. If the device cannot do so, it MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_UNSUPP.

If properties of the endpoint (obtained with a PROBE request) are compatible with properties of other endpoints already attached to the requested domain, then the device SHOULD attach the endpoint. Otherwise the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_UNSUPP.

A device that does not reject the request MUST attach the endpoint.

#### 5.13.6.4 DETACH request

```
struct virtio_iommu_req_detach {
    struct virtio_iommu_req_head head;
    le32 domain;
    le32 endpoint;
    u8 reserved[8];
    struct virtio_iommu_req_tail tail;
};
```

Detach an endpoint from a domain. When this request completes, the endpoint cannot access any mapping from that domain anymore. However the endpoint may then be in bypass mode and access the guest-physical address space.

After all endpoints have been successfully detached from a domain, it ceases to exist and its ID can be reused by the driver for another domain.

##### 5.13.6.4.1 Driver Requirements: DETACH request

The driver SHOULD set *reserved* to zero.

##### 5.13.6.4.2 Device Requirements: DETACH request

The device MUST ignore *reserved*.

If the endpoint identified by *endpoint* doesn't exist, then the device MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_NOENT.

If the domain identified by *domain* doesn't exist, or if the endpoint identified by *endpoint* isn't attached to this domain, then the device MAY set the request *status* to VIRTIO\_IOMMU\_S\_INVALID.

The device MUST ensure that after being detached from a domain, the endpoint cannot access any mapping from that domain.

### 5.13.6.5 MAP request

```
struct virtio_iommu_req_map {
    struct virtio_iommu_req_head head;
    le32 domain;
    le64 virt_start;
    le64 virt_end;
    le64 phys_start;
    le32 flags;
    struct virtio_iommu_req_tail tail;
};

/* Read access is allowed */
#define VIRTIO_IOMMU_MAP_F_READ (1 << 0)
/* Write access is allowed */
#define VIRTIO_IOMMU_MAP_F_WRITE (1 << 1)
/* Accesses are to memory-mapped I/O device */
#define VIRTIO_IOMMU_MAP_F_MMIO (1 << 2)
```

Map a range of virtually-contiguous addresses to a range of physically-contiguous addresses of the same size. After the request succeeds, all endpoints attached to this domain can access memory in the range  $[virt\_start; virt\_end]$  (inclusive). For example, if an endpoint accesses address  $VA \in [virt\_start; virt\_end]$ , the device (or the physical IOMMU) translates the address:  $PA = VA - virt\_start + phys\_start$ . If the access parameters are compatible with *flags* (for instance, the access is write and *flags* are `VIRTIO_IOMMU_MAP_F_READ | VIRTIO_IOMMU_MAP_F_WRITE`) then the IOMMU allows the access to reach *PA*.

The range defined by *virt\_start* and *virt\_end* should be within the limits specified by *input\_range*. Given  $phys\_end = phys\_start + virt\_end - virt\_start$ , the range defined by *phys\_start* and *phys\_end* should be within the guest-physical address space. This includes upper and lower limits, as well as any carving of guest-physical addresses for use by the host. Guest physical boundaries are set by the host in a platform specific way.

Availability and allowed combinations of *flags* depend on the underlying IOMMU architectures. `VIRTIO_IOMMU_MAP_F_READ` and `VIRTIO_IOMMU_MAP_F_WRITE` are usually implemented, although `READ` is sometimes implied by `WRITE`. In addition combinations such as "WRITE and not READ" might not be supported.

The `VIRTIO_IOMMU_MAP_F_MMIO` flag is a memory type rather than a protection flag. It is only available when the `VIRTIO_IOMMU_F_MMIO` feature has been negotiated. Accesses to the mapping are not speculated, buffered, cached, split into multiple accesses or combined with other accesses. It may be used, for example, to map Message Signaled Interrupt doorbells when a `VIRTIO_IOMMU_RESV_MEM_T_MSI` region isn't available. To trigger interrupts the endpoint performs a direct memory write to another peripheral, the IRQ chip.

This request is only available when `VIRTIO_IOMMU_F_MAP_UNMAP` has been negotiated.

#### 5.13.6.5.1 Driver Requirements: MAP request

The driver SHOULD set undefined *flags* bits to zero.

The driver SHOULD NOT send MAP requests on a bypass domain.

*virt\_end* MUST be strictly greater than *virt\_start*.

The driver SHOULD set the `VIRTIO_IOMMU_MAP_F_MMIO` flag when the physical range corresponds to memory-mapped device registers. The physical range SHOULD have a single memory type: either normal memory or memory-mapped I/O.

If it intends to allow read accesses from endpoints attached to the domain, the driver MUST set the `VIRTIO_IOMMU_MAP_F_READ` flag.

If the `VIRTIO_IOMMU_F_MMIO` feature isn't negotiated, the driver MUST NOT use the `VIRTIO_IOMMU_MAP_F_MMIO` flag.

### 5.13.6.5.2 Device Requirements: MAP request

If *virt\_start*, *phys\_start* or (*virt\_end* + 1) is not aligned on the page granularity, the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_RANGE.

If a mapping already exists in the requested range, the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If the device doesn't recognize a *flags* bit, it MUST reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If *domain* does not exist, the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_NOENT.

If the domain is a bypass domain, the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

The device MUST NOT allow writes to a range mapped without the VIRTIO\_IOMMU\_MAP\_F\_WRITE flag. However, if the underlying architecture does not support write-only mappings, the device MAY allow reads to a range mapped with VIRTIO\_IOMMU\_MAP\_F\_WRITE but not VIRTIO\_IOMMU\_MAP\_F\_READ.

### 5.13.6.6 UNMAP request

```
struct virtio_iommu_req_unmap {
    struct virtio_iommu_req_head head;
    le32  domain;
    le64  virt_start;
    le64  virt_end;
    u8    reserved[4];
    struct virtio_iommu_req_tail tail;
};
```

Unmap a range of addresses mapped with VIRTIO\_IOMMU\_T\_MAP. We define here a mapping as a virtual region created with a single MAP request. All mappings covered by the range [*virt\_start*; *virt\_end*] (inclusive) are removed.

The semantics of unmapping are specified in 5.13.6.6.1 and 5.13.6.6.2, and illustrated with the following requests, assuming each example sequence starts with a blank address space. We define two pseudocode functions `map(virt_start, virt_end) -> mapping` and `unmap(virt_start, virt_end)`.

```
(1) unmap(virt_start=0,
          virt_end=4)           -> succeeds, doesn't unmap anything

(2) a = map(virt_start=0,
            virt_end=9);
    unmap(0, 9)                 -> succeeds, unmaps a

(3) a = map(0, 4);
    b = map(5, 9);
    unmap(0, 9)                 -> succeeds, unmaps a and b

(4) a = map(0, 9);
    unmap(0, 4)                 -> fails, doesn't unmap anything

(5) a = map(0, 4);
    b = map(5, 9);
    unmap(0, 4)                 -> succeeds, unmaps a

(6) a = map(0, 4);
    unmap(0, 9)                 -> succeeds, unmaps a

(7) a = map(0, 4);
    b = map(10, 14);
    unmap(0, 14)                -> succeeds, unmaps a and b
```

As illustrated by example (4), partially removing a mapping isn't supported.

This request is only available when VIRTIO\_IOMMU\_F\_MAP\_UNMAP has been negotiated.



#### 5.13.6.6.1 Driver Requirements: UNMAP request

The driver SHOULD set the *reserved* field to zero.

The range, defined by *virt\_start* and *virt\_end*, SHOULD cover one or more contiguous mappings created with MAP requests. The range MAY spill over unmapped virtual addresses.

The first address of a range MUST either be the first address of a mapping or be outside any mapping. The last address of a range MUST either be the last address of a mapping or be outside any mapping.

The driver SHOULD NOT send UNMAP requests on a bypass domain.

#### 5.13.6.6.2 Device Requirements: UNMAP request

If the *reserved* field of an UNMAP request is not zero, the device MAY set the request *status* to `VIRTIO_IOMMU_S_INVALID`, in which case the device MAY perform the UNMAP operation.

If *domain* does not exist, the device SHOULD set the request *status* to `VIRTIO_IOMMU_S_NOENT`.

If the domain is a bypass domain, the device SHOULD reject the request and set *status* to `VIRTIO_IOMMU_S_INVALID`.

If a mapping affected by the range is not covered in its entirety by the range (the UNMAP request would split the mapping), then the device SHOULD set the request *status* to `VIRTIO_IOMMU_S_RANGE`, and SHOULD NOT remove any mapping.

If part of the range or the full range is not covered by an existing mapping, then the device SHOULD remove all mappings affected by the range and set the request *status* to `VIRTIO_IOMMU_S_OK`.

#### 5.13.6.7 PROBE request

If the `VIRTIO_IOMMU_F_PROBE` feature bit is present, the driver sends a `VIRTIO_IOMMU_T_PROBE` request for each endpoint that the virtio-iommu device manages. This probe is performed before attaching the endpoint to a domain.

```
struct virtio_iommu_req_probe {
    struct virtio_iommu_req_head head;
    /* Device-readable */
    le32 endpoint;
    u8 reserved[64];

    /* Device-writable */
    u8 properties[probe_size];
    struct virtio_iommu_req_tail tail;
};
```

**endpoint** has the same meaning as in ATTACH and DETACH requests.

**reserved** is used as padding, so that future extensions can add fields to the device-readable part.

**properties** contains a list of properties of the *endpoint*, filled by the device. The length of the *properties* field is *probe\_size* bytes. Each property is described with a struct `virtio_iommu_probe_property` header, which may be followed by a value of size *length*.

```
struct virtio_iommu_probe_property {
    le16 {
        type      : 12;
        reserved   : 4;
    };
    le16 length;
};
```

The driver allocates a buffer for the PROBE request, large enough to accommodate *probe\_size* bytes of *properties*. It writes *endpoint* and adds the buffer to the request queue. The device fills the *properties* field with a list of properties for this endpoint.

The driver parses the first property by reading *type*, then *length*. If the driver recognizes *type*, it reads and handles the rest of the property. The driver then reads the next property, that is located  $(length + 4)$  bytes after the beginning of the first one, and so on. The driver parses all properties until it reaches an empty property (*type* is 0) or the end of *properties*.

Available property types are described in section 5.13.6.8.

#### 5.13.6.7.1 Driver Requirements: PROBE request

The size of *properties* MUST be *probe\_size* bytes.

The driver SHOULD set field *reserved* of the PROBE request to zero.

If the driver doesn't recognize the *type* of a property, it SHOULD ignore the property.

The driver SHOULD NOT deduce the property length from *type*.

The driver MUST ignore a property whose *reserved* field is not zero.

If the driver ignores a property, it SHOULD continue parsing the list.

#### 5.13.6.7.2 Device Requirements: PROBE request

The device MUST ignore field *reserved* of a PROBE request.

If the endpoint identified by *endpoint* doesn't exist, then the device SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_NOENT.

If the device does not offer the VIRTIO\_IOMMU\_F\_PROBE feature, and if the driver sends a VIRTIO\_IOMMU\_T\_PROBE request, then the device SHOULD NOT write the buffer and SHOULD set the used length to zero.

The device SHOULD set field *reserved* of a property to zero.

The device MUST write the size of a property without the struct `virtio_iommu_probe_property` header, in bytes, into *length*.

When two properties follow each other, the device MUST put the second property exactly  $(length + 4)$  bytes after the beginning of the first one.

If the *properties* list is smaller than *probe\_size*, the device SHOULD NOT write any property. It SHOULD reject the request and set *status* to VIRTIO\_IOMMU\_S\_INVALID.

If the device doesn't fill all *probe\_size* bytes with properties, it SHOULD fill the remaining bytes of *properties* with zeroes.

### 5.13.6.8 PROBE properties

```
#define VIRTIO_IOMMU_PROBE_T_RESV_MEM 1
```

#### 5.13.6.8.1 Property RESV\_MEM

The RESV\_MEM property describes a chunk of reserved virtual memory. It may be used by the device to describe virtual address ranges that cannot be used by the driver, or that are special.

```
struct virtio_iommu_probe_resv_mem {
    struct virtio_iommu_probe_property head;
    u8    subtype;
    u8    reserved[3];
    le64  start;
    le64  end;
};
```

Fields *start* and *end* describe the range of reserved virtual addresses. *subtype* may be one of:

**VIRTIO\_IOMMU\_RESV\_MEM\_T\_RESERVED (0)** These virtual addresses cannot be used in a MAP requests. The region is be reserved by the device, for example, if the platform needs to setup DMA mappings of its own.

**VIRTIO\_IOMMU\_RESV\_MEM\_T\_MSI (1)** This region is a doorbell for Message Signaled Interrupts (MSIs). It is similar to VIRTIO\_IOMMU\_RESV\_MEM\_T\_RESERVED, in that the driver cannot map virtual addresses described by the property.

In addition it provides information about MSI doorbells. If the endpoint doesn't have a VIRTIO\_IOMMU\_RESV\_MEM\_T\_MSI property, then the driver creates an MMIO mapping to the doorbell of the MSI controller.

#### 5.13.6.8.1.1 Driver Requirements: Property RESV\_MEM

The driver SHOULD NOT map any virtual address described by a VIRTIO\_IOMMU\_RESV\_MEM\_T\_RESERVED or VIRTIO\_IOMMU\_RESV\_MEM\_T\_MSI property.

The driver MUST ignore *reserved*.

The driver SHOULD treat any *subtype* it doesn't recognize as if it was VIRTIO\_IOMMU\_RESV\_MEM\_T\_RESERVED.

#### 5.13.6.8.1.2 Device Requirements: Property RESV\_MEM

The device SHOULD set *reserved* to zero.

The device SHOULD NOT present more than one VIRTIO\_IOMMU\_RESV\_MEM\_T\_MSI property per endpoint.

The device SHOULD NOT present multiple RESV\_MEM properties that overlap each other for the same endpoint.

The device SHOULD reject a MAP request that overlaps a RESV\_MEM region.

The device SHOULD NOT allow accesses from the endpoint to RESV\_MEM regions to affect any other component than the endpoint and the driver.

#### 5.13.6.9 Fault reporting

The device can report translation faults and other significant asynchronous events on the event virtqueue. The driver initially populates the queue with device-writeable buffers. When the device needs to report an event, it fills a buffer and notifies the driver. The driver consumes the report and adds a new buffer to the virtqueue.

If no buffer is available, the device can either wait for one to be consumed, or drop the event.

```
struct virtio_iommu_fault {
    u8    reason;
    u8    reserved[3];
    le32  flags;
    le32  endpoint;
    le32  reserved1;
    le64  address;
};

#define VIRTIO_IOMMU_FAULT_F_READ    (1 << 0)
#define VIRTIO_IOMMU_FAULT_F_WRITE  (1 << 1)
#define VIRTIO_IOMMU_FAULT_F_ADDRESS (1 << 8)
```

**reason** The reason for this report. It may have the following values:

**VIRTIO\_IOMMU\_FAULT\_R\_UNKNOWN (0)** An internal error happened, or an error that cannot be described with the following reasons.

**VIRTIO\_IOMMU\_FAULT\_R\_DOMAIN (1)** The endpoint attempted to access *address* without being attached to a domain.

**VIRTIO\_IOMMU\_FAULT\_R\_MAPPING (2)** The endpoint attempted to access *address*, which wasn't mapped in the domain or didn't have the correct protection flags.

**flags** Information about the fault context.

**endpoint** The endpoint causing the fault.

**reserved and reserved1** Should be zero.

**address** If VIRTIO\_IOMMU\_FAULT\_F\_ADDRESS is set, the address causing the fault.

When the fault is reported by a physical IOMMU, the fault reasons may not match exactly the reason of the original fault report. The device does its best to find the closest match.

If the device encounters an internal error that wasn't caused by a specific endpoint, it is unlikely that the driver would be able to do anything else than print the fault and stop using the device, so reporting the fault on the event queue isn't useful. In that case, we recommend using the DEVICE\_NEEDS\_RESET status bit.

#### 5.13.6.9.1 Driver Requirements: Fault reporting

If the *reserved* field is not zero, the driver **MUST** ignore the fault report.

The driver **MUST** ignore *reserved1*.

The driver **MUST** ignore undefined *flags*.

If the driver doesn't recognize *reason*, it **SHOULD** treat the fault as if it was VIRTIO\_IOMMU\_FAULT\_R\_UNKNOWN.

#### 5.13.6.9.2 Device Requirements: Fault reporting

The device **SHOULD** set *reserved* and *reserved1* to zero.

The device **SHOULD** set undefined *flags* to zero.

The device **SHOULD** write a valid endpoint ID in *endpoint*.

The device **MAY** omit setting VIRTIO\_IOMMU\_FAULT\_F\_ADDRESS and writing *address* in any fault report, regardless of the *reason*.

If a buffer is too small to contain the fault report<sup>13</sup>, the device **SHOULD NOT** use multiple buffers to describe it. The device **MAY** fall back to using an older fault report format that fits in the buffer.

## 5.14 Sound Device

The virtio sound card is a virtual audio device supporting input and output PCM streams.

A device is managed by control requests and can send various notifications through dedicated queues. A driver can transmit PCM frames using message-based transport or shared memory.

A small part of the specification reuses existing layouts and values from the High Definition Audio specification (HDA). It allows to provide the same functionality and assist in two possible cases:

1. implementation of a universal sound driver,
2. implementation of a sound driver as part of the High Definition Audio subsystem.

### 5.14.1 Device ID

25

---

<sup>13</sup>This would happen for example if the device implements a more recent version of this specification, whose fault report contains additional fields.

## 5.14.2 Virtqueues

- 0 controlq
- 1 eventq
- 2 txq
- 3 rxq

The control queue is used for sending control messages from the driver to the device.

The event queue is used for sending notifications from the device to the driver.

