
Information Technology – SCSI RDMA Protocol (SRP)

This is an internal working document of T10, a Technical Committee of Accredited Standards Committee INCITS (InterNational Committee for Information Technology Standards). As such this is not a completed standard and has not been approved. The contents may be modified by the T10 Technical Committee. The contents are actively being modified by T10. This document is made available for review and comment only.

Permission is granted to members of INCITS, its technical committees, and their associated task groups to reproduce this document for the purposes of INCITS standardization activities without further permission, provided this notice is included. All other rights are reserved. Any duplication of this document for commercial or for-profit use is strictly prohibited.

T10 Technical editor:

Cris Simpson
Intel Corporation
2111 NE 25th Ave.
Hillsboro, Oregon 97124
USA

Telephone: +1.503.712.4333
Facsimile: +1.503.712.2200
Email: cris.simpson@intel.com

Reference number

ANSI INCITS.***:200x

POINTS OF CONTACT

T10 Chair

John Lohmeyer
LSI Logic
4420 Arrows West Drive
Colorado Springs, CO 80907-3444
USA

Telephone: +1.719.533.7560
Facsimile: +1.719.533.7183
Email: lohmeier@t10.org

INCITS Secretariat

INCITS Secretariat
1250 Eye Street, NW Suite 200
Washington, DC 20005
<http://www.incits.org/>

T10 Web Site

www.t10.org

T10 Reflector

To subscribe send email to majordomo@T10.org with 'subscribe' in message body
To unsubscribe send email to majordomo@T10.org with 'unsubscribe' in message body
Internet address for distribution via T10 reflector: T10@T10.org

Document Distribution

NCITS Online Store
managed by Techstreet
1327 Jones Drive
Ann Arbor, MI 48105

Global Engineering
15 Inverness Way East
Englewood, CO 80112-5704

T10 Vice Chair

George O. Penokie
IBM / Tivoli Systems
3605 Highway 52 North
Rochester, MN 55901
USA

Telephone: +1.507.253.5208
Facsimile: +1.507.253.2880
Email: gop@us.ibm.com

Telephone: +1.202.737.8888
Facsimile: +1.202.638.4922
Email: ncits@itic.org

Web: <http://www.techstreet.com/ncits.html>
Telephone: +1.734.302.7801 or
+1.800.699.9277
Facsimile: +1.734.302.7811

Web: <http://global.ihs.com>
Telephone: +1.303.792.2181 or
+1.800.854.7179
Facsimile: +1.303.792.2192

Revision History

Revision 16a (3 July 2002)

- a) Removed change bars, ~~deleted~~ text, line numbers.
- b) Unmarked ~~new~~ text.
- c) On advice of T10 vice-chair, reduced list of example standards in clause 1.

Revision 16 (2 July 2002)

- a) Completed changes resolving letter ballot comments.
- b) Included I_T nexus loss [01-232r0] definition.
- c) Comment resolution details may be found in T10/01-328r8.

Revision 15 (26 April 2002)

- a) Split SRP working draft (this document) from letter ballot comment resolution document T10/01-328r7.
- b) Reset change bars. Change bars in this document represent changes made since SRPr14.
- c) In-progress changes resolving letter ballot comments.

Revision 14 (17 April 2002)

In-progress changes resolving letter ballot comments.

Revision 13 (4 April 2002)

- a) In-progress changes resolving letter ballot comments.
- b) Removed Annex C.

Revision 12 (11 March 2002)

In-progress changes resolving letter ballot comments.

Revision 11 (14 February 2002)

In-progress changes resolving letter ballot comments.

Special thanks to Ed Gardner, Ophidian Designs, for his work as the original editor for SRP.

Revision 10 (3 October 2001)

- a) [01-289r0] Comments from 24 September 2001 SRP teleconference.
- b) [01-298r1] Comments from 28 September 2001 SRP teleconference.
- c) Reformatted SRP to InfinibandTM annex.

Revision 09 (12 September 2001)

- a) [01-230r2] SRP buffer descriptor rewrite;
- b) [01-250r1] SRP operation overview;
- c) [01-263r0] Comments from August 21 SRP teleconference; and
- d) Uniform use of "SRP target port", "SRP initiator port", "RDMA channel" and "IB channel".

Revision 08 (8 August 2001)

- a) [01-028r6] SRP InfinibandTM annex;
- b) [01-193r1] SRP alias entry designation formats (with extensive editorial changes);
- c) [01-205r1] SRP Initiator Logout proposal;
- d) [01-177r2] SRP model for RDMA communication services; and
- e) [01-172r4] SRP to SAM-2 protocol.

Working Draft

Revision 07 (17 July 2001)

- a) [01-195] Changes from June 19-20 SRP working group minutes; and
- b) Corrections described in June 21 T10 reflector message from Kamran_Tavakoli@adaptec.com.

Revision 06 (14 June 2001)

- a) [01-171r0] SRP_LOGOUT_REJECT, as modified during the May 25 teleconference (see 01-178);
- b) [01-173r1] SRP bidirectional residuals, as modified during the May 25 teleconference (see 01-178);
- c) Other changes approved during the May 25 teleconference (see 01-178);
- d) Reconciled SRP_AER_REQ format to match revised SRP_RSP;
- e) Reconciled SRP_TASK_MGMT format to match current SRP_CMD; and
- f) Editorial changes and minor corrections in response to comments received on previous revisions.

Revision 05 (23 May 2001)

Numerous editorial changes. No intentional technical changes.

Revision 04 (10 May 2001)

Added mode pages, residual count clarification, AER, scatter / gather revision, total transfer length, logout, target / initiator port identifiers in login. Removed VI terminology, target reset, multiple command IUs. Believed to contain all approved changes through May 3 working group other than those listed above.

Revision 03 (29 January 2001)

Added RDMA Communication Model description. Fixed editorial errors in command IUs (restored bytes 4 to 7, three dots).

Revision 02 (4 January 2001)

Incorporates 00-354r2, scatter/gather and IU format changes defined at November 29-30 SRP working group (see 01-009r0), name changed to SRP, partial changes to use non-VI terminology.

Revision 01 (7 July 2000)

First semi-complete draft. Based on 99-316r1, 00-172r0 and 00-240r0. Tags expanded from 16 to 32 bits. TRD COUNT renamed REQUESTLIMIT and expanded to 32 bits. SVP_CMD and SVP_RSP IUs expanded to accomodate these fields and provide additional reserved words. Defined IU maximum size negotiation. Changed order of data transfer descriptor to match the order in Infiniband RDMA transport header.

Revision 00 (17 May 2000)

Partial draft.

Working Draft

**ANSI (r)
INCITS.***:200x**

**American National Standard for Information Systems –
Information Technology –
SCSI RDMA Protocol (SRP)**

Secretariat

InterNational Committee for Information Technology Standards

Approved mm dd yy

American National Standards Institute, Inc.

ABSTRACT

This standard describes the message format and protocol definitions required to transfer commands and data between a SCSI (Small Computer System Interface) initiator port and a SCSI target port using an RDMA communication service.

Working Draft

American National Standard

Approval of an American National Standard requires verification by ANSI that the requirements for due process, consensus, and other criteria for approval have been met by the standards developer. Consensus is established when, in the judgment of the ANSI Board of Standards Review, substantial agreement has been reached by directly and materially affected interests. Substantial agreement means much more than a simple majority, but not necessarily unanimity. Consensus requires that all views and objections be considered, and that effort be made towards their resolution.

The use of American National Standards is completely voluntary; their existence does not in any respect preclude anyone, whether he has approved the standards or not, from manufacturing, marketing, purchasing, or using products, processes, or procedures not conforming to the standards.

The American National Standards Institute does not develop standards and will in no circumstances give interpretation on any American National Standard. Moreover, no person shall have the right or authority to issue an interpretation of an American National Standard in the name of the American National Standards Institute. Requests for interpretations should be addressed to the secretariat or sponsor whose name appears on the title page of this standard.

CAUTION NOTICE: This American National Standard may be revised or withdrawn at any time. The procedures of the American National Standards Institute require that action be taken periodically to reaffirm, revise, or withdraw this standard. Purchasers of American National Standards may receive current information on all standards by calling or writing the American National Standards Institute.

CAUTION: The developers of this standard have requested that holders of patents that may be required for the implementation of the standard, disclose such patents to the publisher. However, neither the developers nor the publisher have undertaken a patent search in order to identify which, if any, patents may apply to this standard. As of the date of publication of this standard and following calls for the identification of patents that may be required for the implementation of the standard, no such claims have been made.

No further patent search is conducted by the developer or the publisher in respect to any standard it processes. No representation is made or implied that licenses are not required to avoid infringement in the use of this standard.

Foreword	xii
Introduction	xiii
1 Scope	1
2 Normative references	2
2.1 Normative references	2
2.2 Approved references	2
2.3 References under development	2
3 Definitions, symbols, abbreviations and conventions	3
3.1 Definitions	3
3.2 Acronyms	4
3.3 Keywords	4
3.4 Conventions	5
3.5 Notation for procedures and functions	6
4 RDMA communication service model	8
4.1 Overview	8
4.2 RDMA Channels	8
4.2.1 Introduction	8
4.2.2 Establishment	9
4.2.3 Disestablishment	11
4.3 Messages	11
4.4 RDMA operations	11
4.4.1 Overview	11
4.4.2 RDMA Write	12
4.4.3 RDMA Read	12
4.5 Ordering and Reliability	12
4.5.1 Ordering and reliability overview	12
4.5.2 Reliability	12
4.5.3 Ordering	13
5 Structure and concepts	14
5.1 Overview of SRP operation	14
5.1.1 RDMA channel establishment and login	14
5.1.2 Single RDMA channel operation	14
5.1.3 Multiple independent RDMA channel operation	14
5.1.4 RDMA channel disconnection	15
5.2 Identifiers	16
5.3 Alias associations	16
5.4 Information unit classes	16
5.5 SRP target port buffer management	16
5.5.1 Buffer management overview	16
5.5.2 SRP requests issued by target port	16
5.5.3 Requests issued by initiator port	16
5.6 Data buffers	18
5.6.1 Memory descriptors	18
5.6.2 Data buffer descriptors	19
6 SRP Information Units	24
6.1 Summary	24

7	SCSI mode parameters	50
7.1	SCSI mode parameter overview and codes	50
7.2	Disconnect-reconnect mode page	50
7.2.1	Valid fields	51
7.2.2	Invalid fields	51
7.3	Protocol specific LUN page	51
7.4	Protocol specific port page	51
Annex A	SRP interface protocol and services	52
A.1	Service interface protocol	52
A.2	SRP services	54
A.3	SAM-2 object mapping	54
A.4	Procedure objects	54
A.5	Application client SCSI command services	55
A.5.1	Application client SCSI command services overview	55
A.5.2	Send SCSI command service	55
A.6	Device server SCSI command services	56
A.6.1	Device server SCSI command services overview	56
A.6.2	Data-out delivery service	56
A.6.3	Data-in delivery service	57
A.7	Task management services	57
A.7.1	Task management functions overview	57
A.7.2	Task management functions	58
A.7.3	ABORT TASK	58
A.7.4	ABORT TASK SET	58
A.7.5	CLEAR ACA	58
A.7.6	CLEAR TASK SET	58
A.7.7	LOGICAL UNIT RESET	58
A.7.8	TARGET RESET	58
A.7.9	WAKEUP	58
Annex B	SRP for the InfiniBand™ Architecture	60
B.1	Overview	60
B.2	Normative references	60
B.3	Definitions and abbreviations	60
B.3.1	Introduction to definitions and abbreviations	60
B.3.2	Definitions	60
B.3.3	Abbreviations	61
B.4	InfiniBand™ Architecture overview	61
B.5	SCSI architecture mapping	64
B.6	Communication management	65
B.6.1	Communication management overview	65
B.6.2	Discovering SRP target ports	65
B.6.3	Establishing a connection	65
B.6.4	Releasing a connection	66
B.6.5	Errors	66
B.6.6	Data-out and data-in operations	66
B.7	InfiniBand™ Architecture protocol requirements	67

Figures

Figure 1 - SCSI document relationships	1
Figure 2 - RDMA communication service example	8
Figure 3 - Example RDMA channel establishment	9
Figure 4 - Memory descriptor mapping	18
Figure 5 - Example indirect data buffer descriptor with no PARTIAL MEMORY DESCRIPTOR LIST field.	22
Figure 6 - Example indirect data buffer descriptor with a PARTIAL MEMORY DESCRIPTOR LIST field.	23
Figure A.1 - SRP reference model	52
Figure A.2 - Model for a four-step confirmed service	53
Figure A.3 - Model for a two-step confirmed service.	53
Figure B.1 - InfiniBand™ Architecture device example	62
Figure B.2 - IB I/O unit example	62
Figure B.3 - SCSI architecture mapping	64

Working Draft

Tables

Table 1 - Memory descriptor	18
Table 2 - Data buffer descriptor formats	19
Table 3 - Supported data buffer descriptor formats	20
Table 4 - Indirect data buffer descriptor	21
Table 5 - SRP requests sent from SRP initiator ports to SRP target ports	24
Table 6 - SRP responses sent from SRP target ports to SRP initiator ports	24
Table 7 - SRP requests sent from SRP target ports to SRP initiator ports	24
Table 8 - SRP responses sent from SRP initiator ports to SRP target ports	24
Table 9 - SRP_LOGIN_REQ request	26
Table 10 - MULTI-CHANNEL ACTION code values	27
Table 11 - SRP_LOGIN_RSP response	28
Table 12 - MULTI-CHANNEL RESULT code values	29
Table 13 - SRP_LOGIN_REJ response	30
Table 14 - SRP_LOGIN_REJ response reason codes	31
Table 15 - SRP_I_LOGOUT request	32
Table 16 - SRP_T_LOGOUT request	33
Table 17 - SRP_T_LOGOUT request reason codes	34
Table 18 - SRP_TSK_MGMT request	35
Table 19 - TASK MANAGEMENT FUNCTION CODES	36
Table 20 - SRP_CMD request	37
Table 21 - TASK ATTRIBUTE	38
Table 22 - SRP_RSP response	40
Table 23 - RESPONSE DATA field	43
Table 24 - RSP_CODE values	43
Table 25 - SRP_CRED_REQ request	44
Table 26 - SRP_CRED_RSP response	45
Table 27 - SRP_AER_REQ request	46
Table 28 - SRP_AER_RSP response	48
Table 29 - SRP mode page codes	50
Table 30 - Disconnect-reconnect mode page	50

Foreword

This foreword is not part of American National Standard INCITS.***-200x.

Suggestions for improvement, requests for interpretation, addenda, or defect reports are welcome. They should be sent to the INCITS Secretariat, c/o Information Technology Industry Council , 1250 Eye Street, NW, Suite 200, Washington, DC 20005.

This standard was processed and approved for submittal to ANSI by the InterNational Committee for Information Technology Standards (INCITS). Committee approval of this standard does not necessarily imply that all committee members voted for approval. At the time it approved this standard, INCITS had the following members:

Karen Higginbottom, Chair

David Michael, Vice-chair

Monica Vago, Secretary

(INCITS Membership to be inserted)

INCITS technical committee T10 on Lower-Level Interfaces, which developed this standard, had the following members:

John B. Lohmeyer, Chair

George O. Penokie, Vice-Chair

Ralph Weber, Secretary

(T10 Membership to be inserted)

Introduction

The Small Computer System Interface (SCSI) command set is widely used and applicable to a wide variety of device types. The transmission of SCSI command set information across an RDMA communication service allows the large body of SCSI application and driver software to be successfully used on the InfiniBand^{TM1} Architecture, the VI Architecture and other interfaces that support RDMA communication service semantics.

The SCSI RDMA Protocol (SRP) standard is divided into the following clauses:

Clause 1 is the scope.

Clause 2 enumerates the normative references that apply to this standard.

Clause 3 describes the definitions, symbols, abbreviations, and conventions used in this standard.

Clause 4 describes the RDMA communication service model.

Clause 5 describes significant concepts of SRP.

Clause 6 describes the information units used by SRP.

Clause 7 defines the SCSI management features for SRP, including the SRP mode pages.

Annex A and Annex B form an integral part of this standard.

1. InfiniBand is a trademark and service mark of the InfiniBand Trade Association.

American National Standard for Information Systems – Information Technology – SCSI RDMA Protocol (SRP)

1 Scope

The SCSI family of standards provides for many different transport protocols that define the rules for exchanging information between different SCSI devices. This standard defines the rules for exchanging information between SCSI devices using an RDMA communication service. Other SCSI transport protocol standards define the rules for exchanging information between SCSI devices using other interconnects.

The set of SCSI standards specifies the interfaces, functions and operations necessary to ensure interoperability between conforming SCSI implementations. This standard is a functional description. Conforming implementations may employ any design technique that does not violate interoperability.

Figure 1 shows the relationship of SCSI transport protocol standards, such as this one, to the other standards and related projects in the SCSI family of standards as of the publication of this standard.

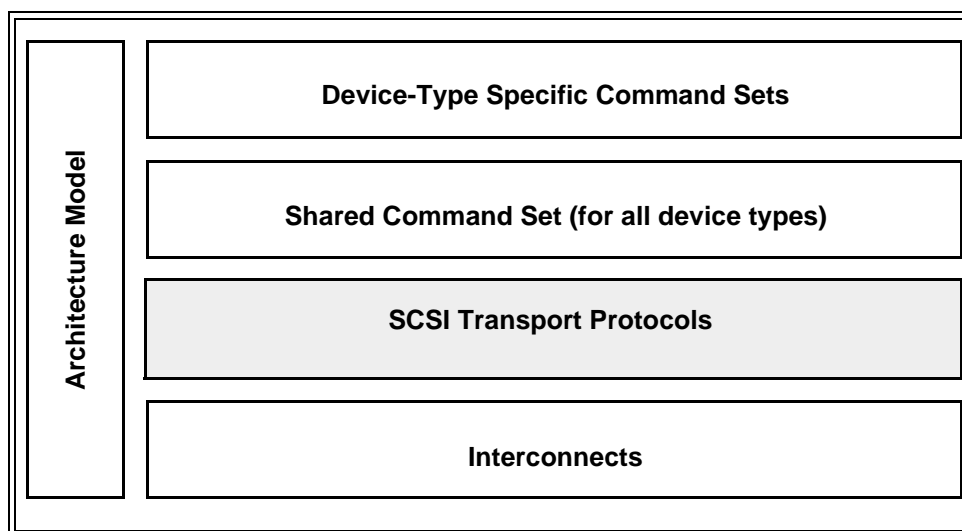


Figure 1 - SCSI document relationships

Figure 1 is intended to show the general relationship of the documents to one another. Figure 1 is not intended to imply a relationship such as a hierarchy, protocol stack or system architecture. It indicates the applicability of a standard to the implementation of a given transport.

At the time this standard was generated, examples of the SCSI general structure included:

Interconnects:

SCSI Parallel Interface - 4	SPI-4	[ANSI INCITS.362-200x]
Serial Storage Architecture Physical Layer 1	SSA-PH	[ANSI X3.293:1996]
Serial Storage Architecture Physical Layer 2	SSA-PH-2	[ANSI NCITS.307:1998]

SCSI Transport Protocols:

Serial Storage Architecture Transport Layer 1	SSA-TL-1	[ANSI X3.295:1996]
Serial Storage Architecture Transport Layer 2	SSA-TL-2	[ANSI NCITS.308:1998]
SCSI Parallel Interface - 2	SPI-2	[ANSI X3.302-1998]

SCSI Fibre Channel Protocol - 2	FCP-2	[T10/1144D]
SCSI Serial Bus Protocol - 2	SBP-2	[ANSI NCITS.325-1998]
Serial Storage Architecture SCSI-3 Protocol	SSA-S3P	[ANSI NCITS.309-1998]
Shared Command Sets:		
SCSI Primary Commands-3	SPC-3	[T10/1416D]
Device-Type Specific Command Sets:		
SCSI-3 Block Commands	SBC	[ANSI NCITS.306-1998]
SCSI-3 Enclosure Services	SES	[ANSI NCITS.305-1998]
SCSI-3 Stream Commands	SSC	[ANSI NCITS.335:2000]
SCSI-3 Medium Changer Commands	SMC	[ANSI NCITS.314:1998]
SCSI Multimedia Command Set - 2	MMC-2	[ANSI NCITS.333:2000]
SCSI Controller Commands - 2	SCC-2	[ANSI NCITS.318:1998]
Object-based Storage Devices Commands	OSD	[T10/1355-D]
Architecture Model:		
SCSI Architecture Model - 2	SAM-2	[T10/1157D]

The term SCSI is used to refer to the family of standards described in this clause.

2 Normative references

2.1 Normative references

The following standards contain provisions that, by reference in the text, constitute provisions of this standard. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this standard are encouraged to investigate the possibility of applying the most recent editions of the standards listed below.

Copies of the following documents may be obtained from ANSI: approved ANSI standards, approved and draft international and regional standards (ISO, IEC, CEN/CENELEC, ITUT), and approved and draft foreign standards (including BSI, JIS, and DIN). For further information, contact ANSI Customer Service Department at +1.212.642.4900 (telephone), +1.212.302.1286 (facsimile) or via the World Wide Web at <http://www.ansi.org>.

Additional availability contact information is provided below as needed.

2.2 Approved references

ISO/IEC 14776-312, SCSI Primary Commands - 2 (SPC-2) [ANSI NCITS.351:200x]

2.3 References under development

At the time of publication, the following referenced standards were still under development. For information on the current status of the document, or regarding availability, contact the relevant standards body or other organization as indicated.

ISO/IEC 14776-412, SCSI Architecture Model - 2 (SAM-2) [T10/1157-D]

ISO/IEC 14776-313, SCSI Primary Commands - 3 (SPC-3) [T10/1416-D]

NOTE 1 - For more information on the current status of a document, contact the INCITS Secretariat at +1.202.737.8888 (phone), +1.202.638.4922 (fax) or via Email at ncits@itic.org. To obtain copies of these documents, contact Global Engineering at 15 Inverness Way, East Englewood, CO 80112-5704 at +1.303.792.2181 (phone), 1.800.854.7179 (phone), or +1.303.792.2192 (fax), or at <http://global.ihs.com>.

Working Draft

3 Definitions, symbols, abbreviations and conventions

3.1 Definitions

3.1.1 acceptance data: Application protocol data communicated from a server consumer to the client consumer when a new RDMA channel is accepted (see 4.2). This protocol uses acceptance data to communicate the SRP_LOGIN_RSP response (see 6.3).

3.1.2 application client: An object that is the source of SCSI commands (see SAM-2).

3.1.3 byte: An 8-bit construct.

3.1.4 channel attributes: Information provided during RDMA channel establishment that identifies the type and characteristics of the desired RDMA channel (see 4.2). The format and interpretation of channel attributes are specific to a particular RDMA communication service.

3.1.5 command: A request describing a unit of work to be performed by a device server (see SAM-2).

3.1.6 command descriptor block (CDB): The structure used to communicate commands from an application client to a device server (see SPC-2).

3.1.7 consumer: An object that communicates with other consumers using an RDMA communication service (see 4.1). In this protocol, a consumer is either an SRP target port or an SRP initiator port.

3.1.8 data-in buffer: The buffer identified by the application client to receive data from the device server during the execution of a command (see SAM-2).

3.1.9 data-out buffer: The buffer identified by the application client to supply data that is sent from the application client to the device server during the execution of a command (see SAM-2).

3.1.10 device server: An object within a logical unit that executes SCSI tasks according to the rules of task management (see SAM-2).

3.1.11 information unit: An organized collection of data specified by this protocol to be transferred as login data, rejection data, acceptance data, or a message on an RDMA channel.

3.1.12 initiator port identifier: A value by which a SCSI initiator port is referenced within a domain (see SAM-2).

3.1.13 logical unit: A target-resident object that implements a device model and processes SCSI commands sent by an application client.

3.1.14 logical unit number (LUN): A 64-bit identifier for a logical unit (see SAM-2).

3.1.15 login data: Data communicated from a client consumer to a server agent or server consumer during RDMA channel establishment (see 4.2) that is meaningful within the client/server application protocol. This protocol uses login data to communicate the SRP_LOGIN_REQ request (see 6.2).

3.1.16 message: A communication sent by one consumer to another using an RDMA channel (see 4.3).

3.1.17 RDMA channel: A communication path between two consumers of an RDMA communication service (see 4.1).

3.1.18 RDMA communication service: The software, protocols, and interconnect that provides message and RDMA operations between pairs of consumers (see clause 4).

3.1.19 RDMA operation: Either an RDMA Read operation or an RDMA Write operation.

3.1.20 RDMA Read operation: An operation by which a requesting consumer may fetch data from memory registered by the other consumer associated with an RDMA channel (see 4.4).

3.1.21 RDMA Write operation: An operation by which a requesting consumer may store data into memory registered by the other consumer associated with an RDMA channel (see 4.4).

Working Draft

3.1.22 rejection data: Application protocol data communicated from a server agent or server consumer to the client consumer when a new RDMA channel is rejected (see 4.2). This protocol uses rejection data to communicate the SRP_LOGIN_REJ response (see 6.4).

3.1.23 sense data: Data returned to an application client in the SENSE DATA field of an SRP_RSP response or an SRP_AER_REQ request. See SAM-2.

3.1.24 server agent: An entity that provides services (e.g., connection management) on behalf of a server consumer.

3.1.25 server identifier: Information provided to an RDMA communication service by a client consumer that allows the RDMA communication service to locate the desired server consumer. The format and interpretation of a server identifier are specific to the RDMA communication service.

3.1.26 SRP initiator port: A SCSI initiator port that uses this protocol to communicate with an SRP target port.

3.1.27 SRP initiator port identifier: A value by which an SRP initiator port is identified to an SRP target port.

3.1.28 SRP target port: A SCSI target port that uses this protocol to communicate with an SRP initiator port.

3.1.29 SRP target port identifier: A value by which an SRP target port is identified within an SRP domain.

3.1.30 status: Response information sent from a device server to an application client upon completion of each command (see SAM-2).

3.1.31 target port identifier: A value by which a SCSI target port is referenced within a domain (see SAM-2).

3.2 Acronyms

CDB	Command Descriptor Block (see 3.1.6)
INCITS	InterNational Committee for Information Technology Standards
LSB	Least significant bit
LUN	Logical Unit Number (see 3.1.14)
MSB	Most significant bit
NCITS	National Committee for Information Technology Standards (now INCITS)
RDMA	Remote Direct Memory Access
SAM-2	SCSI Architecture Model - 2 (see 2.3)
SCSI	The architecture defined by the family of standards described in clause 1
SPC-2	SCSI Primary Commands - 2 (see 2.3)
SRP	SCSI RDMA Protocol (this standard)

3.3 Keywords

3.3.1 expected: A keyword used to describe the behavior of the hardware or software in the design models assumed by this standard. Other hardware and software design models may also be implemented.

3.3.2 ignored: A keyword used to describe an unused bit, byte, word, field or code value. The contents or value of an ignored bit, byte, word, field or code value shall not be examined by the receiving SCSI device and may be set to any value by the transmitting SCSI device.

3.3.3 invalid: A keyword used to describe an illegal or unsupported bit, byte, word, field or code value. Receipt of an invalid bit, byte, word, field or code value shall be reported as an error.

Working Draft

3.3.4 mandatory: A keyword indicating an item that is required to be implemented as defined in this standard.

3.3.5 may: A keyword that indicates flexibility of choice with no implied preference (equivalent to "may or may not").

3.3.6 may not: Keywords that indicate flexibility of choice with no implied preference (equivalent to "may or may not").

3.3.7 obsolete: A keyword indicating that an item was defined in prior SCSI standards but has been removed from this standard.

3.3.8 optional: A keyword that describes features that are not required to be implemented by this standard. However, if any optional feature defined by this standard is implemented, then it shall be implemented as defined in this standard.

3.3.9 reserved: A keyword referring to bits, bytes, words, fields and code values that are set aside for future standardization. A reserved bit, byte, word or field shall be set to zero, or in accordance with a future extension to this standard. Recipients are not required to check reserved bits, bytes, words or fields for zero values. Receipt of reserved code values in defined fields shall be reported as an error.

3.3.10 restricted: A keyword referring to bits, bytes, words, and fields that are set aside for use in other SCSI standards. A restricted bit, byte, word, or field shall be treated as a reserved bit, byte, word or field for the purposes of the requirements defined in this standard.

3.3.11 shall: A keyword indicating a mandatory requirement. Designers are required to implement all such mandatory requirements to ensure interoperability with other products that conform to this standard.

3.3.12 should: A keyword indicating flexibility of choice with a strongly preferred alternative; equivalent to the phrase "it is strongly recommended".

3.4 Conventions

Certain words and terms used in this standard have a specific meaning beyond the normal English meaning. These words and terms are defined either in 3.1 or in the text where they first appear.

Names of commands, statuses, sense keys, additional sense codes and additional sense code qualifiers are in all uppercase (e.g., REQUEST SENSE).

Names of fields and state variables are in small uppercase (e.g. ALLOCATION LENGTH). When a field or state variable name contains acronyms, uppercase letters may be used for readability (e.g. NORMACA). Normal case is used when the contents of a field or state variable are being discussed. Fields or state variables containing only one bit are usually referred to as the NAME bit instead of the NAME field.

Normal case is used for words having the normal English meaning.

Numbers that are not immediately followed by lower-case b or h are decimal values.

Numbers immediately followed by lower-case b (e.g. 0101b) are binary values.

Numbers or upper case letters immediately followed by lower-case h (e.g. FA23h) are hexadecimal values.

Lists sequenced by letters (e.g., a-red, b-blue, c-green) show no ordering relationship between the listed items. Numbered lists (e.g., 1-red, 2-blue, 3-green) show an ordering between the listed items.

If a conflict arises between text, tables or figures, the order of precedence to resolve the conflicts is text; then tables; and finally figures. Not all tables or figures are fully described in the text. Tables show data format and values.

Notes do not constitute any requirements for implementors.

Working Draft

3.5 Notation for procedures and functions

In this standard, the model for functional interfaces between objects is the callable procedure. Such interfaces are specified using the following notation:

[Result =] Procedure Name (IN ([input-1] [,input-2] ...), OUT ([output-1] [,output-2] ...))

Where:

Result: A single value representing the outcome of the procedure or function.

Procedure Name: A descriptive name for the function to be performed.

Input-1, Input-2, ...: A comma-separated list of names identifying caller-supplied input data objects.

Output-1, Output-2, ...: A comma-separated list of names identifying output data objects to be returned by the procedure.

"[...]": Brackets enclosing optional or conditional parameters and arguments.

This notation allows data objects to be specified as inputs and outputs. The following is an example of a procedure specification:

Found = Search (IN (Pattern, Item List), OUT ([Item Found]))

Where:

Found = Flag

Flag, which, if set, indicates that a matching item was located.

Input Arguments:

Pattern = ... /* Definition of Pattern object */

Object containing the search pattern.

Item List = Item<NN> /* Definition of Item List as an array of NN Item objects*/

Contains the items to be searched for a match.

Output Arguments:

Item Found = Item ... /* Item located by the search procedure */

This object is only returned if the search succeeds.

Working Draft

4 RDMA communication service model

4.1 Overview

This protocol is designed to operate using an RDMA communication service. An RDMA communication service provides communication between pairs of consumers using messages for control information and RDMA operations for data transfers. This clause describes an abstract RDMA communication service suitable for supporting this protocol. Annex B describes the mapping of these functions to those provided by the InfiniBand™¹ Architecture.

Figure 2 shows an example system that uses an RDMA communication service. Communication is provided by RDMA channels. An RDMA channel provides communication between two consumers. A single pair of consumers may communicate using many RDMA channels if sufficient resources are available. Some environments may use multiple special purpose RDMA channels between a single pair of consumers (e.g., a pair of consumers may use certain RDMA channels for messages and other RDMA channels for RDMA operations).

The RDMA communication service in figure 2 is comprised of adapters and other unspecified components (e.g. wires, fabric switches). The components of an RDMA communication service are implementation-specific.