The tx queue is used to send PCM frames for output streams.

The rx queue is used to receive PCM frames from input streams.

## 5.14.3 Feature Bits

None currently defined.

## 5.14.4 Device Configuration Layout

```
struct virtio_snd_config {  
    le32 jacks;  
    le32 streams;  
    le32 chmaps;  
};
```

A configuration space contains the following fields:

**jacks** (driver-read-only) indicates a total number of all available jacks.

**streams** (driver-read-only) indicates a total number of all available PCM streams.

**chmaps** (driver-read-only) indicates a total number of all available channel maps.

## 5.14.5 Device Initialization

1. Configure the control, event, tx and rx queues.
2. Read the *jacks* field and send a control request to query information about the available jacks.
3. Read the *streams* field and send a control request to query information about the available PCM streams.
4. Read the *chmaps* field and send a control request to query information about the available channel maps.
5. Populate the event queue with empty buffers.

### 5.14.5.1 Driver Requirements: Device Initialization

- The driver **MUST** populate the event queue with empty buffers of at least the struct `virtio_snd_event` size.
- The driver **MUST NOT** put a device-readable buffers in the event queue.

## 5.14.6 Device Operation

All control messages are placed into the control queue and all notifications are placed into the event queue. They use the following layout structure and definitions:

```

enum {
    /* jack control request types */
    VIRTIO_SND_R_JACK_INFO = 1,
    VIRTIO_SND_R_JACK_REMAP,

    /* PCM control request types */
    VIRTIO_SND_R_PCM_INFO = 0x0100,
    VIRTIO_SND_R_PCM_SET_PARAMS,
    VIRTIO_SND_R_PCM_PREPARE,
    VIRTIO_SND_R_PCM_RELEASE,
    VIRTIO_SND_R_PCM_START,
    VIRTIO_SND_R_PCM_STOP,

    /* channel map control request types */
    VIRTIO_SND_R_CHMAP_INFO = 0x0200,

    /* jack event types */
    VIRTIO_SND_EVT_JACK_CONNECTED = 0x1000,
    VIRTIO_SND_EVT_JACK_DISCONNECTED,

    /* PCM event types */
    VIRTIO_SND_EVT_PCM_PERIOD_ELAPSED = 0x1100,
    VIRTIO_SND_EVT_PCM_XRUN,

    /* common status codes */
    VIRTIO_SND_S_OK = 0x8000,
    VIRTIO_SND_S_BAD_MSG,
    VIRTIO_SND_S_NOT_SUPP,
    VIRTIO_SND_S_IO_ERR
};

/* a common header */
struct virtio_snd_hdr {
    le32 code;
};

/* an event notification */
struct virtio_snd_event {
    struct virtio_snd_hdr hdr;
    le32 data;
};

```

A generic control message consists of a request part and a response part.

A request part has, or consists of, a common header containing the following device-readable field:

**code** specifies a device request type (VIRTIO\_SND\_R\_\*).

A response part has, or consists of, a common header containing the following device-writable field:

**code** indicates a device request status (VIRTIO\_SND\_S\_\*).

The status field can take one of the following values:

- VIRTIO\_SND\_S\_OK - success.
- VIRTIO\_SND\_S\_BAD\_MSG - a control message is malformed or contains invalid parameters.
- VIRTIO\_SND\_S\_NOT\_SUPP - requested operation or parameters are not supported.
- VIRTIO\_SND\_S\_IO\_ERR - an I/O error occurred.

The request part may be followed by an additional device-readable payload, and the response part may be followed by an additional device-writable payload.

An event notification contains the following device-writable fields:

**hdr** indicates an event type (VIRTIO\_SND\_EVT\_\*).

**data** indicates an optional event data.

For all entities involved in the exchange of audio data, the device uses one of the following data flow directions:

```
enum {
    VIRTIO_SND_D_OUTPUT = 0,
    VIRTIO_SND_D_INPUT
};
```

#### 5.14.6.1 Item Information Request

A special control message is used to request information about any kind of configuration item. The request part uses the following structure definition:

```
struct virtio_snd_query_info {
    struct virtio_snd_hdr hdr;
    le32 start_id;
    le32 count;
    le32 size;
};
```

The request contains the following device-readable fields:

**hdr** specifies a particular item request type (VIRTIO\_SND\_R\_\*\_INFO).

**start\_id** specifies the starting identifier for the item (the range of available identifiers is limited by the total number of particular items that is indicated in the device configuration space).

**count** specifies the number of items for which information is requested (the total number of particular items is indicated in the device configuration space).

**size** specifies the size of the structure containing information for one item (used for backward compatibility).

The response consists of the `virtio_snd_hdr` structure (contains the request status code), followed by the device-writable information structures of the item. Each information structure begins with the following common header:

```
struct virtio_snd_info {
    le32 hda_fn_nid;
};
```

The header contains the following field:

**hda\_fn\_nid** indicates a function group node identifier (see [HDA](#), section 7.1.2). This field can be used to link together different types of resources (e.g. jacks with streams and channel maps with streams).

#### 5.14.6.2 Driver Requirements: Item Information Request

- The driver MUST NOT set *start\_id* and *count* such that *start\_id* + *count* is greater than the total number of particular items that is indicated in the device configuration space.
- The driver MUST provide a buffer of `sizeof(struct virtio_snd_hdr) + count * size` bytes for the response.

#### 5.14.6.3 Relationships with the High Definition Audio Specification

The High Definition Audio specification introduces the codec as part of the hardware that implements some of the functionality. The codec architecture and capabilities are described by tree structure of special nodes each of which is either a function module or a function group (see [HDA](#) for details).

The virtio sound specification assumes that a single codec is implemented in the device. Function module nodes are simulated by item information structures, and function group nodes are simulated by the *hda\_fn\_nid* field in each such structure.

#### 5.14.6.4 Jack Control Messages

A jack control request has, or consists of, a common header with the following layout structure:

```
struct virtio_snd_jack_hdr {
    struct virtio_snd_hdr hdr;
    le32 jack_id;
};
```

The header consists of the following device-readable fields:

**hdr** specifies a request type (VIRTIO\_SND\_R\_JACK\_\*).

**jack\_id** specifies a jack identifier from 0 to *jacks* - 1.

##### 5.14.6.4.1 VIRTIO\_SND\_R\_JACK\_INFO

Query information about the available jacks.

The request consists of the virtio\_snd\_query\_info structure (see [Item Information Request](#)). The response consists of the virtio\_snd\_hdr structure, followed by the following jack information structures:

```
/* supported jack features */
enum {
    VIRTIO_SND_JACK_F_REMAP = 0
};

struct virtio_snd_jack_info {
    struct virtio_snd_info hdr;
    le32 features; /* 1 <= VIRTIO_SND_JACK_F_XXX */
    le32 hda_reg_defconf;
    le32 hda_reg_caps;
    u8 connected;

    u8 padding[7];
};
```

The structure contains the following device-writable fields:

**features** indicates a supported feature bit map:

- VIRTIO\_SND\_JACK\_F\_REMAP - jack remapping support.

**hda\_reg\_defconf** indicates a pin default configuration value (see [HDA](#), section 7.3.3.31).

**hda\_reg\_caps** indicates a pin capabilities value (see [HDA](#), section 7.3.4.9).

**connected** indicates the current jack connection status (1 - connected, 0 - disconnected).

##### 5.14.6.4.1.1 Device Requirements: Jack Information

- The device MUST NOT set undefined feature values.
- The device MUST initialize the *padding* bytes to 0.

##### 5.14.6.4.2 VIRTIO\_SND\_R\_JACK\_REMAP

If the VIRTIO\_SND\_JACK\_F\_REMAP feature bit is set in the jack information, then the driver can send a control request to change the association and/or sequence number for the specified jack ID.

The request uses the following structure and layout definitions:

```
struct virtio_snd_jack_remap {
    struct virtio_snd_jack_hdr hdr; /* .code = VIRTIO_SND_R_JACK_REMAP */
    le32 association;
    le32 sequence;
};
```



The request contains the following device-readable fields:

**association** specifies the selected association number.

**sequence** specifies the selected sequence number.

#### 5.14.6.5 Jack Notifications

Jack notifications consist of a `virtio_snd_event` structure, which has the following device-writable fields:

**hdr** indicates a jack event type:

- `VIRTIO_SND_EVT_JACK_CONNECTED` - an external device has been connected to the jack.
- `VIRTIO_SND_EVT_JACK_DISCONNECTED` - an external device has been disconnected from the jack.

**data** indicates a jack identifier from 0 to *jacks* - 1.

#### 5.14.6.6 PCM Control Messages

A PCM control request has, or consists of, a common header with the following layout structure:

```
struct virtio_snd_pcm_hdr {  
    struct virtio_snd_hdr hdr;  
    le32 stream_id;  
};
```

The header consists of the following device-readable fields:

**hdr** specifies request type (`VIRTIO_SND_R_PCM_*`).

**stream\_id** specifies a PCM stream identifier from 0 to *streams* - 1.

##### 5.14.6.6.1 PCM Command Lifecycle

A PCM stream has the following command lifecycle:

###### 1. SET PARAMETERS

The driver negotiates the stream parameters (format, transport, etc) with the device.

Possible valid transitions: set parameters, prepare.

###### 2. PREPARE

The device prepares the stream (allocates resources, etc).

Possible valid transitions: set parameters, prepare, start, release.

###### 3. Output only: the driver transfers data for pre-buffing.

###### 4. START

The device starts the stream (unmute, putting into running state, etc).

Possible valid transitions: stop.

###### 5. The driver transfers data to/from the stream.

###### 6. STOP

The device stops the stream (mute, putting into non-running state, etc).

Possible valid transitions: start, release.

###### 7. RELEASE

The device releases the stream (frees resources, etc).

Possible valid transitions: set parameters, prepare.

#### 5.14.6.6.2 VIRTIO\_SND\_R\_PCM\_INFO

Query information about the available streams.

The request consists of the `virtio_snd_query_info` structure (see [Item Information Request](#)). The response consists of the `virtio_snd_hdr` structure, followed by the following stream information structures:

```
/* supported PCM stream features */
enum {
    VIRTIO_SND_PCM_F_SHMEM_HOST = 0,
    VIRTIO_SND_PCM_F_SHMEM_GUEST,
    VIRTIO_SND_PCM_F_MSG_POLLING,
    VIRTIO_SND_PCM_F_EVT_SHMEM_PERIODS,
    VIRTIO_SND_PCM_F_EVT_XRUNS
};

/* supported PCM sample formats */
enum {
    /* analog formats (width / physical width) */
    VIRTIO_SND_PCM_FMT_IMA_ADPCM = 0, /* 4 / 4 bits */
    VIRTIO_SND_PCM_FMT_MU_LAW, /* 8 / 8 bits */
    VIRTIO_SND_PCM_FMT_A_LAW, /* 8 / 8 bits */
    VIRTIO_SND_PCM_FMT_S8, /* 8 / 8 bits */
    VIRTIO_SND_PCM_FMT_U8, /* 8 / 8 bits */
    VIRTIO_SND_PCM_FMT_S16, /* 16 / 16 bits */
    VIRTIO_SND_PCM_FMT_U16, /* 16 / 16 bits */
    VIRTIO_SND_PCM_FMT_S18_3, /* 18 / 24 bits */
    VIRTIO_SND_PCM_FMT_U18_3, /* 18 / 24 bits */
    VIRTIO_SND_PCM_FMT_S20_3, /* 20 / 24 bits */
    VIRTIO_SND_PCM_FMT_U20_3, /* 20 / 24 bits */
    VIRTIO_SND_PCM_FMT_S24_3, /* 24 / 24 bits */
    VIRTIO_SND_PCM_FMT_U24_3, /* 24 / 24 bits */
    VIRTIO_SND_PCM_FMT_S20, /* 20 / 32 bits */
    VIRTIO_SND_PCM_FMT_U20, /* 20 / 32 bits */
    VIRTIO_SND_PCM_FMT_S24, /* 24 / 32 bits */
    VIRTIO_SND_PCM_FMT_U24, /* 24 / 32 bits */
    VIRTIO_SND_PCM_FMT_S32, /* 32 / 32 bits */
    VIRTIO_SND_PCM_FMT_U32, /* 32 / 32 bits */
    VIRTIO_SND_PCM_FMT_FLOAT, /* 32 / 32 bits */
    VIRTIO_SND_PCM_FMT_FLOAT64, /* 64 / 64 bits */
    /* digital formats (width / physical width) */
    VIRTIO_SND_PCM_FMT_DSD_U8, /* 8 / 8 bits */
    VIRTIO_SND_PCM_FMT_DSD_U16, /* 16 / 16 bits */
    VIRTIO_SND_PCM_FMT_DSD_U32, /* 32 / 32 bits */
    VIRTIO_SND_PCM_FMT_IEC958_SUBFRAME /* 32 / 32 bits */
};

/* supported PCM frame rates */
enum {
    VIRTIO_SND_PCM_RATE_5512 = 0,
    VIRTIO_SND_PCM_RATE_8000,
    VIRTIO_SND_PCM_RATE_11025,
    VIRTIO_SND_PCM_RATE_16000,
    VIRTIO_SND_PCM_RATE_22050,
    VIRTIO_SND_PCM_RATE_32000,
    VIRTIO_SND_PCM_RATE_44100,
    VIRTIO_SND_PCM_RATE_48000,
    VIRTIO_SND_PCM_RATE_64000,
    VIRTIO_SND_PCM_RATE_88200,
    VIRTIO_SND_PCM_RATE_96000,
    VIRTIO_SND_PCM_RATE_176400,
    VIRTIO_SND_PCM_RATE_192000,
    VIRTIO_SND_PCM_RATE_384000
};

struct virtio_snd_pcm_info {
    struct virtio_snd_info hdr;
    le32 features; /* 1 << VIRTIO_SND_PCM_F_XXX */
    le64 formats; /* 1 << VIRTIO_SND_PCM_FMT_XXX */
    le64 rates; /* 1 << VIRTIO_SND_PCM_RATE_XXX */
    u8 direction;
    u8 channels_min;
```

```

    u8 channels_max;

    u8 padding[5];
};

```

The structure contains the following device-writable fields:

**features** indicates a bit map of the supported features, which can be negotiated by setting the stream parameters:

- VIRTIO\_SND\_PCM\_F\_SHMEM\_HOST - supports sharing a host memory with a guest,
- VIRTIO\_SND\_PCM\_F\_SHMEM\_GUEST - supports sharing a guest memory with a host,
- VIRTIO\_SND\_PCM\_F\_MSG\_POLLING - supports polling mode for message-based transport,
- VIRTIO\_SND\_PCM\_F\_EVT\_SHMEM\_PERIODS - supports elapsed period notifications for shared memory transport,
- VIRTIO\_SND\_PCM\_F\_EVT\_XRUNS - supports underrun/overflow notifications.

**formats** indicates a supported sample format bit map.

**rates** indicates a supported frame rate bit map.

**direction** indicates the direction of data flow (VIRTIO\_SND\_D\_\*).

**channels\_min** indicates a minimum number of supported channels.

**channels\_max** indicates a maximum number of supported channels.

Only interleaved samples are supported.

#### 5.14.6.6.2.1 Device Requirements: Stream Information

- The device MUST NOT set undefined feature, format, rate and direction values.
- The device MUST initialize the *padding* bytes to 0.

#### 5.14.6.6.3 VIRTIO\_SND\_R\_PCM\_SET\_PARAMS

Set selected stream parameters for the specified stream ID.

The request uses the following structure and layout definitions:

```

struct virtio_snd_pcm_set_params {
    struct virtio_snd_pcm_hdr hdr; /* .code = VIRTIO_SND_R_PCM_SET_PARAMS */
    le32 buffer_bytes;
    le32 period_bytes;
    le32 features; /* 1 << VIRTIO_SND_PCM_F_XXX */
    u8 channels;
    u8 format;
    u8 rate;

    u8 padding;
};

```

The request contains the following device-readable fields:

**buffer\_bytes** specifies the size of the hardware buffer used by the driver.

**period\_bytes** specifies the size of the hardware period used by the driver.

**features** specifies a selected feature bit map:

- VIRTIO\_SND\_PCM\_F\_SHMEM\_HOST - use shared memory allocated by the host (is a placeholder and MUST NOT be selected at the moment),
- VIRTIO\_SND\_PCM\_F\_SHMEM\_GUEST - use shared memory allocated by the guest (is a placeholder and MUST NOT be selected at the moment),

- VIRTIO\_SND\_PCM\_F\_MSG\_POLLING - suppress available buffer notifications for tx and rx queues (device should poll virtqueue),
- VIRTIO\_SND\_PCM\_F\_EVT\_SHMEM\_PERIODS - enable elapsed period notifications for shared memory transport,
- VIRTIO\_SND\_PCM\_F\_EVT\_XRUNS - enable underrun/overrun notifications.

**channels** specifies a selected number of channels.

**format** specifies a selected sample format (VIRTIO\_SND\_PCM\_FMT\_\*).

**rate** specifies a selected frame rate (VIRTIO\_SND\_PCM\_RATE\_\*).

#### 5.14.6.6.3.1 Device Requirements: Stream Parameters

- If the device has an intermediate buffer, its size MUST be no less than the specified *buffer\_bytes* value.

#### 5.14.6.6.3.2 Driver Requirements: Stream Parameters

- *period\_bytes* MUST be a divider *buffer\_bytes*, i.e.  $buffer\_bytes \% period\_bytes == 0$ .
- The driver MUST NOT set undefined feature, format and rate values.
- The driver MUST NOT set the feature, format, and rate that are not specified in the stream configuration.
- The driver MUST NOT set the *channels* value as less than the *channels\_min* or greater than the *channels\_max* values specified in the stream configuration.
- The driver MUST NOT set the VIRTIO\_SND\_PCM\_F\_SHMEM\_HOST and VIRTIO\_SND\_PCM\_F\_SHMEM\_GUEST bits at the same time.
- The driver MUST initialize the *padding* byte to 0.
- If the bits associated with the shared memory are not selected, the driver MUST use the tx and rx queues for data transfer (see [PCM I/O Messages](#)).

#### 5.14.6.6.4 VIRTIO\_SND\_R\_PCM\_PREPARE

Prepare a stream with specified stream ID.

#### 5.14.6.6.5 VIRTIO\_SND\_R\_PCM\_RELEASE

Release a stream with specified stream ID.

#### 5.14.6.6.5.1 Device Requirements: Stream Release

- The device MUST complete all pending I/O messages for the specified stream ID.
- The device MUST NOT complete the control request while there are pending I/O messages for the specified stream ID.

#### 5.14.6.6.6 VIRTIO\_SND\_R\_PCM\_START

Start a stream with specified stream ID.

#### 5.14.6.6.7 VIRTIO\_SND\_R\_PCM\_STOP

Stop a stream with specified stream ID.

### 5.14.6.7 PCM Notifications

The device can announce support for different PCM events using feature bits in the stream information structure. To enable notifications, the driver must negotiate these features using the set stream parameters request (see 5.14.6.6.3).

PCM stream notifications consist of a `virtio_snd_event` structure, which has the following device-writable fields:

**hdr** indicates a PCM stream event type:

- `VIRTIO_SND_EVT_PCM_PERIOD_ELAPSED` - a hardware buffer period has elapsed, the period size is controlled using the `period_bytes` field.
- `VIRTIO_SND_EVT_PCM_XRUN` - an underflow for the output stream or an overflow for the input stream has occurred.

**data** indicates a PCM stream identifier from 0 to `streams - 1`.

### 5.14.6.8 PCM I/O Messages

An I/O message consists of the header part, followed by the buffer, and then the status part.

```
/* an I/O header */
struct virtio_snd_pcm_xfer {
    le32 stream_id;
};

/* an I/O status */
struct virtio_snd_pcm_status {
    le32 status;
    le32 latency_bytes;
};
```

The header part consists of the following device-readable field:

**stream\_id** specifies a PCM stream identifier from 0 to `streams - 1`.

The status part consists of the following device-writable fields:

**status** contains `VIRTIO_SND_S_OK` if an operation is successful, and `VIRTIO_SND_S_IO_ERR` otherwise.

**latency\_bytes** indicates the current device latency.

Since all buffers in the queue (with one exception) should be of the size `period_bytes`, the completion of such an I/O request can be considered an elapsed period notification.

#### 5.14.6.8.1 Output Stream

In case of an output stream, the header is followed by a device-readable buffer containing PCM frames for writing to the device. All messages are placed into the tx queue.

##### 5.14.6.8.1.1 Device Requirements: Output Stream

- The device MUST NOT complete the I/O request until the buffer is totally consumed.

##### 5.14.6.8.1.2 Driver Requirements: Output Stream

- The driver SHOULD populate the tx queue with `period_bytes` sized buffers. The only exception is the end of the stream.
- The driver MUST NOT place device-writable buffers into the tx queue.

### 5.14.6.8.2 Input Stream

In case of an input stream, the header is followed by a device-writable buffer being populated with PCM frames from the device. All messages are placed into the rx queue.

A used descriptor specifies the length of the buffer that was written by the device. It should be noted that the length value contains the size of the `virtio_snd_pcm_status` structure plus the size of the recorded frames.

#### 5.14.6.8.2.1 Device Requirements: Input Stream

- The device **MUST NOT** complete the I/O request until the buffer is full. The only exception is the end of the stream.

#### 5.14.6.8.2.2 Driver Requirements: Input Stream

- The driver **SHOULD** populate the rx queue with *period\_bytes* sized empty buffers before starting the stream.
- The driver **MUST NOT** place device-readable buffers into the rx queue.

### 5.14.6.9 Channel Map Control Messages

A device can provide one or more channel maps assigned to all streams with the same data flow direction in the same function group.

#### 5.14.6.9.1 VIRTIO\_SND\_R\_CHMAP\_INFO

Query information about the available channel maps.

The request consists of the `virtio_snd_query_info` structure (see [Item Information Request](#)). The response consists of the `virtio_snd_hdr` structure, followed by the following channel map information structures:

```
/* standard channel position definition */
enum {
    VIRTIO_SND_CHMAP_NONE = 0, /* undefined */
    VIRTIO_SND_CHMAP_NA, /* silent */
    VIRTIO_SND_CHMAP_MONO, /* mono stream */
    VIRTIO_SND_CHMAP_FL, /* front left */
    VIRTIO_SND_CHMAP_FR, /* front right */
    VIRTIO_SND_CHMAP_RL, /* rear left */
    VIRTIO_SND_CHMAP_RR, /* rear right */
    VIRTIO_SND_CHMAP_FC, /* front center */
    VIRTIO_SND_CHMAP_LFE, /* low frequency (LFE) */
    VIRTIO_SND_CHMAP_SL, /* side left */
    VIRTIO_SND_CHMAP_SR, /* side right */
    VIRTIO_SND_CHMAP_RC, /* rear center */
    VIRTIO_SND_CHMAP_FLC, /* front left center */
    VIRTIO_SND_CHMAP_FRC, /* front right center */
    VIRTIO_SND_CHMAP_RLC, /* rear left center */
    VIRTIO_SND_CHMAP_RRC, /* rear right center */
    VIRTIO_SND_CHMAP_FLW, /* front left wide */
    VIRTIO_SND_CHMAP_FRW, /* front right wide */
    VIRTIO_SND_CHMAP_FLH, /* front left high */
    VIRTIO_SND_CHMAP_FCH, /* front center high */
    VIRTIO_SND_CHMAP_FRH, /* front right high */
    VIRTIO_SND_CHMAP_TC, /* top center */
    VIRTIO_SND_CHMAP_TFL, /* top front left */
    VIRTIO_SND_CHMAP_TFR, /* top front right */
    VIRTIO_SND_CHMAP_TFC, /* top front center */
    VIRTIO_SND_CHMAP_TRL, /* top rear left */
    VIRTIO_SND_CHMAP_TRR, /* top rear right */
    VIRTIO_SND_CHMAP_TRC, /* top rear center */
    VIRTIO_SND_CHMAP_TFLC, /* top front left center */
    VIRTIO_SND_CHMAP_TFRC, /* top front right center */
    VIRTIO_SND_CHMAP_TSL, /* top side left */
    VIRTIO_SND_CHMAP_TSR, /* top side right */
    VIRTIO_SND_CHMAP_LLFE, /* left LFE */
}
```

```

    VIRTIO_SND_CHMAP_RLFE,      /* right LFE */
    VIRTIO_SND_CHMAP_BC,       /* bottom center */
    VIRTIO_SND_CHMAP_BLC,      /* bottom left center */
    VIRTIO_SND_CHMAP_BRC       /* bottom right center */
};

/* maximum possible number of channels */
#define VIRTIO_SND_CHMAP_MAX_SIZE 18

struct virtio_snd_chmap_info {
    struct virtio_snd_info hdr;
    u8 direction;
    u8 channels;
    u8 positions[VIRTIO_SND_CHMAP_MAX_SIZE];
};

```

The structure contains the following device-writable fields:

**direction** indicates the direction of data flow (VIRTIO\_SND\_D\_\*).

**channels** indicates the number of valid channel position values.

**positions** indicates channel position values (VIRTIO\_SND\_CHMAP\_\*).

#### 5.14.6.9.1.1 Device Requirements: Channel Map Information

- The device MUST NOT set undefined direction values.
- The device MUST NOT set the channels value to more than VIRTIO\_SND\_CHMAP\_MAX\_SIZE.

## 5.15 Memory Device

The virtio memory device provides and manages a memory region in guest physical address space. This memory region is partitioned into memory blocks of fixed size that can either be in the state plugged or unplugged. Once plugged, a memory block can be used like ordinary RAM. The driver selects memory blocks to (un)plug and requests the device to perform the (un)plug.

The device requests the driver to plug a certain amount of memory, by setting the *requested\_size* in the device configuration, which can change at runtime. It is up to the device driver to fulfill this request by (un)plugging memory blocks. Once the *plugged\_size* is greater or equal to the *requested\_size*, requests to plug memory blocks will be rejected by the device.

The device-managed memory region is split into two parts, the usable region and the unusable region. All memory blocks in the unusable region are unplugged and requests to plug them will be rejected. The device will grow the usable region to fit the *requested\_size*. Usually, the usable region is bigger than the *requested\_size* of the device, to give the driver some flexibility when selecting memory blocks to plug.

On initial start, and after a system reset, all memory blocks are unplugged. In corner cases, memory blocks might still be plugged after a system reset, and the driver usually requests to unplug all memory while initializing, before starting to select memory blocks to plug.

The device-managed memory region is not exposed as RAM via other firmware / hw interfaces (e.g., e820 on x86). The driver is responsible for deciding how plugged memory blocks will be used. A common use case is to expose plugged memory blocks to the operating system as system RAM, available for the page allocator.

Some platforms provide memory properties for system RAM that are usually queried and modified using special CPU instructions. Memory properties might be implicitly queried or modified on memory access. Memory properties can include advanced memory protection, access and change indication, or memory usage indication relevant in virtualized environments. <sup>14</sup> The device provides the exact same properties

<sup>14</sup>For example, s390x provides storage keys for each 4 KiB page and may, depending on the configuration, provide storage attributes for each 4 KiB page.

with the exact same semantics for plugged device memory as available for comparable RAM in the same configuration.

### 5.15.1 Device ID

24

### 5.15.2 Virtqueues

0 guest-request

### 5.15.3 Feature bits

**VIRTIO\_MEM\_F\_ACPI\_PXM (0)** The field *node\_id* in the device configuration is valid and corresponds to an ACPI PXM.

**VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE (1)** The driver is not allowed to access unplugged memory.<sup>15</sup>

### 5.15.4 Device configuration layout

All fields of this configuration are always available and read-only for the driver.

```
struct virtio_mem_config {
    le64 block_size;
    le16 node_id;
    le8 padding[6];
    le64 addr;
    le64 region_size;
    le64 usable_region_size;
    le64 plugged_size;
    le64 requested_size;
};
```

**block\_size** is the size and the alignment in bytes of a memory block. Cannot change.

**node\_id** has no meaning without VIRTIO\_MEM\_F\_ACPI\_PXM. With VIRTIO\_MEM\_F\_ACPI\_PXM, this field is valid and corresponds to an ACPI PXM. Cannot change.

**padding** has no meaning and is reserved for future use.

**addr** is the guest physical address of the start of the device-managed memory region in bytes. Cannot change.

**region\_size** is the size of device-managed memory region in bytes. Cannot change.

**usable\_region\_size** is the size of the usable device-managed memory region. Can grow up to *region\_size*. Can only shrink due to VIRTIO\_MEM\_REQ\_UNPLUG\_ALL requests.

**plugged\_size** is the amount of plugged memory in bytes within the usable device-managed memory region.

**requested\_size** is the requested amount of plugged memory within the usable device-managed memory region.

#### 5.15.4.1 Driver Requirements: Device configuration layout

The driver MUST NOT write to device configuration fields.

The driver MUST ignore the value of *padding*.

The driver MUST ignore the value of *node\_id* without VIRTIO\_MEM\_F\_ACPI\_PXM.

<sup>15</sup>On platforms with memory properties that might get modified implicitly on memory access, this feature is expected to be offered by the device.



#### 5.15.4.2 Device Requirements: Device configuration layout

The device MAY change *usable\_region\_size* and *requested\_size*.

The device MUST NOT change *block\_size*, *node\_id*, *addr*, and *region\_size*, except during a system reset.

The device MUST change *plugged\_size* to reflect the size of plugged memory blocks.

The device MUST set *usable\_region\_size* to *requested\_size* or greater.

The device MUST set *block\_size* to a power of two.

The device MUST set *addr*, *region\_size*, *usable\_region\_size*, *plugged\_size*, *requested\_size* to multiples of *block\_size*.

The device MUST set *region\_size* to 0 or greater.

The device MUST NOT shrink *usable\_region\_size*, except when processing an UNPLUG ALL request, or during a system reset.

The device MUST send a configuration update notification when changing *usable\_region\_size* or *requested\_size*, except when processing an UNPLUG ALL request.

The device SHOULD NOT send a configuration update notification when changing *plugged\_size*.

The device MAY send a configuration update notification even if nothing changed.

#### 5.15.5 Device Initialization

On initialization, the driver first discovers the device's virtqueues. It then reads the device configuration.

In case the driver detects that the device still has memory plugged (*plugged\_size* in the device configuration is greater than 0), the driver will either try to re-initialize by issuing STATE requests, or try to unplug all memory before continuing. Special handling might be necessary in case some plugged memory might still be relevant (e.g., system dump, memory still in use after unloading the driver).

##### 5.15.5.1 Driver Requirements: Device Initialization

The driver SHOULD accept VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE if it is offered and the driver supports it.

The driver SHOULD issue UNPLUG ALL requests until successful if the device still has memory plugged and the plugged memory is not in use.

##### 5.15.5.2 Device Requirements: Device Initialization

A device MAY fail to operate further if VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE is not accepted.

The device MUST NOT change the state of memory blocks during device reset.

The device MUST NOT modify memory or memory properties of plugged memory blocks during device reset.

#### 5.15.6 Device Operation

The device notifies the driver about the amount of memory the device wants the driver to consume via the device. These resize requests from the device are communicated via the *requested\_size* in the device configuration. The driver will react by requesting to (un)plug specific memory blocks, to make the *plugged\_size* match the *requested\_size* as close as possible.

The driver sends requests to the device on the guest-request virtqueue, notifies the device, and waits for the device to respond. Requests have a common header, defining the request type, followed by request-specific data. All requests are 24 bytes long and have the layout:

```

struct virtio_mem_req {
    le16 type;
    le16 padding[3];

    union {
        struct virtio_mem_req_plug plug;
        struct virtio_mem_req_unplug unplug;
        struct virtio_mem_req_state state;
    } u;
}

```

Possible request types are:

```

#define VIRTIO_MEM_REQ_PLUG          0
#define VIRTIO_MEM_REQ_UNPLUG        1
#define VIRTIO_MEM_REQ_UNPLUG_ALL    2
#define VIRTIO_MEM_REQ_STATE         3

```

Responses have a common header, defining the response type, followed by request-specific data. All responses are 10 bytes long and have the layout:

```

struct virtio_mem_resp {
    le16 type;
    le16 padding[3];

    union {
        struct virtio_mem_resp_state state;
    } u;
}

```

Possible response types, in general, are:

```

#define VIRTIO_MEM_RESP_ACK          0
#define VIRTIO_MEM_RESP_NACK         1
#define VIRTIO_MEM_RESP_BUSY         2
#define VIRTIO_MEM_RESP_ERROR        3

```

### 5.15.6.1 Driver Requirements: Device Operation

The driver MUST NOT write memory or modify memory properties of unplugged memory blocks.

The driver MUST NOT read memory or query memory properties of unplugged memory blocks outside *usable\_region\_size*.

The driver MUST NOT read memory or query memory properties of unplugged memory blocks inside *usable\_region\_size* via DMA.

If VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE has not been negotiated, the driver SHOULD NOT read memory or query memory properties of unplugged memory blocks inside *usable\_region\_size* via the CPU.

If VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE has been negotiated, the driver MUST NOT read memory or query memory properties of unplugged memory blocks.

The driver MUST NOT request unplug of memory blocks while corresponding memory or memory properties are still in use.

The driver SHOULD initialize memory blocks after plugging them, the content is undefined.

The driver SHOULD react to resize requests from the device (*requested\_size* in the device configuration changed) by (un)plugging memory blocks.

The driver SHOULD only plug memory blocks it can actually use.

The driver MAY not reach the requested size (*requested\_size* in the device configuration), for example, because it cannot free up any plugged memory blocks to unplug them, or it would not be able to make use of unplugged memory blocks after plugging them (e.g., alignment).

### 5.15.6.2 Device Requirements: Device Operation

The device MUST provide the exact same memory properties with the exact same semantics for device memory the platform provides in the same configuration for comparable RAM.

The device MAY modify memory of unplugged memory blocks or reset memory properties of such memory blocks to platform defaults at any time.

The device MUST NOT modify memory or memory properties of plugged memory blocks.

The device MUST allow the driver to read and write memory and to query and modify memory attributes of plugged memory blocks.

If VIRTIO\_MEM\_F\_UNPLUGGED\_INACCESSIBLE has not been negotiated, the device MUST allow the driver to read memory and to query memory properties of unplugged memory blocks inside *usable\_region* - size via the CPU. <sup>16</sup>

The device MAY change the state of memory blocks during system resets.

The device SHOULD unplug all memory blocks during system resets.

### 5.15.6.3 PLUG request

Request to plug consecutive memory blocks that are currently unplugged.

The request-specific data in a PLUG request has the format:

```
struct virtio_mem_req_plug {  
    le64 addr;  
    le16 nb_blocks;  
    le16 padding[3];  
}
```

*addr* is the guest physical address of the first memory block. *nb\_blocks* is the number of consecutive memory blocks

Responses don't have request-specific data defined.

#### 5.15.6.3.1 Driver Requirements: PLUG request

The driver MUST ignore anything except the response type in a response.

#### 5.15.6.3.2 Device Requirements: PLUG request

The device MUST ignore anything except the request type and the request-specific data in a request.

The device MUST ignore the *padding* in the request-specific data in a request.

The device MUST reject requests with VIRTIO\_MEM\_RESP\_ERROR if *addr* is not aligned to the *block\_size* in the device configuration, if *nb\_blocks* is not greater than 0, or if any memory block outside of the usable device-managed memory region is covered by the request.