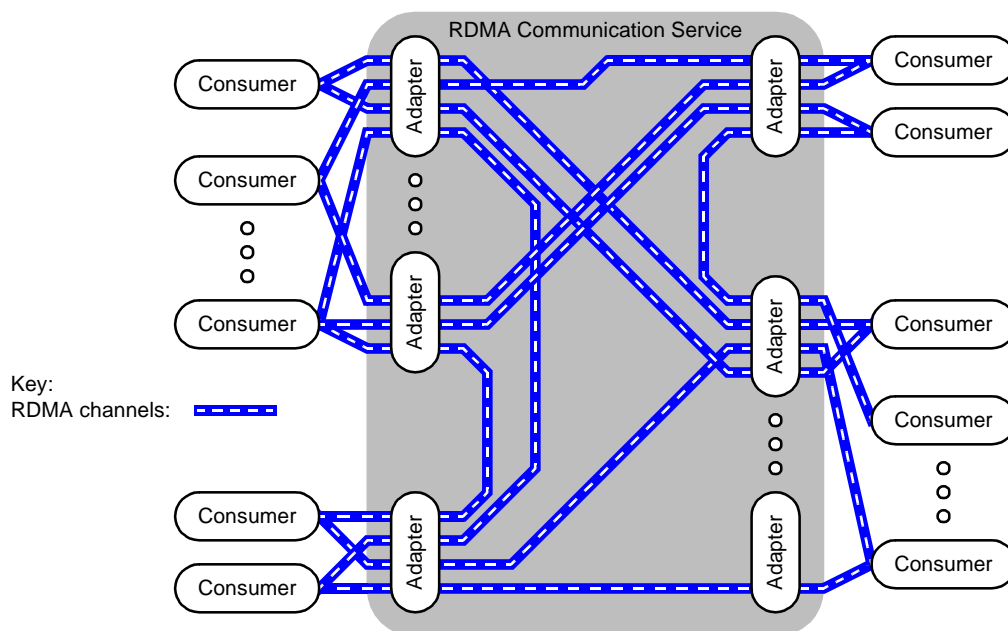


Figure 2 - RDMA communication service example

4.2 RDMA Channels

4.2.1 Introduction

An RDMA channel provides communication between a pair of consumers using messages, RDMA operations, or both. An RDMA channel is a dynamic association, established and destroyed upon request. Establishing an RDMA channel may require obtaining resources to support the RDMA channel, either within the RDMA channel's consumers or within the RDMA communication service or both. The resources associated with an RDMA channel may be released after the RDMA channel is disconnected.

1. InfiniBand is a trademark and service mark of the InfiniBand Trade Association.

4.2.2 Establishment

Figure 3 shows an example of the process by which an RDMA channel is established. A client consumer requests that the RDMA communication service establish an RDMA channel. The request is directed to a server agent and, if successful, resolved to a server consumer. The resulting RDMA channel provides communication between the client consumer and the server consumer.

A client consumer provides a server identifier to establish an RDMA channel. The format and interpretation of a server identifier are specific to the RDMA communication service. A server identifier may specify an individual server consumer or multiple server consumers (e.g., a server identifier may identify an adapter as shown in figure 2, specifying all consumers that implement a specific application protocol and are accessible through that adapter).

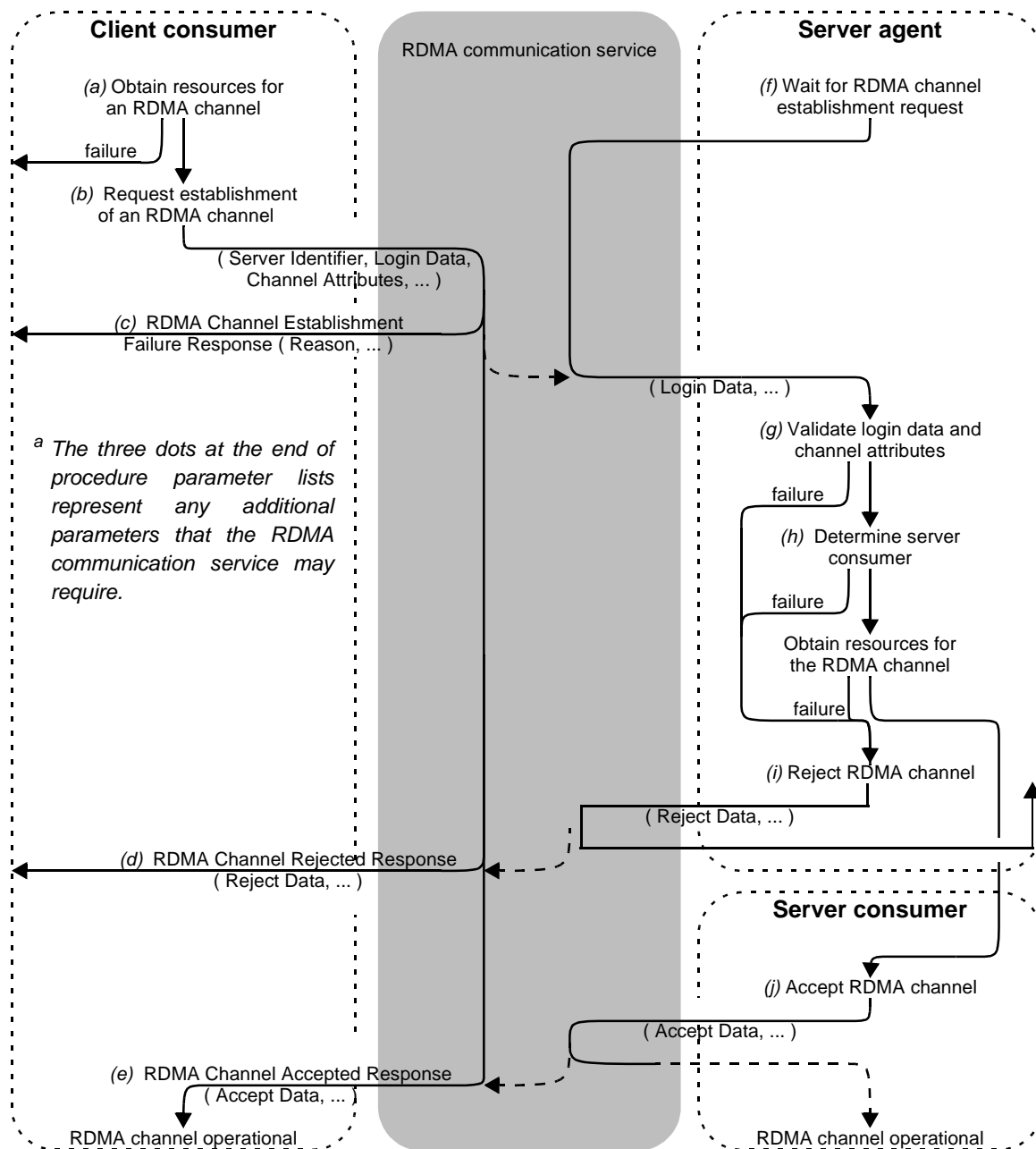


Figure 3 - Example RDMA channel establishment

Working Draft

In the example shown in figure 3, the recipient of an RDMA channel establishment request, identified by a server identifier, is called a server agent. The server agent determines whether an RDMA channel establishment request may be accepted and the server consumer to which it shall be assigned. A server agent may not be a distinct object. Some or all of the actions that figure 3 shows being performed by a server agent may be performed by a server consumer or by the RDMA communication service.

An RDMA communication service may require that the client consumer obtain resources (*Figure 3, a*) before requesting that an RDMA channel be established. After obtaining those resources, the client consumer may request (*b*) that the RDMA communication service establish an RDMA channel. The request includes a server identifier, login data, channel attributes, and any other parameters required by the RDMA communication service. Under this protocol, the client consumer is an SRP initiator port, the server identifier is an interconnect-specific value that enables the RDMA communication service to locate one or more SRP target ports, and the login data contains an SRP_LOGIN_REQ request (see 6.2).

The RDMA communication service returns one of three responses to the client consumer for an RDMA channel establishment request:

- a) An RDMA Channel Establishment failure response (*c*);
- b) An RDMA Channel Rejected response (*d*); or
- c) An RDMA Channel Accepted response (*e*).

An RDMA Channel Establishment Failure response (*c*) indicates that the RDMA channel was not established for some reason internal to the RDMA communication service. An RDMA Channel Establishment Failure response may return an RDMA communication service-specific reason code to identify the cause of the failure as well as other RDMA communication service-specific data.

An RDMA Channel Rejected response (*d*) indicates that the request was rejected by the server agent or server consumer. An RDMA Channel Rejected response may return rejection data provided by the server agent or server consumer. Rejection data may include a reason for rejecting the request. In this protocol, the rejection data includes an SRP_LOGIN_REJ response (see 6.4). An RDMA Channel Rejected response may also return RDMA communication service-specific data.

An RDMA Channel Accepted response (*e*) indicates that the RDMA channel has been successfully established. The client consumer may use the RDMA channel in accordance with the application protocol. An RDMA Channel Accepted response returns acceptance data provided by the server agent or server consumer. In this protocol, the acceptance data includes an SRP_LOGIN_RSP response (see 6.3). An RDMA Channel Accepted Response may also return data specific to the RDMA communication service.

An RDMA communication service may require that a server agent register (*f*) itself prior to receiving connection establishment requests. In figure 3 this is shown as a registration request (e.g., subroutine call) that returns control to the server agent when an RDMA Channel Establishment request is received. The way that a server agent registers with an RDMA communication service is specific to that service or the server.

RDMA Channel Establishment requests that are acceptable to the RDMA communication service are passed (*g*) to the server agent. The server agent receives the login data and RDMA communication service-specific data from the client consumer's request. In this protocol, the login data includes an SRP_LOGIN_REQ request (see 6.2).

The server agent determines whether the RDMA Channel Establishment request may be accepted and determines (*i*) the server consumer to be associated with the RDMA channel. If the request is not accepted the server agent or server consumer instructs the RDMA communication service to reject (*j*) the RDMA channel. The server agent or server consumer provides rejection data and any RDMA communication service-specific data that is required. In this protocol, the rejection data shall contain an SRP_LOGIN_REJ response (see 6.4).

If the RDMA Channel Establishment request is accepted, the server agent or server consumer instructs the RDMA communication service to accept (*k*) the RDMA channel. The server agent or server consumer provides

acceptance data and any RDMA communication service-specific data that is required. In this protocol, the acceptance data shall contain an SRP_LOGIN_RSP response (see 6.3).

4.2.3 Disestablishment

An RDMA channel may be disconnected by the RDMA communication service due to an error or by request from either of the RDMA channel's consumers. The consumers may each be notified that the RDMA channel has been disconnected, allowing the consumers to recover any resources associated with the RDMA channel. The time to deliver a notification may vary depending upon the RDMA communication service, the consumer being notified, and the specific circumstances of the disconnection request.

A disconnection request or error causes an RDMA channel to become non-operational. Operations in progress on an RDMA channel at the time it becomes non-operational and operations requested subsequently may not complete, and the status of those operations may be indeterminate.

4.3 Messages

A message is sent by one consumer associated with an RDMA channel (the sending consumer) to the other consumer associated with the RDMA channel (the receiving consumer). A message contains a payload of some number of data bytes. An RDMA communication service may provide a facility known as solicited message reception notification, by which some messages may be marked by the sender to indicate to the recipient that the marked message is more urgent than a message that is not marked. How such a marked message is handled by the recipient is outside the scope of this standard.

A sending consumer requests that a message be sent by providing the following to an RDMA communication service:

- a) the message's payload length;
- b) the message's payload data;
- c) the RDMA channel to use; and
- d) whether to use normal or solicited message reception notification.

The RDMA communication service attempts to deliver the message to the receiving consumer. If delivery succeeds, the RDMA communication service notifies the receiving consumer that a message has been received, providing the message's length, payload, and the RDMA channel on which the message was received. The RDMA communication service may also provide an indication of whether the sending consumer specified normal or solicited message reception notification.

Sending a message on an RDMA channel when no receive buffer has been provided, or when the provided receive buffer is too small for the message, may result in behavior that is not specified by this standard.

NOTE 2 - Such behavior may include (but is not limited to) disconnecting the RDMA channel, discarding or truncating the message, or delaying delivery of the message until a suitable message receive buffer becomes available. The RDMA communication service may not provide an error indication.

An RDMA communication service may not provide a way for a sending consumer to determine whether a message has been delivered to the receiving consumer (e.g., an acknowledgement may indicate only that a message was received without error).

4.4 RDMA operations

4.4.1 Overview

An RDMA channel may provide RDMA Write operations, RDMA Read operations, or both between its consumers.

A consumer may allow RDMA access by registering some or all of its memory with an RDMA communication service. The RDMA communication service returns a memory handle to identify the registered memory. The consumer may specify that the memory handle is usable for memory access on only a specified RDMA channel

or on a group of RDMA channels. The consumer may impose other access restrictions allowed by the RDMA communication service (e.g. read-only access).

A consumer that has registered memory and obtained a memory handle may communicate the memory handle to another consumer. This may be done using an application protocol contained in message payloads. The other consumer may then use the memory handle to request RDMA operations that access the memory registered by the first consumer.

The registered memory identified by a memory handle is represented as a memory address space. Accessible locations are identified by addresses. An RDMA communication service is not required to provide a way to determine, from a message handle, which memory locations are accessible, the number of locations that are accessible, or the type of access allowed.

4.4.2 RDMA Write

An RDMA Write operation allows a requesting consumer to store data into memory registered by another consumer. A requesting consumer provides the following to an RDMA communication service when it requests an RDMA Write operation:

- a) An RDMA channel to use for the operation;
- b) A memory handle that is usable for access on that RDMA channel;
- c) A range of addresses within the memory address space identified by the memory handle; and
- d) Data to be written into the specified range of addresses.

An RDMA communication service is not required to provide a way for a requesting consumer to determine whether the data has been written into the specified range of addresses in registered memory (e.g., an acknowledgment may only signify error-free reception of the data). An RDMA communication service is not required to provide a way for the consumer that registered the memory to determine whether an RDMA Write operation is in progress or has completed.

4.4.3 RDMA Read

An RDMA Read operation allows a requesting consumer to fetch data from memory registered by another consumer. A requesting consumer provides the following to an RDMA communication service when it requests an RDMA Read operation:

- a) An RDMA channel to use for the operation;
- b) A memory handle that is usable for access on that RDMA channel;
- c) A range of addresses within the memory address space identified by the memory handle; and
- d) A buffer into which to place the data read from the specified range of addresses.

The RDMA communication service notifies the requesting consumer after data has been successfully obtained from the specified range of addresses and placed in the requestor's buffer. An RDMA communication service is not required to provide a way for the consumer that registered the memory to determine whether an RDMA Read operation is in progress or has completed.

4.5 Ordering and Reliability

4.5.1 Ordering and reliability overview

This protocol operates using an RDMA communication service having the characteristics described in 4.5.2 and 4.5.3. Use of this protocol with an RDMA communication service having different characteristics is outside the scope of this standard.

4.5.2 Reliability

An RDMA communication service shall deliver each message sent on an RDMA channel to the receiving consumer or destroy the RDMA channel. Each delivered message shall be delivered to the receiving consumer once, without duplication; the RDMA communication service shall discard any duplicates that may result from

Working Draft

retransmission or other mechanisms. Each delivered message shall be delivered to the receiving consumer complete and error-free.

The RDMA communication service shall provide to the sending consumer an indication of the completion status of each RDMA communication service request. This status shall be one of:

- a) successful - The request completed without error.
- b) error - The request was not completed due to an error. The RDMA communication service may provide additional information about the error. This status should be returned immediately when the RDMA channel does not exist or has experienced an error.
- c) timeout - No indication was received, completion status of request is unknown, RDMA communication service has experienced an error. The length of time after which a timeout indication is returned is specific to the RDMA communication service.

4.5.3 Ordering

Messages sent on an RDMA channel shall be delivered to the receiving consumer in the order they were sent. The data for all RDMA Write operations requested on an RDMA channel by a consumer prior to that same consumer sending a message on the same RDMA channel shall be available to the receiving consumer (e.g. stored into registered memory) before the message is delivered to the receiving consumer. The order in which multiple RDMA Writes (without an intervening message) are applied to a particular memory location is specific to the RDMA communication service.

Messages sent on different RDMA channels may be delivered in any order. The data for RDMA Write operations may be stored into registered memory in any order relative to the delivery of messages sent on other RDMA channels. RDMA Write operations requested on different RDMA channels may store data into the same registered memory location in any order.

RDMA Read operations may be processed in any order.

If an RDMA communication service fails to meet the ordering requirements of this subclause on an RDMA channel, it shall destroy the RDMA channel.

5 Structure and concepts

5.1 Overview of SRP operation

5.1.1 RDMA channel establishment and login

SRP initiator ports login with SRP target ports when a new RDMA channel is established for use with this protocol. The login process associates an RDMA channel with a specific SRP initiator port and SRP target port (i.e., an I_T nexus (see SAM-2)) and negotiates parameters that govern the use of that RDMA channel.