The device MUST reject requests with VIRTIO\_MEM\_RESP\_ERROR if any memory block covered by the request is already plugged.

The device MAY reject requests with VIRTIO\_MEM\_RESP\_BUSY if the request can currently not be processed.

The device MUST acknowledge requests with VIRTIO\_MEM\_RESP\_ACK in case all memory blocks were successfully plugged. The device MUST reflect the change in the device configuration *plugged\_size*.

<sup>16</sup>To allow for simplified dumping of memory. The CPU is expected to copy such memory to another location before starting DMA.

#### 5.15.6.4 UNPLUG request

Request to unplug consecutive memory blocks that are currently plugged.

The request-specific data in an UNPLUG request has the format:

```
struct virtio_mem_req_unplug {
    le64 addr;
    le16 nb_blocks;
    le16 padding[3];
}
```

*addr* is the guest physical address of the first memory block. *nb\_blocks* is the number of consecutive memory blocks

Responses don't have request-specific data defined.

##### 5.15.6.4.1 Driver Requirements: UNPLUG request

The driver MUST ignore anything except the response type in a response.

##### 5.15.6.4.2 Device Requirements: UNPLUG request

The device MUST ignore anything except the request type and the request-specific data in a request.

The device MUST ignore the *padding* in the request-specific data in a request.

The device MUST reject requests with VIRTIO\_MEM\_RESP\_ERROR if *addr* is not aligned to the *block\_size* in the device configuration, if *nb\_blocks* is not greater than 0, or if any memory block outside of the usable device-managed memory region is covered by the request.

The device MUST reject requests with VIRTIO\_MEM\_RESP\_ERROR if any memory block covered by the request is already unplugged.

The device MAY reject requests with VIRTIO\_MEM\_RESP\_BUSY if the request can currently not be processed.

The device MUST acknowledge requests with VIRTIO\_MEM\_RESP\_ACK in case all memory blocks were successfully unplugged. The device MUST reflect the change in the device configuration *plugged\_size*.

#### 5.15.6.5 UNPLUG ALL request

Request to unplug all memory blocks the device has currently plugged. If successful, the *plugged\_size* in the device configuration will be 0 and *usable\_region\_size* might have changed.

Requests don't have request-specific data defined, only the request type is relevant. Responses don't have request-specific data defined, only the response type is relevant.

##### 5.15.6.5.1 Driver Requirements: UNPLUG request

The driver MUST ignore any data in a response except the response type.

##### 5.15.6.5.2 Device Requirements: UNPLUG request

The device MUST ignore any data in a request except the request type.

The device MUST ignore the *padding* in the request-specific data in a request.

The device MAY reject requests with VIRTIO\_MEM\_RESP\_BUSY if the request can currently not be processed.

The device MUST acknowledge requests with VIRTIO\_MEM\_RESP\_ACK in case all memory blocks were successfully unplugged.

The device MUST set *plugged\_size* to 0 in case the request is acknowledged with VIRTIO\_MEM\_RESP\_ACK.

The device MAY modify *usable\_region\_size* before responding with VIRTIO\_MEM\_RESP\_ACK.

#### 5.15.6.6 STATE request

Request the state (plugged, unplugged, mixed) of consecutive memory blocks.

The request-specific data in a STATE request has the format:

```
struct virtio_mem_req_state {
    le64 addr;
    le16 nb_blocks;
    le16 padding[3];
};
```

*addr* is the guest physical address of the first memory block. *nb\_blocks* is the number of consecutive memory blocks.

The request-specific data in a STATE response has the format:

```
struct virtio_mem_resp_state {
    le16 type;
};
```

Whereby *type* defines one of three different state types:

```
#define VIRTIO_MEM_STATE_PLUGGED      0
#define VIRTIO_MEM_STATE_UNPLUGGED    1
#define VIRTIO_MEM_STATE_MIXED        2
```

##### 5.15.6.6.1 Driver Requirements: STATE request

The driver MUST ignore anything except the response type and the request-specific data in a response.

The driver MUST ignore the request-specific data in a response in case the response type is not VIRTIO\_MEM\_RESP\_ACK.

##### 5.15.6.6.2 Device Requirements: STATE request

The device MUST ignore anything except the request type and the request-specific data in a request.

The device MUST ignore the *padding* in the request-specific data in a request.

The device MUST reject requests with VIRTIO\_MEM\_RESP\_ERROR if *addr* is not aligned to the *block\_size* in the device configuration, if *nb\_blocks* is not greater than 0, or if any memory block outside of the usable device-managed memory region is covered by the request.

The device MUST acknowledge requests with VIRTIO\_MEM\_RESP\_ACK, supplying the state of the memory blocks.

The device MUST set the state type in the response to VIRTIO\_MEM\_STATE\_PLUGGED if all requested memory blocks are plugged. The device MUST set the state type in the response to VIRTIO\_MEM\_STATE\_UNPLUGGED if all requested memory blocks are unplugged. Otherwise, the device MUST set state type in the response to VIRTIO\_MEM\_STATE\_MIXED.

## 5.16 I2C Adapter Device

virtio-i2c is a virtual I2C adapter device. It provides a way to flexibly organize and use the host I2C controlled devices from the guest. By attaching the host ACPI I2C controlled nodes to the virtual I2C adapter device, the guest can communicate with them without changing or adding extra drivers for these controlled I2C devices.

### 5.16.1 Device ID

34

### 5.16.2 Virtqueues

0 requestq

### 5.16.3 Feature bits

**VIRTIO\_I2C\_F\_ZERO\_LENGTH\_REQUEST (0)** The device supports zero-length I2C request and *VIRTIO\_I2C\_FLAGS\_M\_RD* flag. It is mandatory to implement this feature.

**Note:** The *VIRTIO\_I2C\_FLAGS\_M\_RD* flag was not present in the initial implementation of the specification and the direction of the transfer (read or write) was inferred from the permissions (read-only or write-only) of the buffer itself. There is no need of having backwards compatibility for the older specification and so the *VIRTIO\_I2C\_FLAGS\_FAIL\_NEXT* feature is made mandatory. The driver should abort negotiation with the device, if the device doesn't offer this feature.

### 5.16.4 Device configuration layout

None currently defined.

### 5.16.5 Device Initialization

1. The virtqueue is initialized.

### 5.16.6 Device Operation

#### 5.16.6.1 Device Operation: Request Queue

The driver queues requests to the virtqueue, and they are used by the device. The request is the representation of segments of an I2C transaction. Each request is of the form:

```
struct virtio_i2c_out_hdr {
    le16 addr;
    le16 padding;
    le32 flags;
};
```

```
struct virtio_i2c_in_hdr {
    u8 status;
};
```

```
struct virtio_i2c_req {
    struct virtio_i2c_out_hdr out_hdr;
    u8 buf[];
    struct virtio_i2c_in_hdr in_hdr;
};
```

The *addr* of the request is the address of the I2C controlled device. For 7-bit address mode (A0 ... A6) and 10-bit address mode (A0 ... A9), the format of *addr* is defined as follows:

Bits	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
7-bit address	0	0	0	0	0	0	0	0	A6	A5	A4	A3	A2	A1	A0	0
10-bit address	A7	A6	A5	A4	A3	A2	A1	A0	1	1	1	1	0	A9	A8	0

The *padding* is used to pad to full dword.

The *flags* of the request is defined as follows:

**VIRTIO\_I2C\_FLAGS\_FAIL\_NEXT(0)** is used to group the requests. For a group requests, a driver clears this bit on the final request and sets it on the other requests. If this bit is set and a device fails to process the current request, it needs to fail the next request instead of attempting to execute it.

**VIRTIO\_I2C\_FLAGS\_M\_RD(1)** is used to mark the request as READ or WRITE. If **VIRTIO\_I2C\_FLAGS\_M\_RD** bit is set in the *flags*, then the request is called a read request. If **VIRTIO\_I2C\_FLAGS\_M\_RD** bit is unset in the *flags*, then the request is called a write request.

Other bits of *flags* are currently reserved as zero for future feature extensibility.

The *buf* is optional and will not be present for a zero-length request, like the SMBus "Quick" command. The *buf* contains one segment of an I2C transaction being read from or written to the device, based on the value of the **VIRTIO\_I2C\_FLAGS\_M\_RD** bit in the *flags* field.

The final *status* byte of the request is written by the device: either **VIRTIO\_I2C\_MSG\_OK** for success or **VIRTIO\_I2C\_MSG\_ERR** for error.

```
#define VIRTIO_I2C_MSG_OK      0
#define VIRTIO_I2C_MSG_ERR    1
```

The virtio I2C protocol supports write-read requests, i.e. an I2C write segment followed by a read segment (usually, the write segment provides the number of an I2C controlled device register to be read), by grouping a list of requests together using the **VIRTIO\_I2C\_FLAGS\_FAIL\_NEXT** flag.

#### 5.16.6.2 Device Operation: Operation Status

*addr*, *flags*, and "length of *buf*" are determined by the driver, while *status* is determined by the processing of the device. A driver, for a write request, puts the data to be written to the device into the *buf*, while a device, for a read request, puts the data read from device into the *buf* according to the request from the driver.

A driver may send one request or multiple requests to the device at a time. The requests in the virtqueue are both queued and processed in order.

If a driver sends multiple requests at a time and a device fails to process some of them, then a device needs to set the *status* of the first failed request to be **VIRTIO\_I2C\_MSG\_ERR**. For the remaining requests in the same group with the first failed one, a driver needs to treat them as **VIRTIO\_I2C\_MSG\_ERR**, no matter what *status* of them, a device needs to fail them instead of attempting to execute them according to the **VIRTIO\_I2C\_FLAGS\_FAIL\_NEXT** bit.

#### 5.16.6.3 Driver Requirements: Device Operation

A driver **MUST** accept the **VIRTIO\_I2C\_F\_ZERO\_LENGTH\_REQUEST** feature and **MUST** abort negotiation with the device, if the device doesn't offer this feature.

A driver **MUST** set *addr* and *flags* before sending the request.

A driver **MUST** set the reserved bits of *flags* to be zero.

A driver **MUST NOT** send the *buf*, for a zero-length request.

A driver **MUST NOT** use *buf*, for a read request, if the final *status* returned from the device is **VIRTIO\_I2C\_MSG\_ERR**.

A driver **MUST** set the **VIRTIO\_I2C\_FLAGS\_M\_RD** flag for a read operation, where the buffer is write-only for the device.

A driver **MUST NOT** set the **VIRTIO\_I2C\_FLAGS\_M\_RD** flag for a write operation, where the buffer is read-only for the device.

A driver **MUST** queue the requests in order if multiple requests are going to be sent at a time.

If multiple requests are sent at a time and some of them returned from the device have the *status* being **VIRTIO\_I2C\_MSG\_ERR**, a driver **MUST** treat the first failed request and the remaining requests in the same group with the first failed one as **VIRTIO\_I2C\_MSG\_ERR**.



#### 5.16.6.4 Device Requirements: Device Operation

A device **MUST** offer the `VIRTIO_I2C_F_ZERO_LENGTH_REQUEST` feature and **MUST** reject any driver that doesn't negotiate this feature.

A device **SHOULD** keep consistent behaviors with the hardware as described in [I2C](#).

A device **MUST NOT** change the value of `addr`, and reserved bits of `flags`.

A device **MUST** not change the value of the `buf` for a write request.

A device **MUST** place one I2C segment of the "length of `buf`", for the read request, into the `buf` according the driver's request.

A device **MUST** guarantee the requests in the virtqueue being processed in order if multiple requests are received at a time.

If multiple requests are received at a time and processing of some of the requests fails, a device **MUST** set the `status` of the first failed request to be `VIRTIO_I2C_MSG_ERR` and **MAY** set the `status` of the remaining requests in the same group with the first failed one to be `VIRTIO_I2C_MSG_ERR`.

## 5.17 SCMI Device

An SCMI device implements the Arm System Control and Management Interface (SCMI). SCMI can be used for sensors, power state management, clock management and performance management among other things.

This section relies on definitions from the [SCMI specification](#).

Virtio SCMI device and driver are mapped to SCMI platform and agent respectively. The device is visible to a particular SCMI agent. The device allows a guest to communicate as an SCMI agent using one or more SCMI protocols. The default SCMI protocols are defined in the [SCMI specification](#). Virtio provides a transport medium for exchanging SCMI messages between the SCMI agent and platform. The virtio SCMI transport allows the queueing of multiple messages and responses.

SCMI FastChannels are not supported.

### 5.17.1 Device ID

32

### 5.17.2 Virtqueues

**0** cmdq

**1** eventq

The cmdq is used by the driver to send commands to the device. The device replies with responses (not delayed responses) over the cmdq.

The eventq is used by the device to send notifications and delayed responses. The eventq only exists if `VIRTIO_SCMI_F_P2A_CHANNELS` was negotiated.

### 5.17.3 Feature bits

**VIRTIO\_SCMI\_F\_P2A\_CHANNELS (0)** Device implements some SCMI notifications, or delayed responses.

**VIRTIO\_SCMI\_F\_SHARED\_MEMORY (1)** Device implements any SCMI statistics shared memory region.

`VIRTIO_SCMI_F_P2A_CHANNELS` is used to determine the existence of the eventq. The eventq is required for SCMI notifications and delayed responses.

`VIRTIO_SCMI_F_SHARED_MEMORY` is used to determine whether the device provides any SCMI statistics shared memory region. SCMI statistics shared memory regions are defined by some SCMI protocols.

The SCMI protocols provide the `PROTOCOL_MESSAGE_ATTRIBUTES` commands to inquire about the particular SCMI notifications and delayed responses implemented by the device. The SCMI protocols provide additional commands to detect other features implemented by the device.

#### 5.17.3.1 Device Requirements: Feature bits

The device **MUST** offer `VIRTIO_SCMI_F_P2A_CHANNELS` if the device can implement at least one SCMI notification, or delayed response.

The device **MUST** offer `VIRTIO_SCMI_F_SHARED_MEMORY` if the device can implement at least one SCMI statistics shared memory region.

#### 5.17.4 Device configuration layout

There is no configuration data for the device.

#### 5.17.5 Device Initialization

The [general requirements on device initialization](#) apply.

#### 5.17.6 Device Operation

The SCMI transport used for the device puts each SCMI message into a dedicated virtio buffer. The driver uses the `cmdq` for transmitting SCMI commands and receiving the corresponding SCMI responses. The device uses the `eventq` for transmitting SCMI notifications and delayed responses. Each message includes an SCMI protocol header and payload, as defined by the [SCMI specification](#).

##### 5.17.6.1 cmdq Operation

Each buffer in the `cmdq` holds a single SCMI command once the buffer has been made available. When the buffer has been marked as used, it contains the SCMI response. An arbitrary number of such SCMI messages can be in transit at the same time. Conceptually, each SCMI message in the `cmdq` uses its own SCMI A2P (agent to platform) channel.

The SCMI response is in the same virtio buffer as the corresponding SCMI command. The response contains the return values which SCMI specifies for each command, whether synchronous or asynchronous. Delayed responses are distinct SCMI messages transmitted over the `eventq`.

Buffers in the `cmdq` contain both the request and the response. A request has the following layout:

```
struct virtio_scmi_request {
    le32 hdr;
    u8 params[<actual parameters size>];
};
```

The `virtio_scmi_request` fields are interpreted as follows:

**hdr** (device-readable) contains the SCMI message header

**params** (device-readable) comprises the SCMI message parameters

A `cmdq` response has the following layout:

```
struct virtio_scmi_response {
    le32 hdr;
    u8 ret_values[<actual return values size>];
};
```

The `virtio_scmi_response` fields are interpreted as follows:

**hdr** (device-writable) contains the SCMI message header

**ret\_values** (device-writable) comprises the SCMI message return values

If VIRTIO\_SCMI\_F\_P2A\_CHANNELS was not negotiated, the device responds to SCMI commands as if no SCMI notifications or delayed responses were implemented.

#### 5.17.6.1.1 Device Requirements: cmdq Operation

The device MAY process available commands out of order and in parallel.

The device MUST process all available commands eventually, even in the case of bursts of multiple command messages.

If the *status* field in the *virtio\_scmi\_response ret\_values* has a value other than SUCCESS, the device MUST set the size of *ret\_values* to the size of the *status* field.

If the driver requests an SCMI notification or a delayed response and there are currently NOT enough available buffers in the eventq, the device SHOULD still return SCMI status code SUCCESS.

If VIRTIO\_SCMI\_F\_P2A\_CHANNELS was not negotiated, the device MUST deny any request for an SCMI notification or a delayed response by returning SCMI status code NOT\_SUPPORTED.

If VIRTIO\_SCMI\_F\_P2A\_CHANNELS was not negotiated, the device MUST NOT indicate in the PROTOCOL\_MESSAGE\_ATTRIBUTES return values that any SCMI notification, or delayed response, is implemented.

#### 5.17.6.1.2 Driver Requirements: cmdq Operation

Before sending a command, the driver MUST wait for responses to all commands whose completion the driver considers prerequisites to executing the command.

With every command message, the driver MUST provide enough device-writable memory to enable the device to return corresponding return values.

If VIRTIO\_SCMI\_F\_P2A\_CHANNELS was not negotiated, the driver MUST NOT request any SCMI notification, nor any delayed response.

#### 5.17.6.2 Setting Up eventq Buffers

The driver has to populate the eventq before the device can use it.

##### 5.17.6.2.1 Driver Requirements: Setting Up eventq Buffers

If VIRTIO\_SCMI\_F\_P2A\_CHANNELS was negotiated, the driver SHOULD populate the eventq with buffers.

The driver MUST NOT put device-readable descriptors into the eventq.

The driver MUST NOT put into the eventq any buffer which is smaller than the largest SCMI P2A (platform to agent) message which the driver will request.

#### 5.17.6.3 eventq Operation

Each buffer in the eventq holds (once the buffer is marked as used) either a single SCMI notification, or a single SCMI delayed response. An arbitrary number of such SCMI messages can be in transit at the same time. Conceptually, each SCMI message transmitted over the eventq uses its own SCMI P2A (platform to agent) channel. Buffers in the eventq have the following layout:

```
struct virtio_scmi_event_msg {
    /* start of device-writable data */
    le32 hdr;
    u8 payload[<actual payload size>];
};
```

**hdr** (device-writable) contains the SCMI message header

**payload** (device-writable) comprises the SCMI message payload

#### 5.17.6.3.1 Device Requirements: eventq Operation

If the device intends to send a notification and there are no available buffers in the eventq, the device MAY drop the notification, or send a corresponding notification later, once enough buffers become available.

The device MAY send the notification later if the events which cause the notification take place in quick succession.

If the device sends the notification later, the device MAY send the notification with updated data, unless the specific SCMI protocol disallows this.

If the device intends to send a notification and there are available buffers, but one of the buffers is too small to fit the notification, the device MAY omit the notification.

If the device intends to send a delayed response and there are no available buffers in the eventq, the device MUST send the corresponding delayed response once enough buffers become available.

If the *status* field in a delayed response *payload* has a value other than SUCCESS, the device MUST set the size of *payload* to the size of the *status* field.

#### 5.17.6.4 Shared Memory Operation

Various SCMI protocols define statistics shared memory regions (for statistics and sensor values).

##### 5.17.6.4.1 Device Requirements: Shared Memory Operation

If VIRTIO\_SCMI\_F\_SHARED\_MEMORY was negotiated, the device MAY implement an SCMI statistics shared memory region using a virtio shared memory region.

If the device implements a shared memory region, the device MUST assign the corresponding shmid as per the following table:

SCMI statistics shared memory region	Virtio shmid
Reserved (invalid)	0
Power state statistics shared memory region	1
Performance domain statistics shared memory region	2
Sensor Values Shared Memory	3
Reserved for future use	4 to 0x7F
Vendor-specific statistics shared memory regions	0x80 to 0xFF
Reserved for future use	0x100 and greater

## 5.18 GPIO Device

The Virtio GPIO device is a virtual General Purpose Input/Output device that supports a variable number of named I/O lines, which can be configured in input mode or in output mode with logical level low (0) or high (1).

### 5.18.1 Device ID

41

### 5.18.2 Virtqueues

0 requestq

1 eventq

The *eventq* virtqueue is available only if the *VIRTIO\_GPIO\_F\_IRQ* feature is offered by the device.

### 5.18.3 Feature bits

**VIRTIO\_GPIO\_F\_IRQ (0)** The device supports interrupts on GPIO lines.

### 5.18.4 Device configuration layout

GPIO device uses the following configuration structure layout:

```
struct virtio_gpio_config {
    le16 ngpio;
    u8 padding[2];
    le32 gpio_names_size;
};
```

**ngpio** is the total number of GPIO lines supported by the device.

**padding** has no meaning and is reserved for future use. This is set to zero by the device.

**gpio\_names\_size** is the size of the gpio-names memory block in bytes, which can be fetched by the driver using the **VIRTIO\_GPIO\_MSG\_GET\_LINE\_NAMES** message. The device sets this to 0 if it doesn't support names for the GPIO lines.

### 5.18.5 Device Initialization

- The driver configures and initializes the *requestq* virtqueue.
- The driver configures and initializes the *eventq* virtqueue, if the **VIRTIO\_GPIO\_F\_IRQ** feature has been negotiated.

### 5.18.6 Device Operation: requestq

The driver uses the *requestq* virtqueue to send messages to the device. The driver sends a pair of buffers, request (filled by driver) and response (to be filled by device later), to the device. The device in turn fills the response buffer and sends it back to the driver.

```
struct virtio_gpio_request {
    le16 type;
    le16 gpio;
    le32 value;
};
```

All the fields of this structure are set by the driver and read by the device.

**type** is the GPIO message type, i.e. one of **VIRTIO\_GPIO\_MSG\_\*** values.

**gpio** is the GPIO line number, i.e.  $0 \leq gpio < ngpio$ .

**value** is a message specific value.

```
struct virtio_gpio_response {
    u8 status;
    u8 value;
};

/* Possible values of the status field */
#define VIRTIO_GPIO_STATUS_OK          0x0
#define VIRTIO_GPIO_STATUS_ERR        0x1
```

All the fields of this structure are set by the device and read by the driver.

**status** of the GPIO message, **VIRTIO\_GPIO\_STATUS\_OK** on success and **VIRTIO\_GPIO\_STATUS\_ERR** on failure.

**value** is a message specific value.

Following is the list of messages supported by the virtio gpio specification.

```

/* GPIO message types */
#define VIRTIO_GPIO_MSG_GET_LINE_NAMES      0x0001
#define VIRTIO_GPIO_MSG_GET_DIRECTION      0x0002
#define VIRTIO_GPIO_MSG_SET_DIRECTION      0x0003
#define VIRTIO_GPIO_MSG_GET_VALUE          0x0004
#define VIRTIO_GPIO_MSG_SET_VALUE          0x0005
#define VIRTIO_GPIO_MSG_SET_IRQ_TYPE       0x0006

/* GPIO Direction types */
#define VIRTIO_GPIO_DIRECTION_NONE         0x00
#define VIRTIO_GPIO_DIRECTION_OUT          0x01
#define VIRTIO_GPIO_DIRECTION_IN           0x02

/* GPIO interrupt types */
#define VIRTIO_GPIO_IRQ_TYPE_NONE          0x00
#define VIRTIO_GPIO_IRQ_TYPE_EDGE_RISING  0x01
#define VIRTIO_GPIO_IRQ_TYPE_EDGE_FALLING 0x02
#define VIRTIO_GPIO_IRQ_TYPE_EDGE_BOTH    0x03
#define VIRTIO_GPIO_IRQ_TYPE_LEVEL_HIGH   0x04
#define VIRTIO_GPIO_IRQ_TYPE_LEVEL_LOW    0x08

```

### 5.18.6.1 requestq Operation: Get Line Names

The driver sends this message to receive a stream of zero-terminated strings, where each string represents the name of a GPIO line, present in increasing order of the GPIO line numbers. The names of the GPIO lines are optional and may be present only for a subset of GPIO lines. If missing, then a zero-byte must be present for the GPIO line. If present, the name string must be zero-terminated and the name must be unique within a GPIO Device. The names of the GPIO lines are encoded in 7-bit ASCII.

These names of the GPIO lines should be most meaningful producer names for the system, such as name indicating the usage. For example "MMC-CD", "Red LED Vdd" and "ethernet reset" are reasonable line names as they describe what the line is used for, while "GPIO0" is not a good name to give to a GPIO line.

Here is an example of how the gpio names memory block may look like for a GPIO device with 10 GPIO lines, where line names are provided only for lines 0 ("MMC-CD"), 5 ("Red LED Vdd") and 7 ("ethernet reset").

```

u8 gpio_names[] = {
    'M', 'M', 'C', '-', 'C', 'D', 0,
    0,
    0,
    0,
    0,
    'R', 'e', 'd', ' ', 'L', 'E', 'D', ' ', 'V', 'd', 'd', 0,
    0,
    'E', 't', 'h', 'e', 'r', 'n', 'e', 't', ' ', 'r', 'e', 's', 'e', 't', 0,
    0,
    0
};

```

The device sets the *gpio\_names\_size* to a non-zero value if this message is supported by the device, else it must be set to zero.

This message type uses different layout for the response structure, as the device needs to return the *gpio\_names* array.

```

struct virtio_gpio_response_N {
    u8 status;
    u8 value[N];
};

```

The driver must allocate the *value[N]* buffer in the *struct virtio\_gpio\_response\_N* for N bytes, where N = *gpio\_names\_size*.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_-GET_LINE_NAMES</i>	0	0

  

Response	status	value[N]	Where N is
	<i>VIRTIO_GPIO_STATUS_*</i>	gpio-names	<i>gpio_names_size</i>

#### 5.18.6.2 requestq Operation: Get Direction

The driver sends this message to request the device to return a line's configured direction.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_-GET_DIRECTION</i>	line number	0

  

Response	status	value
	<i>VIRTIO_GPIO_STATUS_*</i>	0 = none, 1 = output, 2 = input

#### 5.18.6.3 requestq Operation: Set Direction

The driver sends this message to request the device to configure a line's direction. The driver can either set the direction to *VIRTIO\_GPIO\_DIRECTION\_IN* or *VIRTIO\_GPIO\_DIRECTION\_OUT*, which also activates the line, or to *VIRTIO\_GPIO\_DIRECTION\_NONE*, which deactivates the line.

The driver should set the value of the GPIO line, using the *VIRTIO\_GPIO\_MSG\_SET\_VALUE* message, before setting the direction of the line to output to avoid any undesired behavior.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_SET_DIRECTION</i>	line number	0 = none, 1 = output, 2 = input

  

Response	status	value
	<i>VIRTIO_GPIO_STATUS_*</i>	0

#### 5.18.6.4 requestq Operation: Get Value

The driver sends this message to request the device to return current value sensed on a line.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_-GET_VALUE</i>	line number	0

  

Response	status	value
	<i>VIRTIO_GPIO_STATUS_*</i>	0 = low, 1 = high

#### 5.18.6.5 requestq Operation: Set Value

The driver sends this message to request the device to set the value of a line. The line may already be configured for output or may get configured to output later, at which point this output value must be used by the device for the line.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_SET_VALUE</i>	line number	0 = low, 1 = high



Response	status	value
	<i>VIRTIO_GPIO_STATUS_*</i>	0

#### 5.18.6.6 requestq Operation: Set IRQ Type

This request is allowed only if the *VIRTIO\_GPIO\_F\_IRQ* feature has been negotiated.

The interrupt configuration is divided into two steps, enabling or disabling of the interrupt at the device and masking or unmasking of the interrupt for delivery at the driver. This request only pertains to enabling or disabling of the interrupt at the device, the masking and unmasking of the interrupt is handled by a separate request that takes place over the *eventq* virtqueue.

The driver sends the *VIRTIO\_GPIO\_MSG\_SET\_IRQ\_TYPE* message over the *requestq* virtqueue to enable or disable interrupt for a GPIO line at the device.

The driver sends this message with trigger type set to any valid value other than *VIRTIO\_GPIO\_IRQ\_TYPE\_NONE*, to enable the interrupt for a GPIO line, this doesn't unmask the interrupt for delivery at the driver though. For edge trigger type, the device should latch the interrupt events from this point onward and notify it to the driver once the interrupt is unmasked. For level trigger type, the device should notify the driver once the interrupt signal on a line is sensed and the interrupt is unmasked for the line.

The driver sends this message with trigger type set to *VIRTIO\_GPIO\_IRQ\_TYPE\_NONE*, to disable the interrupt for a GPIO line. The device should discard any latched interrupt event associated with it. In order to change the trigger type of an already enabled interrupt, the driver should first disable the interrupt and then re-enable it with appropriate trigger type.

The interrupts are masked at initialization and the driver unmask them by queuing a pair of buffers, of type *virtio\_gpio\_irq\_request* and *virtio\_gpio\_irq\_response*, over the *eventq* virtqueue for a GPIO line. A separate pair of buffers must be queued for each GPIO line, the driver wants to configure for interrupts. Once an already enabled interrupt is unmasked by the driver, the device can notify the driver of an active interrupt signal on the GPIO line. This is done by updating the *struct virtio\_gpio\_irq\_response* buffer's *status* with *VIRTIO\_GPIO\_IRQ\_STATUS\_VALID* and returning the updated buffers to the driver. The interrupt is masked automatically at this point until the buffers are available again at the device.

The interrupt for a GPIO line should not be unmasked before being enabled by the driver.

If the interrupt is disabled by the driver, by setting the trigger type to *VIRTIO\_GPIO\_IRQ\_TYPE\_NONE*, or the interrupt is unmasked without being enabled first, the device should return any unused pair of buffers for the GPIO line, over the *eventq* virtqueue, after setting the *status* field to *VIRTIO\_GPIO\_IRQ\_STATUS\_INVALID*. This also masks the interrupt.

Request	type	gpio	value
	<i>VIRTIO_GPIO_MSG_SET_IRQ_TYPE</i>	line number	one of <i>VIRTIO_GPIO_IRQ_TYPE_*</i>

Response	status	value
	<i>VIRTIO_GPIO_STATUS_*</i>	0

#### 5.18.6.7 requestq Operation: Message Flow

- The driver queues *struct virtio\_gpio\_request* and *virtio\_gpio\_response* buffers to the *requestq* virtqueue, after filling all fields of the *struct virtio\_gpio\_request* buffer as defined by the specific message type.
- The driver notifies the device of the presence of buffers on the *requestq* virtqueue.
- The device, after receiving the message from the driver, processes it and fills all the fields of the *struct virtio\_gpio\_response* buffer (received from the driver). The *status* must be set to *VIRTIO\_GPIO\_STATUS\_OK* on success and *VIRTIO\_GPIO\_STATUS\_ERR* on failure.
- The device puts the buffers back on the *requestq* virtqueue and notifies the driver of the same.
- The driver fetches the buffers and processes the response received in the *virtio\_gpio\_response* buffer.

- The driver can send multiple messages in parallel for same or different GPIO line.

#### 5.18.6.8 Driver Requirements: requestq Operation

- The driver MUST send messages on the *requestq* virtqueue.
- The driver MUST queue both *struct virtio\_gpio\_request* and *virtio\_gpio\_response* for every message sent to the device.
- The *struct virtio\_gpio\_request* buffer MUST be filled by the driver and MUST be read-only for the device.
- The *struct virtio\_gpio\_response* buffer MUST be filled by the device and MUST be writable by the device.
- The driver MAY send multiple messages for same or different GPIO lines in parallel.
- The driver MUST NOT send IRQ messages if the *VIRTIO\_GPIO\_F\_IRQ* feature has not been negotiated.
- The driver MUST NOT send IRQ messages for a GPIO line configured for output.
- The driver MUST set the IRQ trigger type to *VIRTIO\_GPIO\_IRQ\_TYPE\_NONE* once it is done using the GPIO line configured for interrupts.
- In order to change the trigger type of an already enabled interrupt, the driver MUST first disable the interrupt and then re-enable it with appropriate trigger type.

#### 5.18.6.9 Device Requirements: requestq Operation

- The device MUST set all the fields of the *struct virtio\_gpio\_response* before sending it back to the driver.
- The device MUST set all the fields of the *struct virtio\_gpio\_config* on receiving a configuration request from the driver.
- The device MUST set the *gpio\_names\_size* field as zero in the *struct virtio\_gpio\_config*, if it doesn't implement names for individual GPIO lines.
- The device MUST set the *gpio\_names\_size* field, in the *struct virtio\_gpio\_config*, with the size of *gpio-names* memory block in bytes, if the device implements names for individual GPIO lines. The strings MUST be zero-terminated and an unique (if available) within the GPIO device.
- The device MUST process multiple messages, for the same GPIO line, sequentially and respond to them in the order they were received on the virtqueue.
- The device MAY process messages, for different GPIO lines, out of order and in parallel, and MAY send message's response to the driver out of order.
- The device MUST discard all state information corresponding to a GPIO line, once the driver has requested to set its direction to *VIRTIO\_GPIO\_DIRECTION\_NONE*.
- The device MUST latch an edge interrupt if the interrupt is enabled but still masked.
- The device MUST NOT latch an level interrupt if the interrupt is enabled but still masked.
- The device MUST discard any latched interrupt for a GPIO line, once interrupt is disabled for the same.

#### 5.18.7 Device Operation: eventq

The *eventq* virtqueue is used by the driver to unmask the interrupts and used by the device to notify the driver of newly sensed interrupts. In order to unmask interrupt on a GPIO line, the driver queues a pair of buffers, *struct virtio\_gpio\_irq\_request* (filled by driver) and *struct virtio\_gpio\_irq\_response* (to be filled by device later), to the *eventq* virtqueue. A separate pair of buffers must be queued for each GPIO line, the driver wants to configure for interrupts. The device, on sensing an interrupt, returns the pair of buffers for the respective GPIO line, which also masks the interrupts. The driver can queue the buffers again to unmask the interrupt.

```
struct virtio_gpio_irq_request {
    le16 gpio;
};
```

This structure is filled by the driver and read by the device.

**gpio** is the GPIO line number, i.e.  $0 \leq \text{gpio} < \text{ngpio}$ .

```
struct virtio_gpio_irq_response {
    u8 status;
};

/* Possible values of the interrupt status field */
#define VIRTIO_GPIO_IRQ_STATUS_INVALID    0x0
#define VIRTIO_GPIO_IRQ_STATUS_VALID     0x1
```

This structure is filled by the device and read by the driver.

**status** of the interrupt event, *VIRTIO\_GPIO\_IRQ\_STATUS\_VALID* on interrupt and *VIRTIO\_GPIO\_IRQ\_STATUS\_INVALID* to return the buffers back to the driver after interrupt is disabled.