SRP initiator ports and SRP target ports shall be determined by their role during RDMA channel establishment. An object that requests RDMA channel establishment as a client consumer (see 4.2) shall be an SRP initiator port. An object that accepts RDMA channel establishment as a server consumer (see 4.2) shall be an SRP target port.

Login occurs during RDMA channel establishment. An SRP initiator port shall provide an SRP_LOGIN_REQ request (see 6.2) as the login data when establishing a new RDMA channel. If an SRP target port accepts a new RDMA channel it shall provide an SRP_LOGIN_RSP response (see 6.3) as the acceptance data. If an SRP target port does not accept a new RDMA channel it shall provide an SRP_LOGIN_REJ response (see 6.4) as the rejection data when rejecting the new RDMA channel.

The SRP_LOGIN_REQ request (see 6.2) contains an SRP initiator port identifier and an SRP target port identifier. An SRP target port shall not accept a new RDMA channel unless its SRP target port identifier is identical to the value in the SRP_LOGIN_REQ request. If an SRP target port accepts a new RDMA channel, it shall treat all communication on that RDMA channel as being with the SRP initiator port identified by the SRP initiator port identifier specified in the SRP_LOGIN_REQ request.

5.1.2 Single RDMA channel operation

An SRP initiator port may specify single RDMA channel operation during login. If an SRP target port accepts such a login, it shall:

- a) Attempt to send an SRP_T_LOGOUT request (see 6.6) on any established RDMA channel that specified the same SRP initiator port identifier. The reason code shall indicate that the RDMA channel was disconnected due to a MULTI-CHANNEL ACTION code in a new SRP_LOGIN_REQ request (see 6.2);
- b) Request disconnection of any established RDMA channel (see 5.1.4) that specified the same SRP initiator port identifier; and
- c) Reject any other RDMA channel establishment requests it has received that specified the same SRP initiator port identifier and that the SRP target port has not yet accepted.

Following acceptance of a login specifying single RDMA channel operation, that single RDMA channel shall be used for all communication between the specified SRP initiator port and SRP target port. Subsequent logins specifying other modes of operation may allow communication using multiple RDMA channels.

When an I_T nexus is using this protocol in the single RDMA channel mode, these events are I_T nexus loss notification events (see SAM-2, SPC-3):

- a) Sending or receiving an SRP_I_LOGOUT request or an SRP_T_LOGOUT request;
- b) Requesting that an RDMA channel be disconnected; or
- c) Receiving notification that an RDMA channel has been disconnected.

5.1.3 Multiple independent RDMA channel operation

An SRP initiator port may specify multiple independent RDMA channel operation during login. An SRP target port shall not accept such a login if doing so would require disconnecting an established RDMA channel with the same SRP initiator port, and shall return the SRP_T_LOGOUT request reason code RDMA CHANNEL LIMIT REACHED FOR THIS INITIATOR.

Working Draft

Following acceptance of a login specifying multiple independent RDMA channel operation, one or more RDMA channels may be used for communication between the SRP initiator port and the SRP target port. All such RDMA channels are associated with the single I_T nexus defined by the SRP initiator port identifier and the SRP target port identifier.

When multiple independent RDMA channels are used, operation of each SRP request is confined to a single RDMA channel. The sender of an SRP request chooses an RDMA channel to use for sending the SRP request. The sender of an SRP response shall use the same RDMA channel as the SRP request for sending the SRP response. All RDMA operations associated with the SRP request shall also use the same RDMA channel as the SRP request.

While each SRP request is confined to a single RDMA channel, SCSI tasks and task management functions may be conveyed on independent RDMA channels associated with the same I_T nexus. SCSI tasks and task management functions interact as specified by SAM-2, SPC-2 and other SCSI command standards (e.g., within an I_T nexus, a SCSI task sent on one RDMA channel may be aborted by an ABORT TASK sent on a different RDMA channel.)

An RDMA communication service may not provide any ordering relationship between SRP requests, SRP responses and RDMA operations that use different RDMA channels. If ordering is important for a sequence of SRP requests, they should be sent using the same RDMA channel.

When an I_T nexus is using this protocol in the multiple independent RDMA channel mode, these events are I_T nexus loss notification events (see SAM-2, SPC-3) when they occur with respect to the last (or only) channel associated with the I_T nexus:

- a) Sending or receiving an SRP_I_LOGOUT request or an SRP_T_LOGOUT request;
- b) Requesting that an RDMA channel be disconnected; or
- c) Receiving notification that an RDMA channel has been disconnected.

5.1.4 RDMA channel disconnection

RDMA channel disconnection may cause (see 5.1.2 and 5.1.3) an I_T nexus loss notification event as described in SAM-2 and SPC-3.

An SRP initiator port should send an SRP_I_LOGOUT request (see 6.5) and wait for the RDMA communication service status indication (see 4.5.2) before requesting that an RDMA channel be disconnected.

After requesting that an RDMA channel be disconnected, after being notified that an RDMA channel has been disconnected, or upon receiving an SRP_T_LOGOUT request (see 6.6), an SRP initiator port shall:

- a) Discard any outstanding request received from an SRP target port on that RDMA channel, without returning a response;
- b) Not send any further messages on that RDMA channel;
- c) Discard any subsequent messages received on that RDMA channel; and
- d) For any outstanding SCSI tasks sent on that RDMA channel, indicate to the application client that the task has terminated with a service delivery system failure.

An SRP target port should send an SRP_T_LOGOUT request (see 6.6) and wait for the RDMA communication service status indication (see 4.5.2) before requesting that an RDMA channel be disconnected.

After requesting that an RDMA channel be disconnected, after being notified that an RDMA channel has been disconnected, or upon receiving an SRP_I_LOGOUT request (see 6.5), an SRP target port shall:

- a) Abort all outstanding SCSI tasks that were contained in SRP_CMD requests (see 6.8) received on that RDMA channel, without returning a response;
- b) Discard any other outstanding requests received from an SRP initiator port on that RDMA channel, without returning a response;
- c) Not send any further messages on that RDMA channel;

Working Draft

- d) Discard any subsequent messages received on that RDMA channel; and
- e) Not alter previously established conditions, including MODE SELECT parameters, reservations, ACA, and CA as a result of the disconnection.

5.2 Identifiers

Initiator port identifiers and target port identifiers (see SAM-2) for this protocol are 16 bytes in length.

5.3 Alias associations

There are no events specific to this protocol that clear alias associations (see SPC-2).

5.4 Information unit classes

Each SRP information unit is classified as an SRP request (see table 5 and table 7) or an SRP response (see table 6 and table 8). SRP requests convey SCSI commands, task management requests and RDMA channel management requests. SRP responses convey SCSI command and task management service responses and RDMA channel management responses. RDMA channel management requests may be issued by SRP target ports or SRP initiator ports.

In normal operation, SRP requests and SRP responses occur in pairs. Each SRP request elicits a single corresponding SRP response from the SRP device receiving the SRP request. An SRP request communicates the start of a remote procedure call; the corresponding SRP response communicates the remote procedure call's completion.

An SRP response shall not be returned:

- a) for an SRP_CMD request if the associated task is aborted;
- b) for an SRP_T_LOGOUT request (see 6.6);
- c) for an SRP_I_LOGOUT request (see 6.5); and
- d) for outstanding SRP requests received on an RDMA channel when an SRP device becomes aware of a failure preventing further communication on that RDMA channel. In this case, the device shall abort all outstanding SRP requests received on that RDMA channel.

In all other cases an SRP device shall return a single SRP response for each SRP request it receives.

SRP responses shall be sent on the RDMA channel on which the corresponding SRP request was received.

5.5 SRP target port buffer management

5.5.1 Buffer management overview

SRP target port buffer management allows an SRP target device to limit the number of SRP requests that may be sent on an RDMA channel. SRP devices may use SRP target port buffer management to manage internal and RDMA channel-related resources.

SRP responses are not subject to buffer management; they may be sent at any time. An SRP device may limit the number of SRP responses it may receive by limiting the number of SRP requests it has outstanding.

5.5.2 SRP requests issued by target port

SRP target ports shall limit themselves to one outstanding SRP request (see table 7) per RDMA channel. Upon sending an SRP request, an SRP target port shall not send another SRP request on the same RDMA channel until after it receives the SRP response (see table 8) for the previous SRP request.

5.5.3 Requests issued by initiator port

This protocol uses a credit-based buffer management algorithm to limit the number of SRP requests (see table 5) that an SRP initiator port may send to an SRP target port. The algorithm uses a field, REQUEST LIMIT DELTA, that is present in most information units sent by an SRP target port to an SRP initiator port (not in SRP_LOGIN_REJ or SRP_T_LOGOUT), and a state variable, REQUEST LIMIT.

Working Draft

Most information units containing a REQUEST LIMIT DELTA field do not generate a confirmation that the SRP initiator port has received the information unit and processed the contents of the request limit delta field. The SRP_CRED_REQ (see 6.10) and SRP_AER_REQ (see 6.12) requests do generate a confirmation through the SRP_CRED_RSP (see 6.11) and SRP_AER_RSP (see 6.13) responses respectively. An SRP initiator port shall process the REQUEST LIMIT DELTA fields of information units received on an RDMA channel in the order that they are received. An SRP initiator port shall process the REQUEST LIMIT DELTA field of a request before sending that request's response. (e.g., An SRP initiator port shall process the REQUEST LIMIT DELTA field of an SRP_CRED_REQ request before sending the SRP_CRED_RSP.) The following rules specify the buffer management algorithm for SRP requests sent by SRP initiator ports:

- a) The Request Limit and Request Limit Delta variables are both signed, two's complement 32-bit integers. SRP initiator ports shall implement a separate copy of the request limit variable for each RDMA channel;
- b) Upon successful completion of RDMA channel establishment an SRP initiator port shall initialize the RDMA channel's Request Limit variable to the value of the REQUEST LIMIT DELTA field received in the SRP_LOGIN_RSP response (see 6.3). Except for providing an SRP_LOGIN_REQ request (see 6.2) when requesting RDMA channel establishment, the SRP initiator port shall not send any SRP information units on the RDMA channel prior to initializing the Request Limit variable;
- c) An SRP initiator port may send an SRP request on an RDMA channel when the value of the RDMA channel's Request Limit variable is greater than zero. An SRP initiator port shall not send an SRP request on any RDMA channel whose Request Limit variable has a value less than or equal to zero; the results of doing so are vendor-specific. To ensure that task management requests may be sent, an SRP initiator port may choose to send commands only when the value of the Request Limit variable is greater than one;
- d) An SRP initiator port shall decrement an RDMA channel's Request Limit variable by one whenever it sends an SRP request on that RDMA channel;
- e) An SRP initiator port shall add (two's complement addition) the value of the REQUEST LIMIT DELTA field to an RDMA channel's Request Limit variable whenever it receives an information unit on that RDMA channel; and
- f) An SRP target port shall not specify a positive value of the REQUEST LIMIT DELTA field that may cause the SRP initiator port's Request Limit variable to exceed 2^{30} ; and
- g) An SRP target port shall not specify a negative value for the REQUEST LIMIT DELTA field in an information unit that may cause the SRP initiator port's Request Limit variable to drop below -2^{31} .

5.6 Data buffers

5.6.1 Memory descriptors

A memory descriptor is a 16-byte structure that identifies a memory segment (see table 1). Figure 4 illustrates the mapping of a memory descriptor to a memory segment.

Table 1 - Memory descriptor

Bit Byte	7	6	5	4	3	2	1	0
0	(MSB)							
...	VIRTUAL ADDRESS							
7								
	(LSB)							
8	(MSB)							
...	MEMORY HANDLE							
11								
	(LSB)							
12	(MSB)							
...	DATA LENGTH							
15								
	(LSB)							

The VIRTUAL ADDRESS field contains an unsigned integer value that identifies the byte address within the memory region of the first byte of the memory segment.

The MEMORY HANDLE field contains an SRP initiator port-specific value that identifies the region that contains the memory segment. The SRP target port shall supply this value with any RDMA operation that accesses the memory segment.

The DATA LENGTH field contains an unsigned integer value that identifies the length of the memory segment in bytes. The interpretation of a memory descriptor where the sum of the VIRTUAL ADDRESS and DATA LENGTH fields exceeds 2^{64} is vendor-specific.

An SRP target port may use a memory descriptor for either RDMA Read operations or RDMA Write operations but not both. SRP target ports shall issue only the appropriate type of RDMA operation for a memory descriptor,

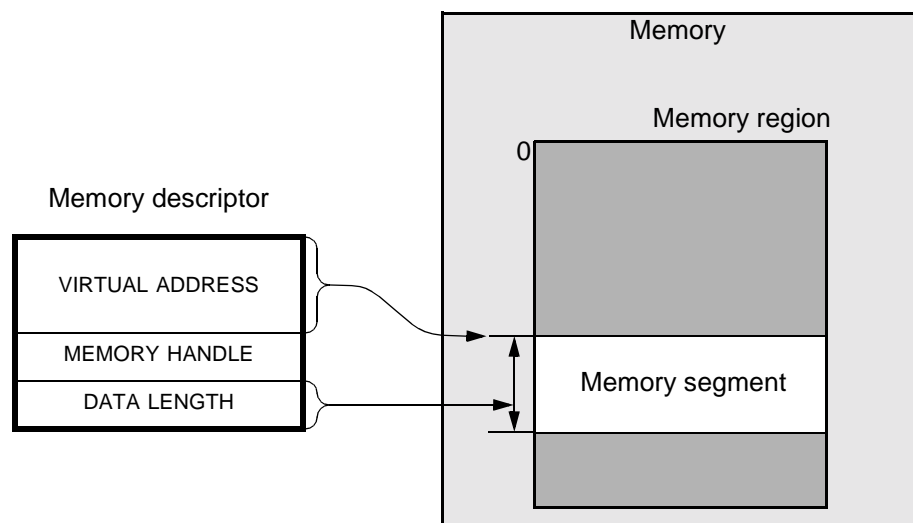


Figure 4 - Memory descriptor mapping

Working Draft

depending on whether the descriptor was a data-in or data-out descriptor, and shall ensure that each RDMA operation is wholly contained within the memory segment by using the following rules:

- a) The RDMA operation's starting virtual address shall be greater than or equal to the memory descriptor's VIRTUAL ADDRESS and less than the sum of the memory descriptor's VIRTUAL ADDRESS and DATA LENGTH; and
- b) The sum of the RDMA operation's VIRTUAL ADDRESS and DATA LENGTH shall be greater than the memory descriptor's VIRTUAL ADDRESS and less than or equal to the sum of the memory descriptor's VIRTUAL ADDRESS and DATA LENGTH.

5.6.2 Data buffer descriptors

5.6.2.1 Overview

An SRP_CMD request (see 6.8) may contain a data-out buffer descriptor, a data-in buffer descriptor, both, or neither, depending upon the data transfer(s) requested by the SCSI command. The format of each data buffer descriptor is specified by a format code value. In an SRP_CMD request with both data-in and data-out buffer descriptors, there is no requirement that both buffer descriptors be of the same format. Some data buffer descriptor formats use the contents of a count field (i.e., SRP_CMD request DATA-OUT BUFFER DESCRIPTOR COUNT or DATA-IN BUFFER DESCRIPTOR COUNT) to further describe the data buffer descriptor format. Table 2 defines data buffer descriptor format code values.

Table 2 - Data buffer descriptor formats

Data buffer descriptor format code	Reference	format code value ^a	buffer descriptor length (bytes) ^c
NO DATA BUFFER DESCRIPTOR PRESENT	5.6.2.3	0h	0
DIRECT DATA BUFFER DESCRIPTOR	5.6.2.4	1h	16
INDIRECT DATA BUFFER DESCRIPTOR	5.6.2.5	2h	20+16×count ^b
^a The format code value for a data-out buffer descriptor is specified by the DATA-OUT BUFFER DESCRIPTOR FORMAT field of an SRP_CMD request (see 6.8). The format code value for a data-in buffer descriptor is specified by the DATA-IN BUFFER DESCRIPTOR FORMAT field of an SRP_CMD request (see 6.8).			
^b The count field for a data-out buffer descriptor is the DATA-OUT BUFFER DESCRIPTOR COUNT field of an SRP_CMD request (see 6.8). The count field for a data-in buffer descriptor is the DATA-IN BUFFER DESCRIPTOR COUNT field of an SRP_CMD request (see 6.8).			
^c The length of a data buffer descriptor is determined from its format code value and the contents of its count field.			