#### 5.18.7.1 eventq Operation: Message Flow

- The virtio-gpio driver is requested by a client driver to enable interrupt for a GPIO line and configure it to a particular trigger type.
- The driver sends the *VIRTIO\_GPIO\_MSG\_SET\_IRQ\_TYPE* message, over the *requestq* virtqueue, and the device configures the GPIO line for the requested trigger type and enables the interrupt. The interrupt is still masked for delivery though. The device shall latch the interrupt from now onward for edge trigger type.
- The driver unmask the interrupt by queuing a pair of buffers to the *eventq* virtqueue for the GPIO line. The driver can do this before enabling the interrupt as well, though the interrupt must be both unmasked and enabled to get delivered at the driver.
- The driver notifies the device of the presence of new buffers on the *eventq* virtqueue. The interrupt is fully configured at this point.
- The device, on sensing an active interrupt on the GPIO line, finds the matching buffers (based on GPIO line number) from the *eventq* virtqueue and update its *struct virtio\_gpio\_irq\_response* buffer's *status* with *VIRTIO\_GPIO\_IRQ\_STATUS\_VALID* and returns the pair of buffers to the device. This results in masking the interrupt as well.
- The device notifies the driver of the presence of returned buffers on the *eventq* virtqueue.
- If the GPIO line is configured for level interrupts, the device ignores any further interrupt signals on this GPIO line, until the interrupt is unmasked again by the driver by making the buffers available to the device. Once the interrupt is unmasked again and the interrupt on the line is still active, the device shall notify the driver again.
- If the GPIO line is configured for edge interrupts, the device latches the interrupt received for this GPIO line, until the interrupt is unmasked again by the driver by making the buffers available to the device. Once the interrupt is unmasked again and an interrupt was latched earlier, the device shall notify the driver again.
- The driver on receiving the notification from the device, processes the interrupt. The interrupt is masked at the device until the buffers are queued again by the driver.
- In a typical guest operating system kernel, the virtio-gpio driver notifies the client driver, that is associated with this GPIO line, to process the event. In the case of a level triggered interrupt, the client driver shall fully process and acknowledge the event at its source to return the line to its inactive state before the interrupt is unmasked again to avoid a spurious interrupt.
- Once the interrupt is handled, the driver may queue a pair of buffers for the GPIO line to unmask the interrupt again.

- The driver can also disable the interrupt by sending the `VIRTIO_GPIO_MSG_SET_IRQ_TYPE` message, with `VIRTIO_GPIO_IRQ_TYPE_NONE` trigger type. In that case, the device shall return the unused pair of buffers for the GPIO line after setting the `status` field with `VIRTIO_GPIO_IRQ_STATUS_INVALID`.

#### 5.18.7.2 Driver Requirements: eventq Operation

- The driver MUST both enable and unmask the interrupt in order to get notified for the same.
- The driver MUST enable the interrupt before unmasking it.
- To unmask the interrupt, the driver MUST queue a separate pair of buffers to the `eventq` virtqueue for each GPIO line.
- The driver MUST NOT add multiple pairs of buffers for the same GPIO line on the `eventq` virtqueue.

#### 5.18.7.3 Device Requirements: eventq Operation

- The device MUST NOT send an interrupt event to the driver for a GPIO line unless the interrupt has been both unmasked and enabled by the driver.
- On receiving `VIRTIO_GPIO_MSG_SET_IRQ_TYPE` message, with `VIRTIO_GPIO_IRQ_TYPE_NONE` trigger type, the device MUST return the buffers, if they were received earlier, after setting the `status` field to `VIRTIO_GPIO_IRQ_STATUS_INVALID`.

## 5.19 PMEM Device

The virtio pmem device is a persistent memory (NVDIMM) device that provides a virtio based asynchronous flush mechanism. This avoids the need for a separate page cache in the guest and keeps the page cache only in the host. Under memory pressure, the host makes use of efficient memory reclaim decisions for page cache pages of all the guests. This helps to reduce the memory footprint and fits more guests in the host system.

The virtio pmem device provides access to byte-addressable persistent memory. The persistent memory is a directly accessible range of system memory. Data written to this memory is made persistent by separately sending a flush command. Writes that have been flushed are preserved across device reset and power failure.

### 5.19.1 Device ID

27

### 5.19.2 Virtqueues

0 req\_vq

### 5.19.3 Feature bits

**VIRTIO\_PMEM\_F\_SHMEM\_REGION (0)** The guest physical address range will be indicated as a shared memory region.

### 5.19.4 Device configuration layout

```
struct virtio_pmem_config {
    le64 start;
    le64 size;
};
```

**start** contains the physical address of the first byte of the persistent memory region, if `VIRTIO_PMEM_F_SHMEM_REGION` has not been negotiated.

**size** contains the length of this address range, if VIRTIO\_PMEM\_F\_SHMEM\_REGION has not been negotiated.

### 5.19.5 Device Initialization

The device indicates the guest physical address to the driver in one of two ways:

1. As a physical address, using `virtio_pmem_config`.
2. As a shared memory region.

The driver determines the start address and size of the persistent memory region in preparation for reading or writing data.

The driver initializes `req_vq` in preparation for making flush requests.

#### 5.19.5.1 Device Requirements: Device Initialization

If VIRTIO\_PMEM\_F\_SHMEM\_REGION has been negotiated, the device **MUST** indicate the guest physical address as a shared memory region. The device **MUST** use shared memory region ID 0. The device **SHOULD** set *start* and *size* to zero.

If VIRTIO\_PMEM\_F\_SHMEM\_REGION has not been negotiated, the device **MUST** indicate the guest physical address as a physical address. The device **MUST** set *start* to the absolute address and *size* to the size of the address range, in bytes.

#### 5.19.5.2 Driver Requirements: Device Initialization

If VIRTIO\_PMEM\_F\_SHMEM\_REGION has been negotiated, the driver **MUST** query shared memory ID 0 for the physical address ranges, and **MUST NOT** use *start* or *stop*.

If VIRTIO\_PMEM\_F\_SHMEM\_REGION has not been negotiated, the driver **MUST** read the physical address ranges from *start* and *stop*.

### 5.19.6 Driver Operations

Requests have the following format:

```
struct virtio_pmem_req {
    le32 type;
};
```

*type* is the request command type.

Possible request types are:

```
#define VIRTIO_PMEM_REQ_TYPE_FLUSH    0
```

### 5.19.7 Device Operations

#### 5.19.7.1 Device Requirements: Device Operation: Virtqueue flush

The device **MUST** ensure that all writes completed before a flush request persist across device reset and power failure before completing the flush request.

#### 5.19.7.2 Device Operations

```
struct virtio_pmem_resp {
    le32 ret;
};
```

*ret* is the value which the device returns after command completion.

### 5.19.7.3 Device Requirements: Device Operation: Virtqueue return

The device MUST return "0" for success and "-1" for failure.

## 5.19.8 Possible security implications

There could be potential security implications depending on how memory mapped backing device is used. By default device emulation is done with SHARED memory mapping. There is a contract between driver and device to access shared memory region for read or write operations.

If a malicious driver or device maps the same memory region, the attacker can make use of known side channel attacks to predict the current state of data. If both attacker and victim somehow execute same shared code after a flush or evict operation, with difference in execution timing attacker could infer another device's data.

## 5.19.9 Countermeasures

### 5.19.9.1 With SHARED mapping

If a device's backing region is shared between multiple devices, this may act as a metric for side channel attacks. As a counter measure every device should have its own (not shared with another driver) SHARED backing memory.

### 5.19.9.2 With PRIVATE mapping

There maybe be chances of side channels attack with PRIVATE memory mapping similar to SHARED with read-only shared mappings. PRIVATE is not used for virtio pmem making this usecase irrelevant.

### 5.19.9.3 Workload specific mapping

When using SHARED mappings with a workload that is a single application inside the driver where the risk in sharing data is very low or nonexistent, the device sharing the same backing region with a SHARED mapping can be used as a valid configuration.

### 5.19.9.4 Prevent cache eviction

Don't allow device shared region eviction from driver filesystem trim or discard like commands with virtio pmem. This rules out any possibility of evict-reload cache side channel attacks if backing region is shared (SHARED) between mutiple devices. Though if we use per device backing file with shared mapping this countermeasure is not required.

---

## 6 Reserved Feature Bits

Currently these device-independent feature bits are defined:

**VIRTIO\_F\_INDIRECT\_DESC (28)** Negotiating this feature indicates that the driver can use descriptors with the `VIRTQ_DESC_F_INDIRECT` flag set, as described in [2.7.5.3 Indirect Descriptors](#) and [2.8.7 Indirect Flag: Scatter-Gather Support](#).

**VIRTIO\_F\_EVENT\_IDX(29)** This feature enables the *used\_event* and the *avail\_event* fields as described in [2.7.7](#), [2.7.8](#) and [2.8.10](#).

**VIRTIO\_F\_VERSION\_1(32)** This indicates compliance with this specification, giving a simple way to detect legacy devices or drivers.

**VIRTIO\_F\_ACCESS\_PLATFORM(33)** This feature indicates that the device can be used on a platform where device access to data in memory is limited and/or translated. E.g. this is the case if the device can be located behind an IOMMU that translates bus addresses from the device into physical addresses in memory, if the device can be limited to only access certain memory addresses or if special commands such as a cache flush can be needed to synchronise data in memory with the device. Whether accesses are actually limited or translated is described by platform-specific means. If this feature bit is set to 0, then the device has same access to memory addresses supplied to it as the driver has. In particular, the device will always use physical addresses matching addresses used by the driver (typically meaning physical addresses used by the CPU) and not translated further, and can access any address supplied to it by the driver. When clear, this overrides any platform-specific description of whether device access is limited or translated in any way, e.g. whether an IOMMU may be present.

**VIRTIO\_F\_RING\_PACKED(34)** This feature indicates support for the packed virtqueue layout as described in [2.8 Packed Virtqueues](#).

**VIRTIO\_F\_IN\_ORDER(35)** This feature indicates that all buffers are used by the device in the same order in which they have been made available.

**VIRTIO\_F\_ORDER\_PLATFORM(36)** This feature indicates that memory accesses by the driver and the device are ordered in a way described by the platform.

If this feature bit is negotiated, the ordering in effect for any memory accesses by the driver that need to be ordered in a specific way with respect to accesses by the device is the one suitable for devices described by the platform. This implies that the driver needs to use memory barriers suitable for devices described by the platform; e.g. for the PCI transport in the case of hardware PCI devices.

If this feature bit is not negotiated, then the device and driver are assumed to be implemented in software, that is they can be assumed to run on identical CPUs in an SMP configuration. Thus a weaker form of memory barriers is sufficient to yield better performance.

**VIRTIO\_F\_SR\_IOV(37)** This feature indicates that the device supports Single Root I/O Virtualization. Currently only PCI devices support this feature.

**VIRTIO\_F\_NOTIFICATION\_DATA(38)** This feature indicates that the driver passes extra data (besides identifying the virtqueue) in its device notifications. See [2.9 Driver Notifications](#).

**VIRTIO\_F\_NOTIF\_CONFIG\_DATA(39)** This feature indicates that the driver uses the data provided by the device as a virtqueue identifier in available buffer notifications. As mentioned in section [2.9](#), when the driver is required to send an available buffer notification to the device, it sends the virtqueue number to be notified. The method of delivering notifications is transport specific. With the PCI transport, the device can optionally provide a per-virtqueue value for the driver to use in driver notifications, instead



of the virtqueue number. Some devices may benefit from this flexibility by providing, for example, an internal virtqueue identifier, or an internal offset related to the virtqueue number.

This feature indicates the availability of such value. The definition of the data to be provided in driver notification and the delivery method is transport specific. For more details about driver notifications over PCI see [4.1.5.2](#).

**VIRTIO\_F\_RING\_RESET(40)** This feature indicates that the driver can reset a queue individually. See [2.6.1](#).

## 6.1 Driver Requirements: Reserved Feature Bits

A driver **MUST** accept `VIRTIO_F_VERSION_1` if it is offered. A driver **MAY** fail to operate further if `VIRTIO_F_VERSION_1` is not offered.

A driver **SHOULD** accept `VIRTIO_F_ACCESS_PLATFORM` if it is offered, and it **MUST** then either disable the IOMMU or configure the IOMMU to translate bus addresses passed to the device into physical addresses in memory. If `VIRTIO_F_ACCESS_PLATFORM` is not offered, then a driver **MUST** pass only physical addresses to the device.

A driver **SHOULD** accept `VIRTIO_F_RING_PACKED` if it is offered.

A driver **SHOULD** accept `VIRTIO_F_ORDER_PLATFORM` if it is offered. If `VIRTIO_F_ORDER_PLATFORM` has been negotiated, a driver **MUST** use the barriers suitable for hardware devices.

If `VIRTIO_F_SR_IOV` has been negotiated, a driver **MAY** enable virtual functions through the device's PCI SR-IOV capability structure. A driver **MUST NOT** negotiate `VIRTIO_F_SR_IOV` if the device does not have a PCI SR-IOV capability structure or is not a PCI device. A driver **MUST** negotiate `VIRTIO_F_SR_IOV` and complete the feature negotiation (including checking the `FEATURES_OK device status` bit) before enabling virtual functions through the device's PCI SR-IOV capability structure. After once successfully negotiating `VIRTIO_F_SR_IOV`, the driver **MAY** enable virtual functions through the device's PCI SR-IOV capability structure even if the device or the system has been fully or partially reset, and even without re-negotiating `VIRTIO_F_SR_IOV` after the reset.

A driver **SHOULD** accept `VIRTIO_F_NOTIF_CONFIG_DATA` if it is offered.

## 6.2 Device Requirements: Reserved Feature Bits

A device **MUST** offer `VIRTIO_F_VERSION_1`. A device **MAY** fail to operate further if `VIRTIO_F_VERSION_1` is not accepted.

A device **SHOULD** offer `VIRTIO_F_ACCESS_PLATFORM` if its access to memory is through bus addresses distinct from and translated by the platform to physical addresses used by the driver, and/or if it can only access certain memory addresses with said access specified and/or granted by the platform. A device **MAY** fail to operate further if `VIRTIO_F_ACCESS_PLATFORM` is not accepted.

If `VIRTIO_F_IN_ORDER` has been negotiated, a device **MUST** use buffers in the same order in which they have been available.

A device **MAY** fail to operate further if `VIRTIO_F_ORDER_PLATFORM` is offered but not accepted. A device **MAY** operate in a slower emulation mode if `VIRTIO_F_ORDER_PLATFORM` is offered but not accepted.

It is **RECOMMENDED** that an add-in card based PCI device offers both `VIRTIO_F_ACCESS_PLATFORM` and `VIRTIO_F_ORDER_PLATFORM` for maximum portability.

A device **SHOULD** offer `VIRTIO_F_SR_IOV` if it is a PCI device and presents a PCI SR-IOV capability structure, otherwise it **MUST NOT** offer `VIRTIO_F_SR_IOV`.

## 6.3 Legacy Interface: Reserved Feature Bits

Transitional devices **MAY** offer the following:



**VIRTIO\_F\_NOTIFY\_ON\_EMPTY (24)** If this feature has been negotiated by driver, the device **MUST** issue a used buffer notification if the device runs out of available descriptors on a virtqueue, even though notifications are suppressed using the `VIRTQ_AVAIL_F_NO_INTERRUPT` flag or the `used_event` field.

**Note:** An example of a driver using this feature is the legacy networking driver: it doesn't need to know every time a packet is transmitted, but it does need to free the transmitted packets a finite time after they are transmitted. It can avoid using a timer if the device notifies it when all the packets are transmitted.

Transitional devices **MUST** offer, and if offered by the device transitional drivers **MUST** accept the following:

**VIRTIO\_F\_ANY\_LAYOUT (27)** This feature indicates that the device accepts arbitrary descriptor layouts, as described in Section [2.7.4.3 Legacy Interface: Message Framing](#).

**UNUSED (30)** Bit 30 is used by qemu's implementation to check for experimental early versions of virtio which did not perform correct feature negotiation, and **SHOULD NOT** be negotiated.

---

## 7 Conformance

This chapter lists the conformance targets and clauses for each; this also forms a useful checklist which authors are asked to consult for their implementations!

### 7.1 Conformance Targets

Conformance targets:

**Driver** A driver **MUST** conform to four conformance clauses:

- Clause [7.2](#).
- One of clauses [7.2.1](#), [7.2.2](#) or [7.2.3](#).
- One of clauses [7.2.4](#), [7.2.5](#), [7.2.6](#), [7.2.7](#), [7.2.8](#), [7.2.9](#), [7.2.10](#), [7.2.11](#), [7.2.12](#), [7.2.13](#), [7.2.14](#), [7.2.15](#), [7.2.16](#), [7.2.17](#), [7.2.18](#), [7.2.19](#), [7.2.20](#) or [7.2.21](#).
- Clause [7.4](#).

**Device** A device **MUST** conform to four conformance clauses:

- Clause [7.3](#).
- One of clauses [7.3.1](#), [7.3.2](#) or [7.3.3](#).
- One of clauses [7.3.4](#), [7.3.5](#), [7.3.6](#), [7.3.7](#), [7.3.8](#), [7.3.9](#), [7.3.10](#), [7.3.11](#), [7.3.12](#), [7.3.13](#), [7.3.14](#), [7.3.15](#), [7.3.16](#), [7.3.17](#), [7.3.18](#), [7.3.19](#), [7.3.20](#) or [7.3.21](#).
- Clause [7.4](#).