5.6.2.2 Supported data buffer descriptor formats

The REQUIRED BUFFER FORMATS field (see table 3) of the SRP_LOGIN_REQ request (see 6.2) indicates the data buffer descriptor formats (see table 2) that an SRP initiator port may specify in requests sent on an RDMA channel. An SRP initiator port shall set the REQUIRED BUFFER FORMATS field to indicate all data buffer descriptor formats that the SRP initiator port may specify in SRP_CMD requests (see 6.8) sent on that RDMA channel. An SRP initiator port shall not issue an SRP_CMD request (see 6.8) indicating a data buffer descriptor format that was not indicated in the REQUIRED BUFFER FORMATS field value for that RDMA channel. SRP target ports are not required to check SRP_CMD requests for data buffer descriptor formats that were not indicated in the REQUIRED BUFFER FORMATS field value. If a target port does detect that an initiator has specified a descriptor format not indicated in the REQUIRED BUFFER FORMATS field, the target port shall send an SRP_T_LOGOUT request (see 6.6) with the reason code UNSUPPORTED FORMAT CODE VALUE SPECIFIED IN DATA-OUT BUFFER

DESCRIPTOR FORMAT FIELD or the reason code UNSUPPORTED FORMAT CODE VALUE SPECIFIED IN DATA-IN BUFFER DESCRIPTOR FORMAT FIELD, as appropriate.

An SRP target port may accept an RDMA channel establishment request and return an SRP_LOGIN_RSP response (see 6.3) if the SRP target port is able to support all of the data buffer descriptor formats indicated in the REQUIRED BUFFER FORMATS field on that RDMA channel. An SRP target port shall reject the RDMA channel establishment request and return an SRP_LOGIN_REJ response (see 6.4) with reason ONE OR MORE REQUESTED DATA BUFFER DESCRIPTOR FORMATS ARE NOT SUPPORTED if the SRP target port is unable to support one or more of the data buffer descriptor formats indicated in the REQUIRED BUFFER FORMATS field on that RDMA channel.

An SRP target port shall indicate the data buffer descriptor formats that it supports in the SUPPORTED BUFFER FORMATS field (see table 3) of the SRP_LOGIN_RSP response (see 6.3) or the SRP_LOGIN_REJ response (see 6.4). All SRP target ports shall support the DIRECT DATA BUFFER DESCRIPTOR format and may support other data buffer descriptor formats.

Table 3 defines the contents of both the REQUIRED BUFFER FORMATS field and the SUPPORTED BUFFER FORMATS field.

Table 3 - Supported data buffer descriptor formats

Byte	Bit	7	6	5	4	3	2	1	0
0		Reserved							
1		Reserved					IDBD	DDBD	Reserved

An SRP initiator port sets the IDBD (indirect data buffer descriptor) bit to one in a SRP_LOGIN_REQ request (see 6.2) if it requires that the target port support the INDIRECT DATA BUFFER DESCRIPTOR format.

The target port shall set the IDBD bit to one in an SRP_LOGIN_RSP response (see 6.3) or in an SRP_LOGIN_REJ response (see 6.4) if the SRP target port supports the INDIRECT DATA BUFFER DESCRIPTOR format. The IDBD bit shall be set to zero in an SRP_LOGIN_RSP response or in an SRP_LOGIN_REJ response if the SRP target port does not support the INDIRECT DATA BUFFER DESCRIPTOR format.

An SRP initiator port sets the DDBD (direct data buffer descriptor) bit to one in a SRP_LOGIN_REQ request (see 6.2) if it requires that the target port support the DIRECT DATA BUFFER DESCRIPTOR format.

The target port shall set the DDBD bit to one in an SRP_LOGIN_RSP response (see 6.3) or in an SRP_LOGIN_REJ response (see 6.4).

The length of requests sent by an SRP initiator port, as determined by the data buffer descriptor formats, shall be limited to the MAXIMUM INITIATOR TO TARGET IU LENGTH field returned in the SRP_LOGIN_RSP response (see 6.3).

5.6.2.3 No data buffer descriptor present

The NO DATA BUFFER DESCRIPTOR PRESENT format code value specifies that the corresponding data buffer descriptor field is not present. The contents of the count field in the SRP_CMD request (i.e., DATA-OUT BUFFER DESCRIPTOR COUNT or DATA-IN BUFFER DESCRIPTOR COUNT) are reserved. SRP target ports shall ignore the contents of the count field.

5.6.2.4 Direct data buffer descriptor format

The DIRECT DATA BUFFER DESCRIPTOR format code value specifies that the corresponding data buffer descriptor field is as defined in table 2 and contains a direct data buffer descriptor. The contents of the count field in the SRP_CMD request are reserved. SRP target ports shall ignore the contents of the count field.

Working Draft

A direct data buffer descriptor contains a single memory descriptor (see table 1). The memory descriptor identifies the data buffer, which is a single memory segment within a memory region's virtual address space. If a direct data buffer descriptor defines a data-out buffer, the SRP target port shall issue only RDMA Read operations using the memory descriptor contained in the direct data buffer descriptor. If a direct data buffer descriptor defines a data-in buffer, the SRP target port shall issue only RDMA Write operations using the memory descriptor contained in the direct data buffer descriptor. The SRP target port shall use the contents of the DATA LENGTH field of the memory descriptor as the length of the data-out buffer or data-in buffer.

5.6.2.5 Indirect data buffer descriptor format

The INDIRECT DATA BUFFER DESCRIPTOR format code value specifies that the corresponding data buffer descriptor field contains an indirect data buffer descriptor (see table 2).

An indirect data buffer is comprised of one or more memory segments. The memory segments may be discontinuous. The memory segments may be spread among several memory regions. The indirect data buffer is the concatenation of the memory segments listed in the indirect data buffer descriptor. Each memory segment may have any length supported by the RDMA communication service, including a length of zero bytes (see figure 5).

Table 4 shows the format of an indirect data buffer descriptor.

Table 4 - Indirect data buffer descriptor

Bit Byte	7	6	5	4	3	2	1	0
0	INDIRECT TABLE MEMORY DESCRIPTOR							
...								
15								
16	(MSB)	TOTAL LENGTH						
...								
19		(LSB)						
20	PARTIAL MEMORY DESCRIPTOR LIST							
...								
19+16xn								

^a The value 'n' is the value contained in the data buffer descriptor's count field. The count field for a data-out buffer descriptor is the DATA-OUT BUFFER DESCRIPTOR COUNT field of an SRP_CMD request (see 6.8). The count field for a data-in buffer descriptor is the DATA-IN BUFFER DESCRIPTOR COUNT field of an SRP_CMD request (see 6.8).

The INDIRECT TABLE MEMORY DESCRIPTOR is a memory descriptor (see table 1) that specifies a memory segment containing an indirect table. An indirect table is a list of one or more memory descriptors. The memory segments specified by the memory descriptors in the indirect table form the indirect data buffer. The value of the DATA LENGTH field of the INDIRECT TABLE MEMORY DESCRIPTOR represents the length, in bytes, of the indirect table, and is the number of memory descriptors in the indirect table multiplied by sixteen (the length, in bytes, of a memory descriptor). SRP target port behavior when the DATA LENGTH field of the INDIRECT TABLE MEMORY DESCRIPTOR contains any other value is vendor-specific.

The TOTAL LENGTH field value is the sum of the DATA LENGTH field values of the memory descriptors in the indirect table. The SRP target port shall use either the TOTAL LENGTH field value or the sum of the DATA LENGTH field values as the length of the data-out buffer or data-in buffer. If the value of the TOTAL LENGTH field is not equal to the sum of the values of the DATA LENGTH fields, SRP target port behavior is vendor-specific.

The PARTIAL MEMORY DESCRIPTOR LIST field is only present when the SRP_CMD information unit's data buffer descriptor's count field (i.e., DATA-OUT BUFFER DESCRIPTOR COUNT or DATA-IN BUFFER DESCRIPTOR COUNT) contains a non-zero value. The PARTIAL MEMORY DESCRIPTOR LIST field contains a list of 'n' memory descriptors that are copies of the first 'n' memory descriptors in the indirect table. The value 'n' is the value contained in the associated count field; SRP target port behavior when the PARTIAL MEMORY DESCRIPTOR LIST field contains any other value is vendor-specific.

5.6.2.5.1 SRP target port indirect data restrictions

An SRP target port shall issue only RDMA Read operations to the indirect table.

If an indirect data buffer descriptor specifies a data-out buffer, the SRP target port shall issue only RDMA Read operations using the memory descriptors contained in the indirect table or the PARTIAL MEMORY DESCRIPTOR LIST field value.

If an indirect data buffer descriptor specifies a data-in buffer, the SRP target port shall only issue RDMA Write operations using the memory descriptors contained in the indirect table or the PARTIAL MEMORY DESCRIPTOR LIST field value.

5.6.2.5.2 Examples of indirect data buffers

Figure 5 illustrates an indirect data buffer descriptor that does not contain a PARTIAL MEMORY DESCRIPTOR LIST field. Memory is shown containing four memory segments: the indirect table, memory segment 1, memory segment 2 and memory segment 3. The mapping of each memory descriptor to its memory segment has been shown as a single arrow. For details of this mapping see 5.6.1 and figure 4. Figure 5 does not show the memory regions in which the memory segments reside. All four segments may be in a single memory region, or may be in different memory regions.

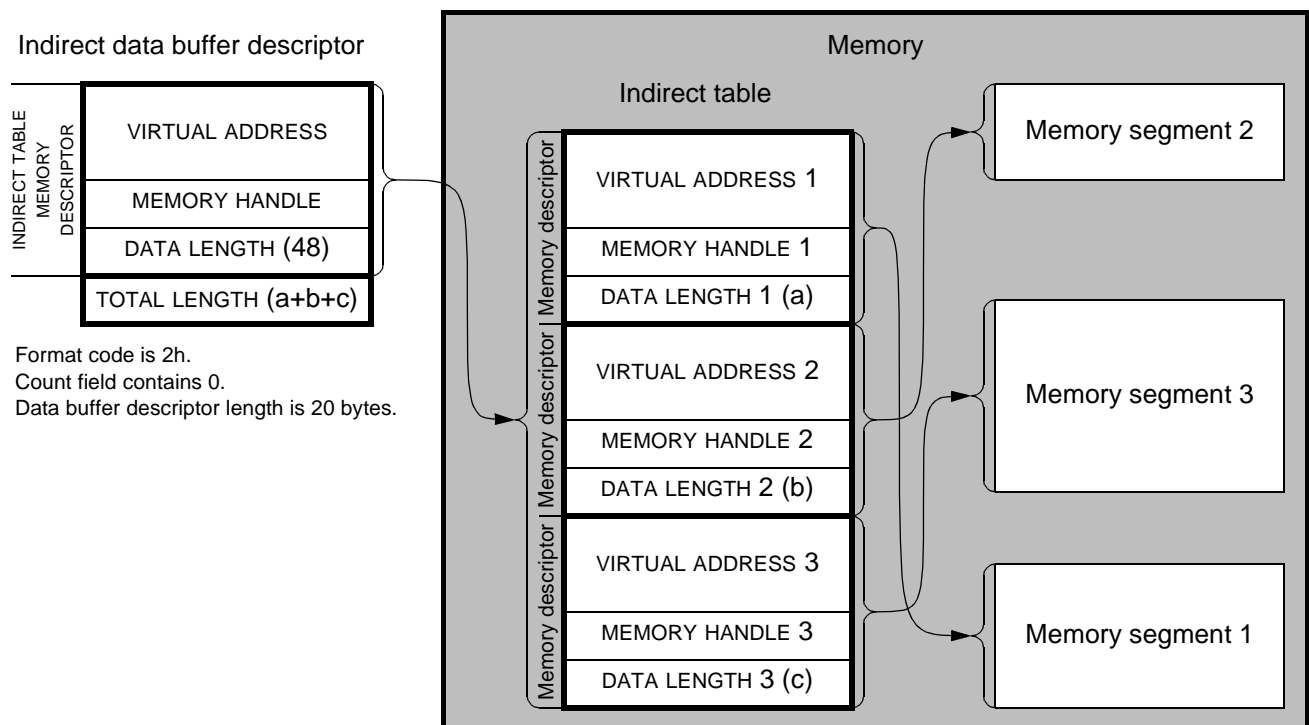
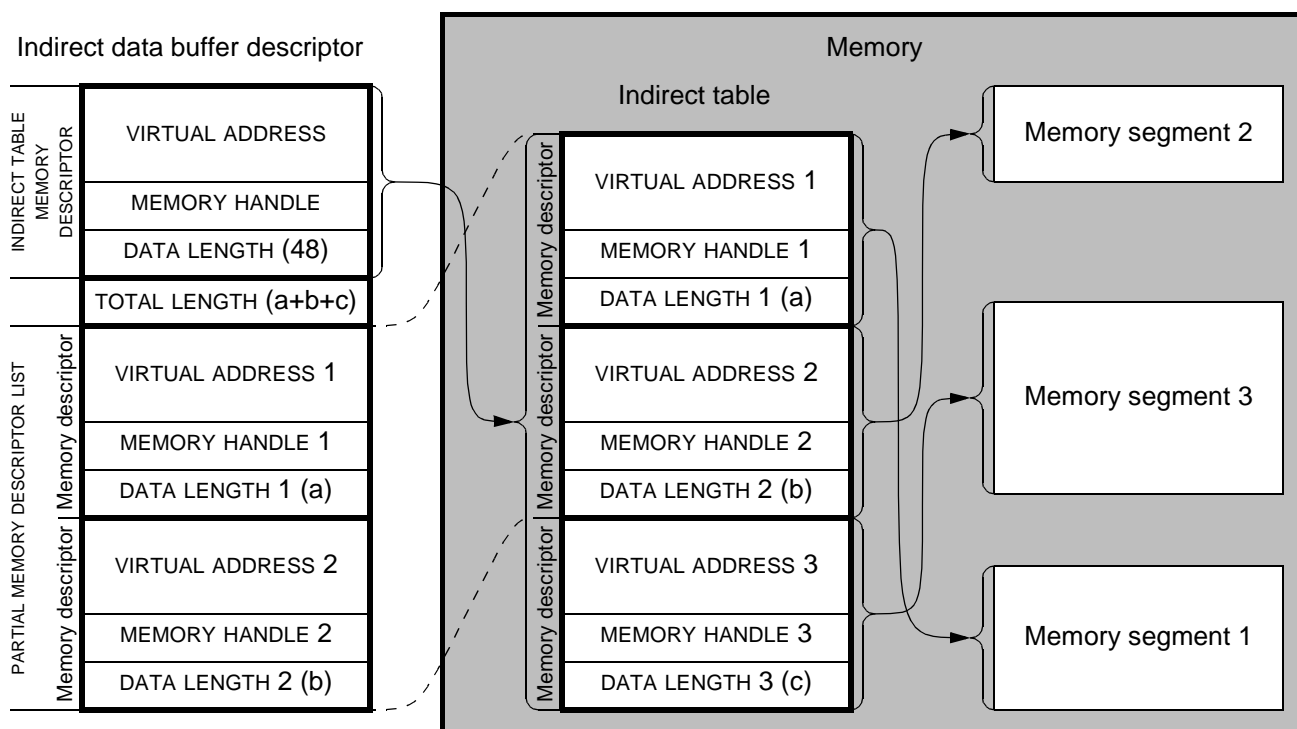


Figure 5 - Example indirect data buffer descriptor with no PARTIAL MEMORY DESCRIPTOR LIST field

In the example shown in figure 5 the data buffer descriptor format code value is 2h and the count field contains zero. The indirect data buffer descriptor is 20 bytes long. The data buffer is comprised of three memory segments: memory segment 1, memory segment 2 and memory segment 3. A separate memory segment

contains the indirect table, a list of three memory descriptors specifying memory segments 1 through 3. The INDIRECT TABLE MEMORY DESCRIPTOR field value of the indirect data buffer descriptor specifies the memory segment containing the indirect table. The DATA LENGTH field of the INDIRECT TABLE MEMORY DESCRIPTOR field value contains 48 (i.e. the length of the indirect table). The TOTAL LENGTH field of the data buffer descriptor contains the sum of the DATA LENGTH field values of the memory descriptors in the indirect table (i.e. the sum of DATA LENGTH 1, DATA LENGTH 2 and DATA LENGTH 3). This sum is the total length of the data buffer.

Figure 6 illustrates the same example as in figure 5 except with a PARTIAL MEMORY DESCRIPTOR LIST field. The data buffer, indirect table, INDIRECT TABLE MEMORY DESCRIPTOR field value and TOTAL LENGTH field value are all identical to the example in figure 5. The data buffer descriptor format code is 2h, the same as in figure 5. However the count field contains the value 2, indicating that the PARTIAL MEMORY DESCRIPTOR LIST field is present and contains two memory descriptors. Those two memory descriptors are copies of the first two memory descriptors in the indirect table. The third memory descriptor is only present in the indirect table. The indirect data buffer descriptor is 52 bytes long.



Format code is 2h.

Count field contains 2.

Data buffer descriptor length is 52 bytes.

Figure 6 - Example indirect data buffer descriptor with a PARTIAL MEMORY DESCRIPTOR LIST field

6 SRP Information Units

6.1 Summary

The information units used by SRP and their characteristics are shown in table 5, table 6, table 7 and table 8.

All SRP initiator ports shall support sending the information units listed in table 5 and table 8, and shall support receiving the information units listed in table 6 and table 7.

All SRP target ports shall support sending the information units listed in table 6 and table 7, and shall support receiving the information units listed in table 5 and table 8.

Table 5 - SRP requests sent from SRP initiator ports to SRP target ports

Information unit	Reference	TYPE value	Length (bytes)	Description
SRP_LOGIN_REQ	6.2	00h	64	Login request
SRP_TSK_MGMT	6.7	01h	64	SCSI task management function
SRP_CMD	6.8	02h	48 minimum	SCSI command
SRP_I_LOGOUT	6.5	03h	16	SRP initiator port logout notification

Table 6 - SRP responses sent from SRP target ports to SRP initiator ports

Information unit	Reference	TYPE value	Length (bytes)	Description
SRP_LOGIN_RSP	6.3	C0h	52	Login successful response
SRP_RSP	6.9	C1h	36 ^a minimum	SCSI status or service response
SRP_LOGIN_REJ	6.4	C2h	32	Login failure response
^a 36 bytes is not sufficient to return any sense data.				

Table 7 - SRP requests sent from SRP target ports to SRP initiator ports

Information unit	Reference	TYPE value	Length (bytes)	Description
SRP_T_LOGOUT	6.6	80h	16	SRP target port logout
SRP_CRED_REQ	6.10	81h	52	SRP target port credit adjustment request
SRP_AER_REQ	6.12	82h	36 ^a minimum	Asynchronous event report request
^a 36 bytes is not sufficient to return any sense data.				

Table 8 - SRP responses sent from SRP initiator ports to SRP target ports

Information unit	Reference	TYPE value	Length (bytes)	Description
SRP_CRED_RSP	6.11	41h	64	Response to SRP target port credit adjustment request
SRP_AER_RSP	6.13	42h	16	Asynchronous event report response

Working Draft

Byte 0 of each SRP information unit contains a TYPE code. The TYPE code value uniquely identifies the information unit and its format. The length of an information unit is indicated by its TYPE code and selected fields within the information unit. If an SRP target port receives an SRP information unit with an invalid TYPE code, or whose length is incorrect for the information unit's type code, the SRP target port shall send an SRP_T_LOGOUT request (see 6.6) and disconnect the RDMA channel.

Bytes 8 through 15 of each information unit contain a TAG value, which provides a mechanism for matching SRP requests with their corresponding SRP responses. Each SRP request information unit (see table 5 and table 7) shall contain a TAG value that is unique among all the outstanding SRP requests from a particular initiator port or target port. Each SRP response shall contain a copy of the TAG value from the corresponding SRP request. Responders are not required to check whether the TAG values of outstanding SRP requests are unique; if TAG values are not unique, responder behavior is unpredictable.

6.2 SRP_LOGIN_REQ request

An SRP_LOGIN_REQ request (see table 9) conveys SRP login parameters from an SRP initiator port to an SRP target port. The SRP_LOGIN_REQ request shall be sent as login data during RDMA channel establishment .

Table 9 - SRP_LOGIN_REQ request

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (00h)							
1	Reserved							
...								
7								
8	(MSB)	TAG						
...								
15	(LSB)							
16	(MSB)	REQUESTED MAXIMUM INITIATOR TO TARGET IU LENGTH						
...								
19	(LSB)							
20	Reserved							
...								
23								
24	REQUIRED BUFFER FORMATS							
25								
26	Reserved	AESOLNT	CRSOLNT	LOSOLNT	Reserved		MULTI-CHANNEL ACTION	
27	Reserved							
28	Reserved							
...								
31								
32	INITIATOR PORT IDENTIFIER							
...								
47								
48	TARGET PORT IDENTIFIER							
...								
63								

The TAG field is defined in 6.1.

The REQUESTED MAXIMUM INITIATOR TO TARGET IU LENGTH field specifies the maximum length in bytes of any information unit that the SRP initiator port sends on this RDMA channel. This value shall be greater than 63.

The REQUIRED BUFFER FORMATS field is defined in 5.6.2.2.

The asynchronous event solicited notification bit (AESOLNT) specifies whether an SRP_AER_REQ request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification. See 6.12.

The credit request solicited notification bit (CRSOLNT) specifies whether an SRP_CRED_REQ request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification. (See 6.10)

The logout solicited notification bit (LOSOLNT) specifies whether an SRP_T_LOGOUT request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification. (See 6.6)

The MULTI-CHANNEL ACTION field (see table 10) indicates how an SRP target port handles existing RDMA channels associated with the same I_T nexus.

Table 10 - MULTI-CHANNEL ACTION code values

MULTI-CHANNEL ACTION	Description
00b	Single RDMA channel operation (see 5.1.2)
01b	Multiple independent RDMA channel operation (see 5.1.3)
10b -11b	Reserved

The INITIATOR PORT IDENTIFIER field and the TARGET PORT IDENTIFIER field specify the I_T nexus that shall be associated with this RDMA channel.

6.3 SRP_LOGIN_RSP response

An SRP_LOGIN_RSP response (see table 11) indicates successful RDMA channel establishment and conveys SRP login parameters from an SRP target port to an SRP initiator port. An SRP_LOGIN_RSP response shall be sent as acceptance data during RDMA channel establishment (see 4.2).

Table 11 - SRP_LOGIN_RSP response

Bit Byte	7	6	5	4	3	2	1	0	
0	TYPE (C0h)								
1	Reserved								
2									
3									
4	(MSB)								
...	REQUEST LIMIT DELTA								
7								(LSB)	
8	(MSB)								
...	TAG								
15								(LSB)	
16	(MSB)								
...	MAXIMUM INITIATOR TO TARGET IU LENGTH								
19								(LSB)	
20	(MSB)								
...	MAXIMUM TARGET TO INITIATOR IU LENGTH								
23								(LSB)	
24	SUPPORTED BUFFER FORMATS								
25									
26	Reserved			SOLNTSUP		Reserved		MULTI-CHANNEL RESULT	
27	Reserved								
28									
...	Reserved								
51									

The REQUEST LIMIT DELTA field is defined in 5.5.

The TAG field shall contain the same value as the TAG field in the SRP_LOGIN_REQ request (see 6.2).

MAXIMUM INITIATOR TO TARGET IU LENGTH specifies the maximum length in bytes of any information unit that the SRP target port is able to receive on this RDMA channel. This value shall be 64 or larger and greater than or equal to the value of REQUESTED MAXIMUM INITIATOR TO TARGET IU LENGTH specified in the SRP_LOGIN_REQ request (see 6.2). The SRP initiator port shall not send any information unit on this RDMA channel longer than this value.

MAXIMUM TARGET TO INITIATOR IU LENGTH specifies the maximum length in bytes of any information unit that the SRP target port may send on this RDMA channel. This value shall be 52 or larger. The SRP target port shall not send any information unit on this RDMA channel longer than this value.

The SUPPORTED BUFFER FORMATS field is defined in 5.6.2.2.

Working Draft

The MULTI-CHANNEL RESULT field (see table 12) indicates how the SRP target port treated existing RDMA channels associated with the same I_T nexus.

Table 12 - MULTI-CHANNEL RESULT code values

MULTI-CHANNEL RESULT	Description
00b	No existing RDMA channels were associated with the same I_T nexus.
01b	One or more existing RDMA channels were terminated.
10b	One or more existing RDMA channels continue to operate independently.
11b	Reserved

The solicited notification supported bit (SOLNTSUP) indicates whether the SRP target port supports solicited message reception notification for messages sent from the SRP target port to an SRP initiator port (see 4.3). If the SOLNTSUP bit is one, the SRP target port supports solicited message reception notification. If the SOLNTSUP bit is zero, the SRP target port only supports normal message reception notification.

6.4 SRP_LOGIN_REJ response

An SRP_LOGIN_REJ response (see table 13) indicates that an RDMA channel could not be established. An SRP_LOGIN_REJ response shall be sent as rejection data (see 4.2).

Table 13 - SRP_LOGIN_REJ response

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (C2h)							
1	Reserved							
2								
3								
4	(MSB)	REASON						
...								
7	(LSB)							
8	(MSB)	TAG						
...								
15	(LSB)							
16	Reserved							
...								
23								
24	SUPPORTED BUFFER FORMATS							
25	Reserved							
26								
...								
31								

The REASON field indicates the reason that the RDMA channel could not be established. This field is defined in table 14.

Table 14 - SRP_LOGIN_REJ response reason codes

REASON code	Description
0000 0000h - 0000 FFFFh	Reserved
0001 0000h	Unable to establish RDMA channel, no reason specified
0001 0001h	Insufficient RDMA channel resources
0001 0002h	REQUESTED MAXIMUM INITIATOR TO TARGET IU LENGTH value too large
0001 0003h	Unable to associate RDMA channel with specified I_T nexus
0001 0004h	One or more requested data buffer descriptor formats are not supported
0001 0005h	SRP target port does not support multiple RDMA channels per I_T nexus
0001 0006h	RDMA channel limit reached for this initiator
0001 0007h - FFFF FFFFh	Reserved

The TAG field shall contain the same value as the TAG field in the SRP_LOGIN_REQ request (see 6.2).

The SUPPORTED BUFFER FORMATS field is defined in 5.6.2.2.

6.5 SRP_I_LOGOUT request

An SRP_I_LOGOUT request (see table 15) is sent by an SRP initiator port to notify the SRP target port that the SRP initiator port is disconnecting the RDMA channel. An SRP_I_LOGOUT request shall be sent as a 16-byte message with normal message reception notification (see 4.3).

Table 15 - SRP_I_LOGOUT request

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (03h)							
1	Reserved							
2								
7								
8	(MSB)	TAG						(LSB)
...								
15								

The TAG field is defined in 6.1.

After sending an SRP_I_LOGOUT request, an SRP initiator port shall wait for the RDMA communication service status indication (see 4.5.2), then request that the RDMA channel be disconnected and perform the actions specified in 5.1.4.

Upon receiving an SRP_I_LOGOUT request an SRP target port shall perform the actions specified in 5.1.4. The SRP target port shall not send an SRP response to an SRP_I_LOGOUT request.

6.6 SRP_T_LOGOUT request

An SRP_T_LOGOUT request (see table 16) is sent by a SRP target port to notify the SRP initiator port that the SRP target port is disconnecting the RDMA channel. An SRP_T_LOGOUT request shall be sent as a 16-byte message.

Table 16 - SRP_T_LOGOUT request

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (80h)							
1	Reserved							SOLNT
2								
3								
4	(MSB)							
...		REASON						
7								(LSB)
8	(MSB)							
...		TAG						
15								(LSB)

The solicited notification (SOLNT) bit indicates whether the SRP initiator port specified normal or solicited message reception notification for SRP_T_LOGOUT requests during login (see 6.2). The SOLNT bit shall contain the value that was specified in the LOSOLNT bit of the SRP_LOGIN_REQ request.

If the solicited notification (SOLNT) bit is one and the SRP target port supports solicited message reception notification (see 6.3), the SRP target port shall send the SRP_T_LOGOUT response with solicited message reception notification (see 4.3). If the SOLNT bit is zero, the SRP target port should send the SRP_T_LOGOUT response with normal message reception notification. An SRP initiator port shall not validate the SOLNT bit against whether an SRP_RSP response was actually received with normal or solicited message reception notification.

The REASON field indicates the reason for disconnecting the RDMA channel. This field is defined in table 17

Table 17 - SRP_T_LOGOUT request reason codes

REASON code	Description
0000 0000h	No reason specified
0000 0001h	Inactive RDMA channel (reclaiming resources)
0000 0002h	Invalid information unit TYPE code received by SRP target port
0000 0003h	SRP initiator port sent response ^a with no corresponding SRP target port request ^b outstanding
0000 0004h	RDMA channel disconnected due to MULTI-CHANNEL ACTION code in new SRP_LOGIN_REQ
0000 0005h	Reserved
0000 0006h	Unsupported format code value specified in DATA-OUT BUFFER DESCRIPTOR FORMAT field
0000 0007h	Unsupported format code value specified in DATA-IN BUFFER DESCRIPTOR FORMAT field
0000 0008h	Invalid length for IU type
0000 0005h FFFF_FFFFh	Reserved
^a See table 8.	
^b See table 7.	

The TAG field is defined in 6.1.

After sending an SRP_T_LOGOUT request, an SRP target port shall wait for the RDMA communication service status indication (see 4.5.2), then request that the RDMA channel be disconnected and perform the actions specified in 5.1.4.

Upon receiving an SRP_T_LOGOUT request an SRP initiator port shall perform the actions specified in 5.1.4. The SRP initiator port shall not send an SRP response to an SRP_T_LOGOUT request.

6.7 SRP_TSK_MGMT request

An SRP_TSK_MGMT request conveys a SCSI task management request (table 18). An SRP_TSK_MGMT request shall be sent with normal message reception notification (see 4.3).

Table 18 - SRP_TSK_MGMT request

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (01h)							
1	Reserved					UCSOLNT	SCSOLNT	Reserved
...	Reserved							
7	Reserved							
8	(MSB)							
...	TAG							
15								(LSB)
16	Reserved							
...	Reserved							
19	Reserved							
20	(MSB)							
...	LOGICAL UNIT NUMBER							
27								(LSB)
28	Reserved							
29	Reserved							
30	TASK MANAGEMENT FUNCTION							
31	Reserved							
32	(MSB)							
...	TAG OF TASK TO BE MANAGED							
39								(LSB)
40	Reserved							
...	Reserved							
47	Reserved							

The unsuccessful completion solicited notification bit (UCSOLNT) specifies whether an SRP_RSP response reporting unsuccessful completion of the task management request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification. See 6.9.

The successful completion solicited notification bit (SCSOLNT) specifies whether an SRP_RSP response reporting successful completion of the task management request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification. See 6.9.

The TAG field is defined in 6.1.

The LOGICAL UNIT NUMBER field identifies the logical unit to which the task management request is directed. The structure of the LOGICAL UNIT NUMBER field shall be as defined in the SCSI Architecture Model-2 standard. This field is reserved if the task management request is not directed to either an I_T_L or I_T_L_Q nexus.

The TASK MANAGEMENT FUNCTION field is defined in table 19. If TASK MANAGEMENT FUNCTION contains a reserved or restricted value, the task manager shall return an SRP_RSP response (see 6.9) containing GOOD status. The RSP_CODE field shall be set to TASK MANAGEMENT FUNCTION NOT SUPPORTED.

Table 19 - TASK MANAGEMENT FUNCTION codes

Code	Description
00h	Reserved
01h	The task manager shall perform an ABORT TASK function (see SAM-2).
02h	The task manager shall perform an ABORT TASK SET function (see SAM-2).
03h	Reserved
04h	The task manager shall perform a CLEAR TASK SET function (see SAM-2).
05h-07h	Reserved
08h	The task manager shall perform a LOGICAL UNIT RESET function (see SAM-2).
09h-1Fh	Reserved
20h	Restricted.
21h-3Fh	Reserved
40h	The task manager shall perform a CLEAR ACA function (see SAM-2).
41h-FFh	Reserved

If TASK MANAGEMENT FLAGS specifies that an ABORT TASK function shall be performed, the TAG OF TASK TO BE MANAGED field specifies the TAG value from the SRP_CMD request (see 6.8) that contained the task to be aborted. The TAG OF TASK TO BE MANAGED field shall be ignored if TASK MANAGEMENT FLAGS specifies any other function.