### 7.2 Clause 1: Driver Conformance

A driver **MUST** conform to the following normative statements:

- [2.1.1](#)
- [2.2.1](#)
- [2.4.2](#)
- [2.5.1](#)
- [2.7.1](#)
- [2.7.4.2](#)
- [2.7.5.2](#)
- [2.7.5.3.1](#)
- [2.7.7.1](#)
- [2.7.6.1](#)
- [2.7.8.3](#)
- [2.7.10.1](#)
- [2.7.13.3.1](#)

- [2.7.13.4.1](#)
- [3.1.1](#)
- [3.3.1](#)
- [6.1](#)

### 7.2.1 Clause 2: PCI Driver Conformance

A PCI driver MUST conform to the following normative statements:

- [4.1.2.2](#)
- [4.1.3.1](#)
- [4.1.4.1](#)
- [4.1.4.3.2](#)
- [4.1.4.5.2](#)
- [4.1.4.9.2](#)
- [4.1.5.1.2.2](#)
- [4.1.5.4.2](#)

### 7.2.2 Clause 3: MMIO Driver Conformance

An MMIO driver MUST conform to the following normative statements:

- [4.2.2.2](#)
- [4.2.3.1.1](#)
- [4.2.3.4.1](#)

### 7.2.3 Clause 4: Channel I/O Driver Conformance

A Channel I/O driver MUST conform to the following normative statements:

- [4.3.1.4](#)
- [4.3.2.1.2](#)
- [4.3.2.3.1](#)
- [4.3.3.1.2.2](#)
- [4.3.3.2.2](#)
- [4.3.3.3.2](#)

### 7.2.4 Clause 5: Network Driver Conformance

A network driver MUST conform to the following normative statements:

- [5.1.4.2](#)
- [5.1.6.2.1](#)
- [5.1.6.3.1](#)
- [5.1.6.4.2](#)
- [5.1.6.5.1.2](#)
- [5.1.6.5.2.2](#)
- [5.1.6.5.4.1](#)

- [5.1.6.5.6.1](#)
- [5.1.6.5.8.2](#)
- [5.1.6.5.7.2](#)

### 7.2.5 Clause 6: Block Driver Conformance

A block driver MUST conform to the following normative statements:

- [5.2.5.1](#)
- [5.2.6.1](#)

### 7.2.6 Clause 7: Console Driver Conformance

A console driver MUST conform to the following normative statements:

- [5.3.6.1](#)
- [5.3.6.2.2](#)

### 7.2.7 Clause 8: Entropy Driver Conformance

An entropy driver MUST conform to the following normative statements:

- [5.4.6.1](#)

### 7.2.8 Clause 9: Traditional Memory Balloon Driver Conformance

A traditional memory balloon driver MUST conform to the following normative statements:

- [5.5.3.1](#)
- [5.5.6.1](#)
- [5.5.6.3.1](#)
- [5.5.6.5.1](#)
- [5.5.6.6.1](#)
- [5.5.6.7.1](#)

### 7.2.9 Clause 10: SCSI Host Driver Conformance

An SCSI host driver MUST conform to the following normative statements:

- [5.6.4.1](#)
- [5.6.6.1.2](#)
- [5.6.6.3.1](#)

### 7.2.10 Clause 11: Input Driver Conformance

An input driver MUST conform to the following normative statements:

- [5.8.5.1](#)
- [5.8.6.1](#)

### 7.2.11 Clause 12: Crypto Driver Conformance

A Crypto driver MUST conform to the following normative statements:

- [5.9.5.2](#)
- [5.9.6.1](#)
- [5.9.7.2.1.5](#)
- [5.9.7.2.1.7](#)
- [5.9.7.4.1](#)
- [5.9.7.5.1](#)
- [5.9.7.6.1](#)
- [5.9.7.7.1](#)

### 7.2.12 Clause 13: Socket Driver Conformance

A socket driver MUST conform to the following normative statements:

- [5.10.6.3.1](#)
- [5.10.6.4.1](#)
- [5.10.6.7.1](#)

### 7.2.13 Clause 14: File System Driver Conformance

A file system driver MUST conform to the following normative statements:

- [5.11.4.1](#)
- [5.11.6.2.2](#)
- [5.11.6.3.1](#)
- [5.11.6.4.2](#)

### 7.2.14 Clause 15: RPMB Driver Conformance

A RPMB driver MUST conform to the following normative statements:

- [5.12.6.2](#)

### 7.2.15 Clause 16: IOMMU Driver Conformance

An IOMMU driver MUST conform to the following normative statements:

- [5.13.3.1](#)
- [5.13.4.1](#)
- [5.13.5.1](#)
- [5.13.6.1](#)
- [5.13.6.3.1](#)
- [5.13.6.4.1](#)
- [5.13.6.5.1](#)
- [5.13.6.6.1](#)
- [5.13.6.7.1](#)
- [5.13.6.8.1.1](#)

- [5.13.6.9.1](#)

### **7.2.16 Clause 17: Sound Driver Conformance**

A sound driver MUST conform to the following normative statements:

- [5.14.5.1](#)
- [5.14.6.2](#)
- [5.14.6.6.3.2](#)
- [5.14.6.8.1.2](#)
- [5.14.6.8.2.2](#)

### **7.2.17 Clause 18: Memory Driver Conformance**

A memory driver MUST conform to the following normative statements:

- [5.15.4.1](#)
- [5.15.5.1](#)
- [5.15.6.1](#)
- [5.15.6.3.1](#)
- [5.15.6.4.1](#)
- [5.15.6.5.1](#)
- [5.15.6.6.1](#)

### **7.2.18 Clause 19: I2C Adapter Driver Conformance**

An I2C Adapter driver MUST conform to the following normative statements:

- [5.16.6.3](#)

### **7.2.19 Clause 20: SCMI Driver Conformance**

An SCMI driver MUST conform to the following normative statements:

- [5.17.6.1.2](#)
- [5.17.6.2.1](#)

### **7.2.20 Clause 21: GPIO Driver Conformance**

A General Purpose Input/Output (GPIO) driver MUST conform to the following normative statements:

- [5.18.6.8](#)
- [5.18.7.2](#)

### **7.2.21 Clause 22: PMEM Driver Conformance**

A PMEM driver MUST conform to the following normative statements:

- [5.19.5.2](#)

## 7.3 Clause 23: Device Conformance

A device MUST conform to the following normative statements:

- [2.1.2](#)
- [2.2.2](#)
- [2.4.1](#)
- [2.5.2](#)
- [2.7.4.1](#)
- [2.7.5.1](#)
- [2.7.5.3.2](#)
- [2.7.7.2](#)
- [2.7.8.2](#)
- [2.7.10.2](#)
- [2.10.2](#)
- [6.2](#)

### 7.3.1 Clause 24: PCI Device Conformance

A PCI device MUST conform to the following normative statements:

- [4.1.1](#)
- [4.1.2.1](#)
- [4.1.3.2](#)
- [4.1.4.2](#)
- [4.1.4.3.1](#)
- [4.1.4.4.1](#)
- [4.1.4.5.1](#)
- [4.1.4.6.1](#)
- [4.1.4.7.1](#)
- [4.1.4.9.1](#)
- [4.1.4.11.0.1](#)
- [4.1.5.1.2.1](#)
- [4.1.5.3.1](#)
- [4.1.5.4.1](#)

### 7.3.2 Clause 25: MMIO Device Conformance

An MMIO device MUST conform to the following normative statements:

- [4.2.2.1](#)

### 7.3.3 Clause 26: Channel I/O Device Conformance

A Channel I/O device MUST conform to the following normative statements:

- [4.3.1.3](#)
- [4.3.2.1.1](#)
- [4.3.2.2.1](#)
- [4.3.2.3.2](#)
- [4.3.2.6.3.1](#)
- [4.3.3.1.2.1](#)
- [4.3.3.2.1](#)
- [4.3.3.3.1](#)

### 7.3.4 Clause 27: Network Device Conformance

A network device MUST conform to the following normative statements:

- [5.1.4.1](#)
- [5.1.6.2.2](#)
- [5.1.6.3.2](#)
- [5.1.6.4.1](#)
- [5.1.6.5.1.1](#)
- [5.1.6.5.2.1](#)
- [5.1.6.5.4.2](#)
- [5.1.6.5.6.2](#)
- [5.1.6.5.7.3](#)

### 7.3.5 Clause 28: Block Device Conformance

A block device MUST conform to the following normative statements:

- [5.2.5.2](#)
- [5.2.6.2](#)

### 7.3.6 Clause 29: Console Device Conformance

A console device MUST conform to the following normative statements:

- [5.3.5.1](#)
- [5.3.6.2.1](#)

### 7.3.7 Clause 30: Entropy Device Conformance

An entropy device MUST conform to the following normative statements:

- [5.4.6.2](#)



### 7.3.8 Clause 31: Traditional Memory Balloon Device Conformance

A traditional memory balloon device MUST conform to the following normative statements:

- [5.5.3.2](#)
- [5.5.6.2](#)
- [5.5.6.3.2](#)
- [5.5.6.5.2](#)
- [5.5.6.6.2](#)
- [5.5.6.7.2](#)

### 7.3.9 Clause 32: SCSI Host Device Conformance

An SCSI host device MUST conform to the following normative statements:

- [5.6.4.2](#)
- [5.6.5](#)
- [5.6.6.1.1](#)
- [5.6.6.3.2](#)

### 7.3.10 Clause 33: Input Device Conformance

An input device MUST conform to the following normative statements:

- [5.8.5.2](#)
- [5.8.6.2](#)

### 7.3.11 Clause 34: Crypto Device Conformance

A Crypto device MUST conform to the following normative statements:

- [5.9.5.1](#)
- [5.9.7.2.1.6](#)
- [5.9.7.2.1.8](#)
- [5.9.7.4.2](#)
- [5.9.7.5.2](#)
- [5.9.7.6.2](#)
- [5.9.7.7.2](#)

### 7.3.12 Clause 35: Socket Device Conformance

A socket device MUST conform to the following normative statements:

- [5.10.6.3.2](#)
- [5.10.6.4.2](#)

### 7.3.13 Clause 36: File System Device Conformance

A file system device MUST conform to the following normative statements:

- [5.11.4.2](#)
- [5.11.6.2.1](#)

- [5.11.6.4.1](#)

### **7.3.14 Clause 37: RPMB Device Conformance**

An RPMB device MUST conform to the following normative statements:

- [5.12.5](#)
- [5.12.6.1.1](#)
- [5.12.6.1.2](#)
- [5.12.6.1.3](#)
- [5.12.6.1.4](#)
- [5.12.6.1.5](#)
- [5.12.6.3](#)

### **7.3.15 Clause 38: IOMMU Device Conformance**

An IOMMU device MUST conform to the following normative statements:

- [5.13.3.2](#)
- [5.13.4.2](#)
- [5.13.6.2](#)
- [5.13.6.3.2](#)
- [5.13.6.4.2](#)
- [5.13.6.5.2](#)
- [5.13.6.6.2](#)
- [5.13.6.7.2](#)
- [5.13.6.8.1.2](#)
- [5.13.6.9.2](#)

### **7.3.16 Clause 39: Sound Device Conformance**

A sound device MUST conform to the following normative statements:

- [5.14.6.4.1.1](#)
- [5.14.6.6.2.1](#)
- [5.14.6.6.3.1](#)
- [5.14.6.6.5.1](#)
- [5.14.6.8.1.1](#)
- [5.14.6.8.2.1](#)
- [5.14.6.9.1.1](#)

### **7.3.17 Clause 40: Memory Device Conformance**

A memory device MUST conform to the following normative statements:

- [5.15.4.2](#)
- [5.15.5.2](#)
- [5.15.6.2](#)

- [5.15.6.3.2](#)
- [5.15.6.4.2](#)
- [5.15.6.5.2](#)
- [5.15.6.6.2](#)

### 7.3.18 Clause 41: I2C Adapter Device Conformance

An I2C Adapter device MUST conform to the following normative statements:

- [5.16.6.4](#)

### 7.3.19 Clause 42: SCMI Device Conformance

An SCMI device MUST conform to the following normative statements:

- [5.17.3.1](#)
- [5.17.6.1.1](#)
- [5.17.6.3.1](#)
- [5.17.6.4.1](#)

### 7.3.20 Clause 43: GPIO Device Conformance

A General Purpose Input/Output (GPIO) device MUST conform to the following normative statements:

- [5.18.6.9](#)
- [5.18.7.3](#)

### 7.3.21 Clause 44: PMEM Device Conformance

A PMEM device MUST conform to the following normative statements:

- [5.19.5.1](#)
- [5.19.7.1](#)
- [5.19.7.3](#)

## 7.4 Clause 45: Legacy Interface: Transitional Device and Transitional Driver Conformance

A conformant implementation MUST be either transitional or non-transitional, see [1.3.1](#).

An implementation MAY choose to implement OPTIONAL support for the legacy interface, including support for legacy drivers or devices, by conforming to all of the MUST or REQUIRED level requirements for the legacy interface for the transitional devices and drivers.

The requirements for the legacy interface for transitional implementations are located in sections named "Legacy Interface" listed below:

- Section [2.2.3](#)
- Section [2.5.3](#)
- Section [2.5.4](#)
- Section [2.7.2](#)
- Section [2.7.3](#)
- Section [2.7.4.3](#)

- [Section 3.1.2](#)
- [Section 4.1.2.3](#)
- [Section 4.1.4.10](#)
- [Section 4.1.5.1.1.1](#)
- [Section 4.1.5.1.3.1](#)
- [Section 4.2.4](#)
- [Section 4.3.2.1.3](#)
- [Section 4.3.2.2.2](#)
- [Section 4.3.3.1.3](#)
- [Section 4.3.2.6.4](#)
- [Section 5.1.3.2](#)
- [Section 5.1.4.3](#)
- [Section 5.1.6.1](#)
- [Section 5.1.6.5.2.3](#)
- [Section 5.1.6.5.3.1](#)
- [Section 5.1.6.5.6.3](#)
- [Section 5.1.6.5.8.3](#)
- [Section 5.2.3.1](#)
- [Section 5.2.4.1](#)
- [Section 5.2.5.3](#)
- [Section 5.2.6.3](#)
- [Section 5.3.4.1](#)
- [Section 5.3.6.3](#)
- [Section 5.5.3.2.0.1](#)
- [Section 5.5.6.2.1](#)
- [Section 5.5.6.3.3](#)
- [Section 5.6.4.3](#)
- [Section 5.6.6.0.1](#)
- [Section 5.6.6.1.3](#)
- [Section 5.6.6.2.1](#)
- [Section 5.6.6.3.3](#)
- [Section 6.3](#)

---

## Appendix A. virtio\_queue.h

This file is also available at the link [https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/listings/virtio\\_queue.h](https://docs.oasis-open.org/virtio/virtio/v1.1/cs01/listings/virtio_queue.h). All definitions in this section are for non-normative reference only.

```
#ifndef VIRTQUEUE_H
#define VIRTQUEUE_H
/* An interface for efficient virtio implementation.
 *
 * This header is BSD licensed so anyone can use the definitions
 * to implement compatible drivers/servers.
 *
 * Copyright 2007, 2009, IBM Corporation
 * Copyright 2011, Red Hat, Inc
 * All rights reserved.
 *
 * Redistribution and use in source and binary forms, with or without
 * modification, are permitted provided that the following conditions
 * are met:
 * 1. Redistributions of source code must retain the above copyright
 * notice, this list of conditions and the following disclaimer.
 * 2. Redistributions in binary form must reproduce the above copyright
 * notice, this list of conditions and the following disclaimer in the
 * documentation and/or other materials provided with the distribution.
 * 3. Neither the name of IBM nor the names of its contributors
 * may be used to endorse or promote products derived from this software
 * without specific prior written permission.
 * THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS ``AS IS'' AND
 * ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE
 * IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE
 * ARE DISCLAIMED. IN NO EVENT SHALL IBM OR CONTRIBUTORS BE LIABLE
 * FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL
 * DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS
 * OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION)
 * HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT
 * LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY
 * OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF
 * SUCH DAMAGE.
 */
#include <stdint.h>

/* This marks a buffer as continuing via the next field. */
#define VIRTQ_DESC_F_NEXT 1
/* This marks a buffer as write-only (otherwise read-only). */
#define VIRTQ_DESC_F_WRITE 2
/* This means the buffer contains a list of buffer descriptors. */
#define VIRTQ_DESC_F_INDIRECT 4

/* The device uses this in used->flags to advise the driver: don't kick me
 * when you add a buffer. It's unreliable, so it's simply an
 * optimization. */
#define VIRTQ_USED_F_NO_NOTIFY 1
/* The driver uses this in avail->flags to advise the device: don't
 * interrupt me when you consume a buffer. It's unreliable, so it's
 * simply an optimization. */
#define VIRTQ_AVAIL_F_NO_INTERRUPT 1

/* Support for indirect descriptors */
#define VIRTIO_F_INDIRECT_DESC 28

/* Support for avail_event and used_event fields */
#define VIRTIO_F_EVENT_IDX 29
```

```

/* Arbitrary descriptor layouts. */
#define VIRTIO_F_ANY_LAYOUT 27

/* Virtqueue descriptors: 16 bytes.
 * These can chain together via "next". */
struct virtq_desc {
    /* Address (guest-physical). */
    le64 addr;
    /* Length. */
    le32 len;
    /* The flags as indicated above. */
    le16 flags;
    /* We chain unused descriptors via this, too */
    le16 next;
};

struct virtq_avail {
    le16 flags;
    le16 idx;
    le16 ring[];
    /* Only if VIRTIO_F_EVENT_IDX: le16 used_event; */
};

/* le32 is used here for ids for padding reasons. */
struct virtq_used_elem {
    /* Index of start of used descriptor chain. */
    le32 id;
    /* Total length of the descriptor chain which was written to. */
    le32 len;
};

struct virtq_used {
    le16 flags;
    le16 idx;
    struct virtq_used_elem ring[];
    /* Only if VIRTIO_F_EVENT_IDX: le16 avail_event; */
};

struct virtq {
    unsigned int num;

    struct virtq_desc *desc;
    struct virtq_avail *avail;
    struct virtq_used *used;
};

static inline int virtq_need_event(uint16_t event_idx, uint16_t new_idx, uint16_t old_idx)
{
    return (uint16_t)(new_idx - event_idx - 1) < (uint16_t)(new_idx - old_idx);
}

/* Get location of event indices (only with VIRTIO_F_EVENT_IDX) */
static inline le16 *virtq_used_event(struct virtq *vq)
{
    /* For backwards compat, used event index is at *end* of avail ring. */
    return &vq->avail->ring[vq->num];
}

static inline le16 *virtq_avail_event(struct virtq *vq)
{
    /* For backwards compat, avail event index is at *end* of used ring. */
    return (le16 *)&vq->used->ring[vq->num];
}
#endif /* VIRTQUEUE_H */

```

---

## Appendix B. Creating New Device Types

Various considerations are necessary when creating a new device type.