6.8 SRP_CMD request

An SRP_CMD request conveys a SCSI command (see table 20).

Table 20 - SRP_CMD request

Byte	Bit	7	6	5	4	3	2	1	0
0		TYPE (02h)							
1		Reserved					UCSOLNT	SCSOLNT	Reserved
...		Reserved							
4									
5		DATA-OUT BUFFER DESCRIPTOR FORMAT				DATA-IN BUFFER DESCRIPTOR FORMAT			
6		DATA-OUT BUFFER DESCRIPTOR COUNT							
7		DATA-IN BUFFER DESCRIPTOR COUNT							
8	(MSB)	TAG							
...									
15		(LSB)							
16		Reserved							
...									
19									
20	(MSB)	LOGICAL UNIT NUMBER							
...									
27		(LSB)							
28		Reserved							
29		Reserved					TASK ATTRIBUTE		
30		Reserved							
31		ADDITIONAL CDB LENGTH = n						Reserved	
32		CDB							
...									
47									
48		ADDITIONAL CDB							
...									
47+4xn									
48+4xn		DATA-OUT BUFFER DESCRIPTOR							
...									
47+4xn+do ^a									
48+4xn+do ^a		DATA-IN BUFFER DESCRIPTOR							
...									
47+4xn+do+di ^b									

^a The value 'do' is the length in bytes of the DATA-OUT BUFFER DESCRIPTOR field, determined from the format code value contained in the DATA-OUT BUFFER DESCRIPTOR FORMAT field and the count value contained in the DATA-OUT BUFFER DESCRIPTOR COUNT field (see 5.6.2).

^b The value 'di' is the length in bytes of the DATA-IN BUFFER DESCRIPTOR field, determined from the format code value contained in the DATA-IN BUFFER DESCRIPTOR FORMAT field and the count value contained in the DATA-IN BUFFER DESCRIPTOR COUNT field (see 5.6.2).

An SRP_CMD request shall be sent as a message whose length is 48 bytes plus the lengths of the ADDITIONAL CDB, DATA-OUT BUFFER DESCRIPTOR, and DATA-IN BUFFER DESCRIPTOR fields. An SRP_CMD request shall be sent with normal message reception notification (see 4.3).

The unsuccessful completion solicited notification bit (UCSOLNT) specifies whether an SRP_RSP response reporting unsuccessful completion of the task management request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification (see 6.9).

The successful completion solicited notification bit (SCSOLNT) specifies whether an SRP_RSP response reporting successful completion of the task management request should use normal or solicited message reception notification. This bit shall be set to one to request solicited notification, or set to zero to request normal notification (see 6.9).

The DATA-OUT BUFFER DESCRIPTOR FORMAT field specifies the format of the DATA-OUT BUFFER DESCRIPTOR field (see 5.6.2).

The DATA-IN BUFFER DESCRIPTOR FORMAT field specifies the format of the DATA-IN BUFFER DESCRIPTOR field (see 5.6.2).

The DATA-OUT BUFFER DESCRIPTOR COUNT field provides additional information to specify the format of the DATA-OUT BUFFER DESCRIPTOR field (see 5.6.2).

The DATA-IN BUFFER DESCRIPTOR COUNT field provides additional information to specify the format of the DATA-IN BUFFER DESCRIPTOR field (see 5.6.2).

The TAG field is defined in 6.1.

The LOGICAL UNIT NUMBER field specifies the address of the logical unit of the I_T_L_Q nexus for the current task. The structure of the logical unit number field shall be as defined in the SCSI Architecture Model-2 standard. If the addressed logical unit does not exist, the task manager shall follow the SCSI rules for selection of invalid logical units as defined in the SCSI Primary Commands-2 standard.

The TASK ATTRIBUTE field is defined in table 21.

Table 21 - TASK ATTRIBUTE

Codes	Description
000b	Requests that the task be managed according to the rules for a simple task attribute. (See SAM-2)
001b	Requests that the task be managed according to the rules for a head of queue task attribute. (See SAM-2)
010b	Requests that the task be managed according to the rules for an ordered attribute. (See SAM-2)
011b	Reserved
100b	Requests that the task be managed according to the rules for an automatic contingent allegiance task attribute. (See SAM-2)
101b-111b	Reserved

The ADDITIONAL CDB LENGTH field contains the length, in four-byte words, of the ADDITIONAL CDB field.

The CDB and ADDITIONAL CDB fields together contain the CDB to be interpreted by the addressed logical unit. Any bytes between the end of the CDB and the end of the two fields shall be reserved.

Working Draft

The contents of the CDB shall be as defined in the SCSI command standards.

The DATA-OUT BUFFER DESCRIPTOR field specifies the buffer that shall be used for data-out transfers (see 5.6.2).

The DATA-IN BUFFER DESCRIPTOR field specifies the buffer that shall be used for data-in transfers (see 5.6.2).

6.9 SRP_RSP response

An SRP_RSP response (see table 22) conveys an SRP response to an SRP_TSK_MGMT request (see 6.7) or an SRP_CMD request (see 6.8) received by a SRP target port. SRP_RSP responses that contain neither RESPONSE DATA nor SENSE DATA shall be sent as a 36 byte message. SRP_RSP responses that contain either RESPONSE DATA or SENSE DATA shall be sent as the minimum length message containing those fields.

Table 22 - SRP_RSP response

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (C1h)							
1	Reserved							SOLNT
2	Reserved							
3								
4	(MSB)	REQUEST LIMIT DELTA						
...								
7		(LSB)						
8	(MSB)	TAG						
...								
15		(LSB)						
16	Reserved							
17								
18	Reserved	DIUNDER	DIOVER	DOUNDER	DOOVER	SNSVALID	RSPVALID	
19	STATUS							
20	(MSB)	DATA-OUT RESIDUAL COUNT						
...								
23		(LSB)						
24	(MSB)	DATA-IN RESIDUAL COUNT						
...								
27		(LSB)						
28	(MSB)	SENSE DATA LIST LENGTH = n						
...								
31		(LSB)						
32	(MSB)	RESPONSE DATA LIST LENGTH = m						
...								
35		(LSB)						
36	(MSB)	RESPONSE DATA (m bytes long)						
...								
35+m		(LSB)						
36+m	(MSB)	SENSE DATA (n bytes long)						
...								
35+m+n		(LSB)						

The solicited notification (SOLNT) bit indicates whether the SRP initiator port specified normal or solicited message reception notification for this response. If the STATUS field is non-zero or if the RSP_CODE field is present

Working Draft

and non-zero, then the SOLNT bit shall contain the value that was specified in the UCSOLNT bit of the corresponding SRP_CMD or SRP_TSK_MGMT request; otherwise, the SOLNT bit shall contain the value that was specified in the SCSOLNT bit of the corresponding SRP_CMD or SRP_TSK_MGMT request.

If the solicited notification (SOLNT) bit is one and the SRP target port supports solicited message reception notification (see 6.3), the SRP target port shall send the SRP_RSP response with solicited message reception notification (see 4.3); otherwise, the SRP target port should send the SRP_RSP response with normal message reception notification. An SRP initiator port shall not validate the SOLNT bit against whether an SRP_RSP response was actually received with normal or solicited message reception notification.

The REQUEST LIMIT DELTA field is defined in 5.5.

The TAG field shall contain the same value as the TAG field in the SRP_TSK_MGMT request (see 6.7) or SRP_CMD request (see 6.8) for which this SRP_RSP response is a response.

DOUNDER, when set to one, indicates that the DATA-OUT RESIDUAL COUNT field is valid and contains the count of data bytes that were expected to be transferred from the data-out buffer, but were not transferred. The application client should examine the DATA-OUT RESIDUAL COUNT field in the context of the command to determine whether or not an error condition occurred.

DOOVER, when set to one, indicates that the DATA-OUT RESIDUAL COUNT field is valid and contains the count of data bytes that could not be transferred from the data-out buffer because the length of the data-out buffer was not sufficient. The application client should examine the DATA-OUT RESIDUAL COUNT field in the context of the command to determine whether or not an error condition occurred.

DOUNDER and DOOVER, when both set to zero, indicate that the DATA-OUT RESIDUAL COUNT field is not valid; the SRP initiator port shall ignore its contents. The SRP target port shall not set both DOUNDER and DOOVER to one.

DIUNDER, when set to one, indicates that the DATA-IN RESIDUAL COUNT field is valid and contains the count of data bytes that were expected to be transferred to the data-in buffer, but were not transferred. The application client should examine the DATA-IN RESIDUAL COUNT field in the context of the command to determine whether or not an error condition occurred.

DIOVER, when set to one, indicates that the DATA-IN RESIDUAL COUNT field is valid and contains the count of data bytes that could not be transferred to the data-in buffer because the length of the data-in buffer was not sufficient. The application client should examine the DATA-IN RESIDUAL COUNT field in the context of the command to determine whether or not an error condition occurred.

DIUNDER and DIOVER, when both set to zero, indicate that the DATA-IN RESIDUAL COUNT field is not valid; the SRP initiator port shall ignore its contents. The SRP target port shall not set both DIUNDER and DIOVER to one.

SNSVALID, when set to zero, indicates that the contents of the SENSE DATA LIST LENGTH field shall be ignored and that the SENSE DATA field is not present. SNSVALID, when set to one, indicates that the contents of the SENSE DATA LIST LENGTH field specify the number of bytes in the SENSE DATA field.

If sense data is provided, SNSVALID shall be set to one and the SENSE DATA LIST LENGTH field shall specify the number of bytes in the SENSE DATA field.

If returning all the sense data provided would cause the SRP_RSP response to be longer than the value of the MAXIMUM TARGET TO INITIATOR IU LENGTH field indicated in the SRP_LOGIN_RSP response (see 6.3) when the RDMA channel was established, the SRP target port shall return an SRP_RSP response whose length is the value from the MAXIMUM TARGET TO INITIATOR IU LENGTH field. The SENSE DATA field shall be truncated as needed to achieve this length. SENSE DATA LIST LENGTH shall contain the length of the truncated SENSE DATA field.

If no sense data is provided, SNSVALID shall be set to zero. The SRP initiator port shall ignore the SENSE DATA LIST LENGTH field and shall assume a length of zero.

RSPVALID set to zero indicates that the contents of the RESPONSE DATA LIST LENGTH field shall be ignored and the RESPONSE DATA field is not present. RSPVALID set to one indicates that the contents of the RESPONSE DATA LIST LENGTH field specify the number of bytes in the RESPONSE DATA field. RSPVALID set to one also indicates that the contents of the STATUS field shall be ignored by the SRP initiator port.

If response data is provided, RSPVALID shall be set to one and the RESPONSE DATA LIST LENGTH field shall specify the number of bytes in the RESPONSE DATA field (see table 23). The RESPONSE DATA LIST LENGTH field shall contain the value four. Other lengths are reserved for future standardization.

If no response data is provided, RSPVALID shall be set to zero. The SRP initiator port shall ignore the RESPONSE DATA LIST LENGTH field and shall assume a length of zero.

Response data shall be provided in any SRP_RSP response that is sent in response to an SRP_TSK_MGMT request (see 6.7). The information in the RSP_CODE field (see table 24) shall indicate the completion status of the task management function.

Response data shall not be provided in any SRP_RSP response that returns a non-zero status code in the STATUS field.

The STATUS field contains the status of a task that completes. See the SAM-2 standard for a list of status codes.

If either DOUNDER or DOOVER is set to one, the DATA-OUT RESIDUAL COUNT field contains a count of the number of residual data bytes that were not transferred from the data-out buffer for this SCSI command. Upon successful completion of an SRP I/O operation, the residual data-out byte count is normally zero and the DATA-OUT RESIDUAL COUNT value is not valid. Some commands may have a non-zero residual data-out byte count that is not an error. SRP target ports are not required to check the data-out length implied by the contents of the CDB for overrun or underrun before processing a SCSI command.

If DOUNDER is set to one and a transfer that did not fill the entire data-out buffer was performed, the value of DATA-OUT RESIDUAL COUNT is defined as follows:

$$\text{DATA-OUT RESIDUAL COUNT} = (\text{data-out buffer length}) - (\text{highest offset of any data-out byte transmitted} + 1)$$

A condition of DOUNDER set to one may not be an error for some devices and some commands.

If DOOVER is set to one, the transfer was truncated because the data-out transfer required by the SCSI command was longer than the data-out buffer (see 5.6.2). Those bytes that could not be transferred without exceeding the length of the data-out buffer shall not be transferred. The DATA-OUT RESIDUAL COUNT is defined as follows:

$$\text{DATA-OUT RESIDUAL COUNT} = (\text{Transfer length required by command}) - (\text{data-out buffer length})$$

If DOOVER is set to one, the termination state of the SRP I/O operation is not certain. Data may not have been transferred from the data-out buffer and the SCSI status byte may not provide correct command completion information.

If either DIUNDER or DIOVER is set to one, the DATA-IN RESIDUAL COUNT field contains a count of the number of residual data bytes that were not transferred to the data-in buffer for this SCSI command. Upon successful completion of an SRP I/O operation, the residual data-in byte count is normally zero and the DATA-IN RESIDUAL COUNT value is not valid. Some commands (e.g., INQUIRY) may have a non-zero residual data-in byte count that is not an error. SRP target ports are not required to check the data-in length implied by the contents of the CDB for overrun or underrun before processing a SCSI command.

If DIUNDER is set to one and a transfer that did not fill the entire data-in buffer was performed, the value of DATA-IN RESIDUAL COUNT is defined as follows:

$$\text{DATA-IN RESIDUAL COUNT} = (\text{data-in buffer length}) - (\text{highest offset of any data-in byte transmitted} + 1)$$

A condition of DIUNDER set to one may not be an error for some devices and some commands.

If DIOVER is set to one, the transfer was truncated because the data-in transfer required by the SCSI command was longer than the data-in buffer (see 5.6.2). Those bytes that could not be transferred without exceeding the length of the data-in buffer shall not be transferred. The DATA-IN RESIDUAL COUNT is defined as follows:

$$\text{DATA-IN RESIDUAL COUNT} = (\text{Transfer length required by command}) - (\text{data-in buffer length})$$

If DIOVER is set to one, the termination state of the SRP I/O operation is not certain. Data may not have been transferred to the data-in buffer and the SCSI status byte may not provide correct command completion information.

The DATA-OUT RESIDUAL COUNT, DATA-IN RESIDUAL COUNT, SENSE DATA LIST LENGTH, and RESPONSE DATA LIST LENGTH fields shall always be present in the SRP_RSP response, regardless of whether their contents are valid.

The RESPONSE DATA field (see table 23) contains information describing protocol failures detected during processing of an SRP request received by the SRP target port. The RESPONSE DATA field shall be present if the SRP target port detects any of the conditions described by a non-zero RSP_CODE value (see table 24).

Table 23 - RESPONSE DATA field

Bit Byte	7	6	5	4	3	2	1	0
0	Reserved							
1	Reserved							
2	Reserved							
3	RSP_CODE							

The RSP_CODE field is defined in table 24.

Table 24 - RSP_CODE values

Codes	Description
00h	TASK MANAGEMENT FUNCTION COMPLETE.
01h	Reserved
02h	REQUEST FIELDS INVALID
03h	Reserved
04h	TASK MANAGEMENT FUNCTION NOT SUPPORTED
05h	TASK MANAGEMENT FUNCTION FAILED
06h-FFh	Reserved

The SENSE DATA field contains the autosense data specified by the SCSI Primary Commands-2 standard. The proper sense data shall be presented when the SCSI status byte of CHECK CONDITION is presented as specified by the SCSI Primary Commands-2 standard. If no conditions requiring the presentation of SCSI sense data have occurred, the SENSE DATA field shall not be included in the SRP_RSP response and the SNSVALID bit shall be zero. SRP devices shall perform autosense.

6.10 SRP_CRED_REQ request

An SRP target port may use SRP_CRED_REQ requests (see table 25) to adjust an SRP initiator port's REQUEST LIMIT value (see 5.5). An SRP_CRED_REQ request shall be sent as a 16 byte message.

Table 25 - SRP_CRED_REQ request

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (81h)							
1	Reserved							SOLNT
2	Reserved							
3								
4	(MSB)	REQUEST LIMIT DELTA						
...								
7								(LSB)
8	(MSB)	TAG						
...								
15								(LSB)

The solicited notification (SOLNT) bit indicates whether the SRP initiator port specified normal or solicited message reception notification during login (see 6.2) for SRP_CRED_REQ requests. The SOLNT bit shall contain the value that was specified in the CRSOLNT bit of the SRP_LOGIN_REQ request.