### B.1 How Many Virtqueues?

It is possible that a very simple device will operate entirely through its device configuration space, but most will need at least one virtqueue in which it will place requests. A device with both input and output (eg. console and network devices described here) need two queues: one which the driver fills with buffers to receive input, and one which the driver places buffers to transmit output.

### B.2 What Device Configuration Space Layout?

Device configuration space should only be used for initialization-time parameters. It is a limited resource with no synchronization between fields written by the driver, so for most uses it is better to use a virtqueue to update configuration information (the network device does this for filtering, otherwise the table in the config space could potentially be very large).

Remember that configuration fields over 32 bits wide might not be atomically writable by the driver. Therefore, no writeable field which triggers an action ought to be wider than 32 bits.

### B.3 What Device Number?

Device numbers can be reserved by the OASIS committee: email [virtio-dev@lists.oasis-open.org](mailto:virtio-dev@lists.oasis-open.org) to secure a unique one.

Meanwhile for experimental drivers, use 65535 and work backwards.

### B.4 How many MSI-X vectors? (for PCI)

Using the optional MSI-X capability devices can speed up interrupt processing by removing the need to read ISR Status register by guest driver (which might be an expensive operation), reducing interrupt sharing between devices and queues within the device, and handling interrupts from multiple CPUs. However, some systems impose a limit (which might be as low as 256) on the total number of MSI-X vectors that can be allocated to all devices. Devices and/or drivers should take this into account, limiting the number of vectors used unless the device is expected to cause a high volume of interrupts. Devices can control the number of vectors used by limiting the MSI-X Table Size or not presenting MSI-X capability in PCI configuration space. Drivers can control this by mapping events to as small number of vectors as possible, or disabling MSI-X capability altogether.

### B.5 Device Improvements

Any change to device configuration space, or new virtqueues, or behavioural changes, should be indicated by negotiation of a new feature bit. This establishes clarity<sup>1</sup> and avoids future expansion problems.

Clusters of functionality which are always implemented together can use a single bit, but if one feature makes sense without the others they should not be gratuitously grouped together to conserve feature bits.

---

<sup>1</sup>Even if it does mean documenting design or implementation mistakes!

---

## Appendix C. Acknowledgements

The following individuals have participated in the creation of this specification and are gratefully acknowledged:

### Participants

Allen Chia, Oracle  
Amit Shah, Red Hat  
Amos Kong, Red Hat  
Anthony Liguori, IBM  
Bruce Rogers, SUSE  
Bryan Venteicher, NetApp  
Chandra Thyamagondlu, Xilinx  
Chet Ensign, OASIS  
Cornelia Huck, Red Hat  
Cunming, Liang, Intel  
Damjan, Marion, Cisco  
Daniel Kiper, Oracle  
Fang Chen, Huawei  
Fang You, Huawei  
Geoff Brown, M2Mi  
Gerd Hoffmann, Red Hat  
Gershon Janssen, Individual Member  
Grant Likely, ARM  
Haggai Eran, Mellanox  
Halil Pasic, IBM  
James Bottomley, Parallels IP Holdings GmbH  
Jani Kokkonen, Huawei  
Jan Kiszka, Siemens AG  
Jens Freimann, Red Hat  
Jian Zhou, Huawei  
Karen Xie, Xilinx  
Kumar Sanghvi, Xilinx  
Lei Gong, Huawei  
Lior Narkis, Mellanox  
Luiz Capitulino, Red Hat  
Marc-André Lureau, Red Hat  
Mark Gray, Intel  
Michael S. Tsirkin, Red Hat  
Mihai Carabas, Oracle  
Nishank Trivedi, NetApp  
Paolo Bonzini, Red Hat  
Paul Mundt, Huawei  
Pawel Moll, ARM  
Peng Long, Huawei  
Piotr Uminski, Intel  
Qian Xum, Intel  
Richard Sohn, Alcatel-Lucent  
Rusty Russell, IBM



Sasha Levin, Oracle  
Sergey Tverdyshev, Thales e-Security  
Stefan Hajnoczi, Red Hat  
Sundar Mohan, Xilinx  
Tom Lyon, Samya Systems, Inc.  
Victor Kaplansky, Red Hat  
Vijay Balakrishna, Oracle  
Wei Wang, Intel  
Xin Zeng, Intel

The following non-members have provided valuable feedback on this specification and are gratefully acknowledged:

## Reviewers

Aaron Conole, Red Hat  
Adam Tao, Huawei  
Alexander Duyck, Intel  
Andreas Pape, ADITG/ESB  
Andrew Thornton, Google  
Arun Subbarao, LynuxWorks  
Baptiste Reynal, Virtual Open Systems  
Bharat Bhushan, NXP Semiconductors  
Brian Foley, ARM  
Chandra Thyamagondlu, Xilinx  
Changpeng Liu, Intel  
Christian Pinto, Virtual Open Systems  
Christoffer Dall, ARM  
Christoph Hellwig, Individual  
Christian Borntraeger, IBM  
Daniel Marcovitch, Mellanox  
David Alan Gilbert, Red Hat  
David Hildenbrand, Red Hat  
David Riddoch, Solarflare  
Denis V. Lunev, OpenVZ  
Dmitry Fleytman, Red Hat  
Don Wallwork, Broadcom  
Emily Drea, ARM  
Eric Auger, Red Hat  
Fam Zheng, Red Hat  
Francesco Fusco, Red Hat  
Frank Yang, Google  
Gil Savir, Intel  
Gonglei (Arei), Huawei  
Greg Kurz, IBM  
Hannes Reiencke, SUSE  
Ian Campbell, Docker  
Ilya Lesokhin, Mellanox  
Jacques Durand, Fujitsu  
Jakub Jermar, Kernkonzept  
Jan Scheurich, Ericsson  
Jason Baron, Akamai  
Jason Wang, Red Hat  
Jean-Philippe Brucker, ARM  
Jianfeng Tan, intel  
Jonathan Helman, Oracle  
Karandeep Chahal, DDN

Kevin Lo, MSI  
Kevin Tian, Intel  
Kully Dhanoa, Intel  
Laura Novich, Red Hat  
Ladi Prosek, Red Hat  
Lars Ganrot, Napatech  
Longpeng (Mike), Huawei  
Mario Torrecillas Rodriguez, ARM  
Mark Rustad, Intel  
Maxime Coquelin, Red Hat  
Namhyung Kim, LG  
Ola Liljedahl, ARM  
Pankaj Gupta, Red Hat  
Patrick Durusau, OASIS  
Pierre Pfister, Cisco  
Pranavkumar Sawargaonkar, Linaro  
Rauchfuss Holm, Huawei  
Rob Miller, Broadcom  
Roman Kiryanov, Google  
Robin Cover, OASIS  
Roger S Chien, Intel  
Sameeh Jubran, Red Hat / Daynix  
Si-Wei Liu, Oracle  
Sridhar Samudrala, Intel  
Stefan Fritsch, Individual  
Stefano Garzarella, Red Hat  
Steven Luong, Cisco  
Thomas Huth, Red Hat  
Tiwei Bie, Intel  
Tomáš Golembiovský, Red Hat  
Venu Busireddy, Oracle  
Victor Kaplansky, Red Hat  
Vijayabhaskar Balakrishna, Oracle  
Vlad Yasevich, Red Hat  
Yan Vugenfirer, Red Hat / Daynix  
Wei Xu, Red Hat  
Will Deacon, ARM  
Willem de Bruijn, Google  
Yuanhan Liu, Intel  
Yuri Benditovich, Red Hat / Daynix  
Zhi Yong Wu, IBM  
Zhoujian, Huawei

---

## Appendix D. Revision History

The following changes have been made since the previous version of this specification:

Revision	Date	Editor	Changes Made
5c43ad7	27 Feb 2019	Halil Pasic	<p>ccw: be more precise about the semantic of revision 1</p> <p>Revision 1 of the CCW transport is currently defined as virtio 1.0. This could become confusing when we bump the version of the virtio specification to 1.1, in a sense that it could be interpreted like one can not use any features not part of the 1.0 specification.</p> <p>So let us try to avoid confusion regarding the semantic of virtio-ccw revision 1.</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-163">https://issues.oasis-open.org/browse/VIRTIO-163</a></p> <p>Signed-off-by: Halil Pasic &lt;pasic@linux.vnet.ibm.com&gt;</p> <p>Reviewed-by: Cornelia Huck &lt;cohuck@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See <a href="#">4.3.2.1</a>.</p>
d348ac0	27 Feb 2019	Halil Pasic	<p>introduction: simplify the designation of legacy</p> <p>The sentence designating the documents defining what later became known as the legacy virtio interface had the most important piece of information placed in parenthesis.</p> <p>Let's reword this sentence so we avoid using an ambiguous designation based on a relative anchor (i.e. 'earlier drafts of this specification') and just use the absolute anchor (version 1.0).</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-164">https://issues.oasis-open.org/browse/VIRTIO-164</a></p> <p>Signed-off-by: Halil Pasic &lt;pasic@linux.vnet.ibm.com&gt;</p> <p>Reviewed-by: Cornelia Huck &lt;cohuck@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See <a href="#">1.3.1</a>.</p>

Revision	Date	Editor	Changes Made
bef3ff7	07 Mar 2019	Stefan Hajnoczi	<p>virtio-blk: document data[] size constraints</p> <p>The struct virtio_blk_req-&gt;data[] field is a multiple of 512 bytes long for read and write requests. Flush requests don't use data[] at all.</p> <p>The new discard and write zeroes requests being introduced in VIRTIO 1.1 put struct virtio_blk_discard_write_zeroes elements into data[], so it must be a multiple of the struct size.</p> <p>The uint8_t data[][512] pseudo-code makes it look like discard and write zeroes requests must pad to 512 bytes. This wastes memory since struct virtio_blk_discard_write_data is only 16 bytes long.</p> <p>Furthermore, all known implementations wishing to take advantage of this upcoming VIRTIO 1.1 feature do not use 512-byte padding (Linux virtio_blk.ko, QEMU virtio-blk device emulation, the SPDK virtio-blk driver, and the SPDK vhost-user-blk device backend).</p> <p>This patch documents the data[] size constraints clearly in the driver normative section. This is clearer than the current pseudo-code.</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/32">https://github.com/oasis-tcs/virtio-spec/issues/32</a></p> <p>Cc: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Cc: Changpeng Liu &lt;changpeng.liu@intel.com&gt;</p> <p>Cc: Stefano Garzarella &lt;sgarzare@redhat.com&gt;</p> <p>Signed-off-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See 5.2.6.</p>
c5c0ce7	07 Mar 2019	Stefan Hajnoczi	<p>virtio-blk: move virtio_blk_discard_write_zeroes definition</p> <p>struct virtio_blk_discard_write_zeroes is defined alongside struct virtio_blk_req but only discussed later in the text. Move it to where it belongs.</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/32">https://github.com/oasis-tcs/virtio-spec/issues/32</a></p> <p>Suggested-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Signed-off-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See 5.2.6.</p>

Revision	Date	Editor	Changes Made
caffe5c	07 Mar 2019	Stefan Hajnoczi	<p>virtio-blk: describe write zeroes unmap semantics</p> <p>Explain the meaning of the unmap flag. The details are already covered in the device normative section but mentioning it here makes the text easier to understand.</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/32">https://github.com/oasis-tcs/virtio-spec/issues/32</a></p> <p>Suggested-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Signed-off-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See 5.2.6.</p>
5f1e981	07 Mar 2019	Stefan Hajnoczi	<p>virtio-blk: avoid inconsistent "DISCARD" term</p> <p>"discard" (lowercase) is used throughout the text. Remove a lone instance of "DISCARD" (uppercase).</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/32">https://github.com/oasis-tcs/virtio-spec/issues/32</a></p> <p>Suggested-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Signed-off-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See 5.2.6.</p>
31a52d2	07 Mar 2019	Stefan Hajnoczi	<p>virtio-blk: clarify semantics of multi-segment discard/write zeroes commands</p> <p>Describe the failure case and maximum number of segments in a multi-segment discard/write zeroes command.</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/34">https://github.com/oasis-tcs/virtio-spec/issues/34</a></p> <p>Signed-off-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See 5.2.6.1.</p>
90047f5	21 Mar 2019	Michael S. Tsirkin	<p>format: replace "- i.e." with ", i.e.,"</p> <p>This seems to be preferred by native speakers, and seems just as effective as a sentence device.</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-171">https://issues.oasis-open.org/browse/VIRTIO-171</a></p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Reviewed-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p>

Revision	Date	Editor	Changes Made
c7b2503	21 Mar 2019	Michael S. Tsirkin	conformance: add links to crypto and input devices Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-174">https://issues.oasis-open.org/browse/VIRTIO-174</a> Signed-off-by: Michael S. Tsirkin <mst@redhat.com> Reviewed-by: Stefan Hajnoczi <stefanha@redhat.com> See 7.1.
0b5288f	21 Mar 2019	Michael S. Tsirkin	signal start and end of structures consistently Make sure all structs have the format: struct X . ... }; Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-170">https://issues.oasis-open.org/browse/VIRTIO-170</a> Signed-off-by: Michael S. Tsirkin <mst@redhat.com> Reviewed-by: Stefan Hajnoczi <stefanha@redhat.com> Reviewed-by: Cornelia Huck <cohuck@redhat.com> See 5.6.6.2.
69daf06	21 Mar 2019	Michael S. Tsirkin	editorial: explain each structure before use Several structures are listed before they are introduced in some way. Add a sentence before each one so they don't appear prior to any prose. Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-166">https://issues.oasis-open.org/browse/VIRTIO-166</a> Signed-off-by: Michael S. Tsirkin <mst@redhat.com> Reviewed-by: Stefan Hajnoczi <stefanha@redhat.com> Reviewed-by: Cornelia Huck <cohuck@redhat.com> See 5.1.6.5, 2.7.8, 5.9.5, 5.7.4, 5.7.6.7, 5.7.6.7 and 5.10.4.
d608f47	21 Mar 2019	Michael S. Tsirkin	conformance: tweak to match OASIS requirements Number clauses as required by OASIS. Also reference the transitional clause. Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-168">https://issues.oasis-open.org/browse/VIRTIO-168</a> Signed-off-by: Michael S. Tsirkin <mst@redhat.com> Reviewed-by: Stefan Hajnoczi <stefanha@redhat.com> Reviewed-by: Cornelia Huck <cohuck@redhat.com> See 7.

Revision	Date	Editor	Changes Made
0dbd52d	21 Mar 2019	Michael S. Tsirkin	<p>introduction: update link to IEEE 802</p> <p>Looks like all GETIEEE links got broken. Let's just point to their main page.</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-175">https://issues.oasis-open.org/browse/VIRTIO-175</a></p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Reviewed-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Reviewed-by: Jens Freimann &lt;jfreimann@redhat.com&gt;</p> <p>See <a href="#">1.1</a>.</p>
7b361ea	21 Mar 2019	Michael S. Tsirkin	<p>editorial: upgrade links to https</p> <p>Several links have been upgraded and now redirect to the https version. Upgrade our version accordingly.</p> <p>Note that some other links use the status 301 - moved permanently apparently in error (e.g. for a language specific redirect), not updating these.</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-173">https://issues.oasis-open.org/browse/VIRTIO-173</a></p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Reviewed-by: Stefan Hajnoczi &lt;stefanha@redhat.com&gt;</p> <p>Reviewed-by: Jens Freimann &lt;jfreimann@redhat.com&gt;</p> <p>See .</p>
3e49aec	21 Mar 2019	Michael S. Tsirkin	<p>conformance: fix confusion about legacy interface</p> <p>The text describing the legacy interface also obliquely refers to a non-transitional implementation. This seems to cause confusion and there's no good reason to do it here: this section is about legacy interface and transitional devices, it add not value at all. Just drop it.</p> <p>Note: the spec does not make it clear whether description of the legacy interface is normative or not, and in particular, this section is not linked to from any conformance targets. Resolving that is left for later.</p> <p>Fixes: <a href="https://issues.oasis-open.org/browse/VIRTIO-167">https://issues.oasis-open.org/browse/VIRTIO-167</a></p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>Reviewed-by: Cornelia Huck &lt;cohuck@redhat.com&gt;</p> <p>Acked-by: Halil Pasic &lt;pasic@linux.ibm.com&gt;</p> <p>See <a href="#">7.4</a>.</p>

Revision	Date	Editor	Changes Made
4cc8a4d	21 Mar 2019	Michael S. Tsirkin	<p>block: drop duplicate text</p> <p>In version 1.1 draft 01 - Section 5.2.6.4 - second bullet:</p> <p>Duplicated text "errors, data_len, sense_len and residual MUST reside in a single, separate device-writable descriptor" appears +both in the beginning and at the end of the 2nd sentence.</p> <p>The original text:</p> <p>For SCSI commands there are additional constraints. errors, data_len, sense_len and residual MUST reside in a single, separate device-writable descriptor, sense MUST reside in a single separate device-writable descriptor of size 96 bytes, and errors, data_len, sense_len and residual MUST reside a single separate device-writable descriptor. I suggest to delete the 1st one, so in the end result, fields are described in same order as appear in struct virtio_scsi_pc_req.</p> <p>Fixes: <a href="https://github.com/oasis-tcs/virtio-spec/issues/39">https://github.com/oasis-tcs/virtio-spec/issues/39</a></p> <p>Reported-by: Gil Savir &lt;gil.savir@intel.com&gt;</p> <p>Signed-off-by: Michael S. Tsirkin &lt;mst@redhat.com&gt;</p> <p>See <a href="#">5.2.6.4</a>.</p>