If the solicited notification (SOLNT) bit is one and the SRP target port supports solicited message reception notification (see 6.3), the SRP target port shall send the SRP_CRED_REQ request with solicited message reception notification (see 4.3); otherwise the SRP target port should send the SRP_CRED_REQ request with normal message reception notification. An SRP initiator port shall not validate the SOLNT bit against whether an SRP_CRED_REQ request was actually received with normal or solicited message reception notification.

The REQUEST LIMIT DELTA field is defined in 5.5.

The TAG field is defined in 6.1.

6.11 SRP_CRED_RSP response

An SRP_CRED_RSP response (see table 26) is the response to an SRP_CRED_REQ request (see 6.10) received by an SRP initiator port. All SRP initiator ports shall support generating SRP_CRED_RSP responses. An SRP_CRED_RSP response shall be sent as a 16-byte message with normal message reception notification (see 4.3).

Table 26 - SRP_CRED_RSP response

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (41h)							
1	Reserved							
7								
8								
...	(MSB)	TAG						
15								(LSB)

The TAG field shall contain the same value as the TAG field in the SRP_CRED_REQ request (see 6.10) for which this SRP_CRED_RSP response is a response.

6.12 SRP_AER_REQ request

A target port sends an SRP_AER_REQ request (see table 27) to report an asynchronous event. An SRP_AER_REQ request shall be sent as the minimum length message capable of carrying the fields. Parameters managing the use of asynchronous event reporting are contained in the Control mode page (see SPC-2).

Table 27 - SRP_AER_REQ request

Bit Byte	7	6	5	4	3	2	1	0						
0	TYPE (82h)													
1	Reserved							SOLNT						
2	Reserved													
3														
4	(MSB)	REQUEST LIMIT DELTA												
...														
7	(LSB)													
8	(MSB)	TAG												
...														
15	(LSB)													
16	Reserved													
...														
19														
20	(MSB)	LOGICAL UNIT NUMBER												
...														
27	(LSB)													
28	(MSB)	SENSE DATA LIST LENGTH = n												
...														
31	(LSB)													
32	Reserved													
...														
35														
36	(MSB)	SENSE DATA (n bytes long)												
...														
35+n	(LSB)													

The solicited notification (SOLNT) bit indicates whether the SRP initiator port specified normal or solicited message reception notification during login (see 6.2) for SRP_AER_REQ requests. The SOLNT bit shall contain the value that was specified in the CRSOLNT bit of the SRP_LOGIN_REQ request.

If the solicited notification (SOLNT) bit is one and the SRP target port supports solicited message reception notification (see 6.3), the SRP target port shall send the SRP_AER_REQ request with solicited message reception notification (see 4.3); otherwise the SRP target port should send the SRP_AER_REQ request with normal message reception notification. An SRP initiator port shall not validate the SOLNT bit against whether an SRP_AER_REQ request was actually received with normal or solicited message reception notification.

The REQUEST LIMIT DELTA field is defined in 5.5.

Working Draft

The TAG field is defined in 6.1.

The SENSE DATA LIST LENGTH field shall specify the number of bytes in the SENSE DATA field.

If including all the sense data provided would cause the SRP_AER_REQ request to be longer than the value of the MAXIMUM TARGET TO INITIATOR IU LENGTH field indicated in the SRP_LOGIN_RSP response (see 6.3) when the RDMA channel was established, the SRP target port shall send an SRP_AER_REQ request whose length is the MAXIMUM TARGET TO INITIATOR IU LENGTH field value. The SENSE DATA field shall be truncated as needed to achieve this length. SENSE DATA LIST LENGTH shall contain the length of the truncated SENSE DATA field.

The SENSE DATA field contains sense data as specified by the SCSI Primary Commands-2 standard.

6.13 SRP_AER_RSP response

An SRP_AER_RSP response (see table 28) conveys an SRP initiator port's SRP response to an SRP_AER_REQ request (see 6.12). An SRP_AER_RSP response shall be sent as a 16-byte message with normal message reception notification (see 4.3)..

Table 28 - SRP_AER_RSP response

Bit Byte	7	6	5	4	3	2	1	0
0	TYPE (42h)							
1								
2	Reserved							
7								
8	(MSB)	TAG						
...								
15								(LSB)

The TAG field shall contain the same value as the TAG field in the SRP_AER_REQ request (see 6.12) for which this SRP_AER_RSP response is a response.

Working Draft

7 SCSI mode parameters

7.1 SCSI mode parameter overview and codes

This subclause describes the block descriptors and the pages used with MODE SELECT and MODE SENSE commands that influence, control, and report the behavior of the SRP target port. All mode parameters not defined in this standard shall influence the behavior of the SCSI devices as specified in the appropriate command set document. The mode pages are addressed to the device server of a logical unit. The mode pages associated with this protocol are listed in table 29.

Table 29 - SRP mode page codes

Page code	Description	Subclause
02h	Disconnect-reconnect page	7.2
18h	Protocol specific LUN page	7.3
19h	Protocol specific port page	7.4

7.2 Disconnect-reconnect mode page

The disconnect-reconnect page (see table 30) provides the application client the means to tune the performance of the service delivery subsystem. This subclause defines the fields in the disconnect-reconnect mode page of the MODE SENSE or MODE SELECT command that are used by SRP target ports.

Table 30 - Disconnect-reconnect mode page

Bit Byte	7	6	5	4	3	2	1	0
0	PS	RESERVED	PAGE CODE (02h)					
1	PAGE LENGTH (0EH)							
2	BUFFER FULL RATIO							
3	BUFFER EMPTY RATIO							
4	BUS INACTIVITY LIMIT							
5								
6	PHYSICAL DISCONNECT TIME LIMIT							
7								
8	CONNECT TIME LIMIT							
9								
10	(MSB)	MAXIMUM BURST SIZE						(LSB)
11								
12	EMDP	FAIR ARBITRATION			DIMM	DTDC		
13	RESERVED							
14	FIRST BURST SIZE							
15								

The application client passes the fields used to control the SRP target port to a device server by means of a MODE SELECT command. The device server then communicates the field values to the SRP target port. The field values are communicated from the device server to the SRP target port in a vendor-specific manner.

Working Draft

7.2.1 Valid fields

SRP devices shall use only the Disconnect-Reconnect page fields listed in this subclause. If any other fields (see 7.2.2) within the disconnect-reconnect page of the MODE SELECT command contain a non-zero value, the device server shall return CHECK CONDITION status for that MODE SELECT command. The device server shall set the sense key to ILLEGAL REQUEST and set the additional sense code to ILLEGAL FIELD IN PARAMETER LIST.

The MAXIMUM BURST SIZE field indicates the maximum size of an RDMA Read or RDMA Write operation that the device server shall perform. This value is expressed in increments of 512 bytes (e.g., a value of one means 512 bytes, two means 1024 bytes, etc.). The device server may round this value down as defined in SPC-2. A value of zero indicates that the maximum transfer size is limited only to that of the underlying interconnect. The value zero shall be implemented by all SRP devices. The application client and device server may use the value of this parameter to adjust internal maximum buffering requirements. A router between an SRP device and another protocol device (e.g. FCP) may intercept and adjust this value to reflect its own maximum buffering capabilities.

The ENABLE MODIFY DATA POINTERS (EMDP) bit indicates whether the SRP target port may use the random buffer access capability to order RDMA requests for a single SCSI command. If the EMDP bit is set to zero, the SRP target port shall generate RDMA requests with continuously increasing addresses for a single SCSI command. If the EMDP bit is set to one, the SRP target port may issue RDMA requests for a single SCSI command in any order. The EMDP function shall be implemented by SRP devices.

7.2.2 Invalid fields

The BUFFER FULL RATIO field, BUFFER EMPTY RATIO field, BUS INACTIVITY LIMIT field, PHYSICAL DISCONNECT TIME LIMIT field, CONNECT TIME LIMIT, FAIR ARBITRATION field, DISCONNECT IMMEDIATE (DIMM) bit, DATA TRANSFER DISCONNECT CONTROL (DTDC) field, and FIRST BURST SIZE field shall be set to zero by SRP initiator and SRP target ports.

7.3 Protocol specific LUN page

The Protocol Specific LUN page shall not be implemented by SRP target ports.

7.4 Protocol specific port page

The Protocol Specific Port page shall not be implemented by SRP target ports.

Annex A

(normative)

SRP interface protocol and services

A.1 Service interface protocol

This standard describes a SCSI device's behavior in terms of functional levels, service interfaces between levels and peer-to-peer protocols. For a full description of the model used in this standard see SAM-2. Figure A.1 shows the model as it appears from the point of view of this standard.

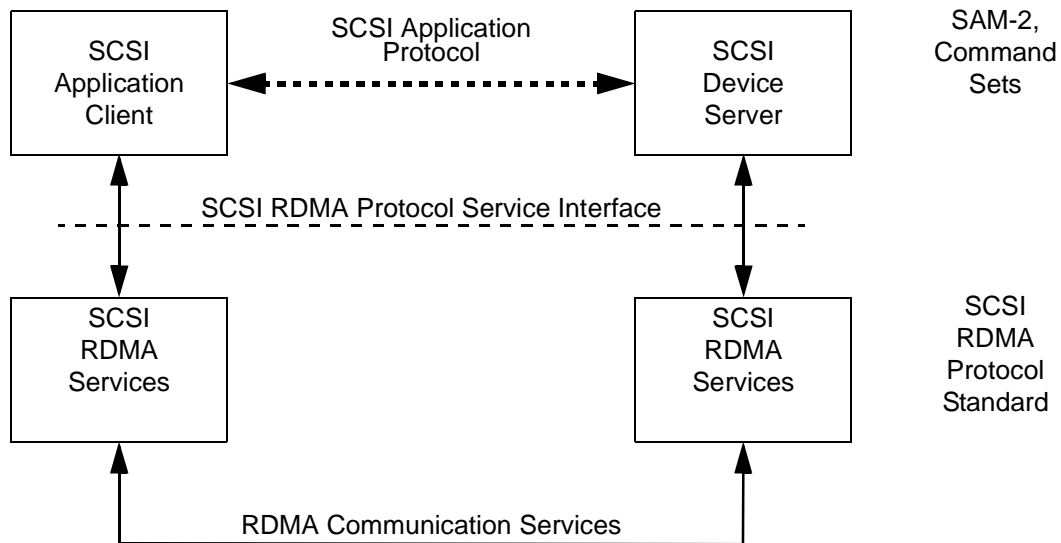


Figure A.1 - SRP reference model

Services between service levels are either four-step confirmed services or two-step confirmed services. A four-step confirmed service consists of a service request, indication, response, and confirmation. A two-step confirmed service consists of a service request and confirmation.

Figure A.2 shows the service and protocol interactions for a four-step confirmed service.

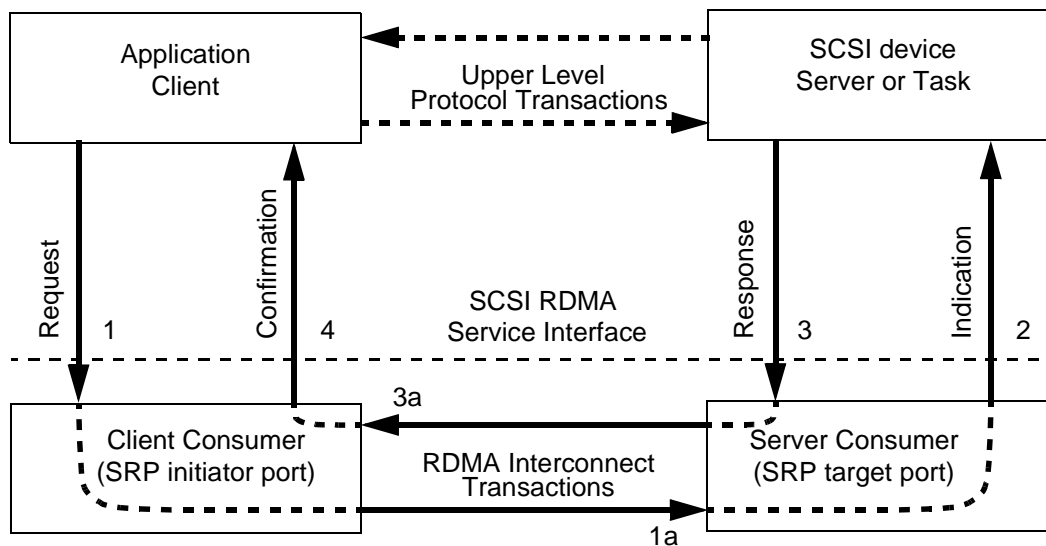


Figure A.2 - Model for a four-step confirmed service

The SCSI RDMA four-step confirmed service protocol consists of the following interactions:

1. A request to the client consumer to invoke a service;
2. An indication from the server consumer notifying the SCSI device server or task manager of an event;
3. A response from the SCSI device server or task manager in reply to an indication;
4. A confirmation from the client consumer upon service completion.

Only application clients shall request a four-step confirmed service be invoked.

Figure A.3 shows the service and protocol interactions for a two-step confirmed service.

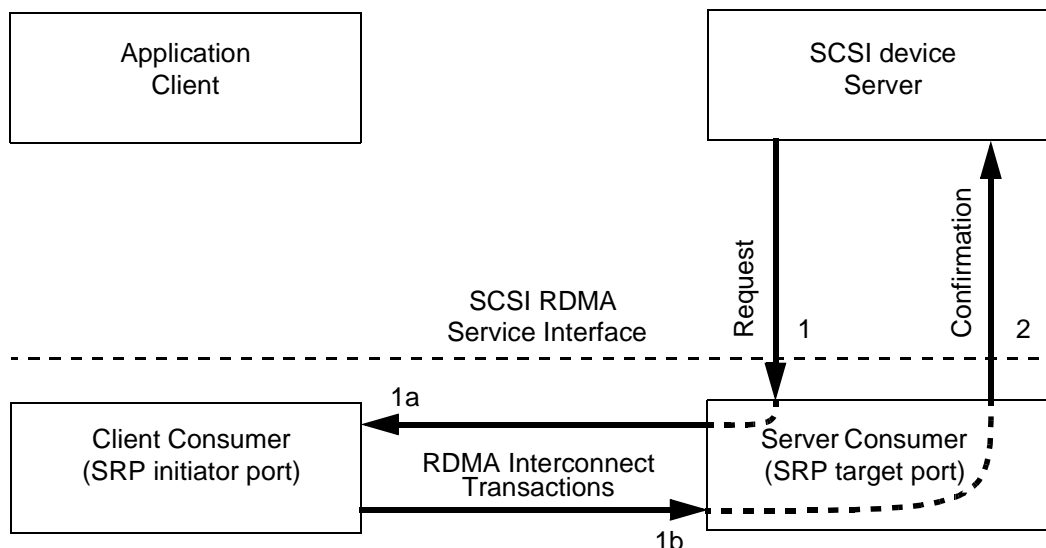


Figure A.3 - Model for a two-step confirmed service

The SCSI RDMA two-step confirmed service interface consists of the following interactions:

1. A request to the server consumer to invoke a service;
2. A confirmation from the server consumer upon service completion.

Only SCSI device servers shall request a two-step confirmed service be invoked.

A.2 SRP services

SRP provides services to enable an application client to request and manage tasks (see SAM-2) and to enable a device server to receive commands and move data to and from an application client. The SRP services are described in terms of the services the SRP initiator port and SRP target port provide.

A.3 SAM-2 object mapping

See table A.1 for the SRP objects corresponding to SAM-2 identifiers and names.

Table A.1 - SAM-2 object mapping

SAM-2 object	SRP object
initiator port identifier	initiator port identifier
initiator port name	
target port identifier	target port identifier
target port name	

A.4 Procedure objects

See table A.2 for a list of the procedure objects used within this standard, the name of the standard where the objects are defined, the standard where the binary contents of the objects are defined, and the routing of the objects. The routing shows:

- a) the source of the object
- b) the final destination of the object, and
- c) the routing of the object.

Table A.2 - Procedure objects

Procedure object	Standard where object format is defined	Object routing
application client buffer offset	SAM-2	DS → targ → init
data-out buffer size	SAM-2	AC → init
data-in buffer size	SAM-2	AC → init
command descriptor block	SAM-2/cmd ^a	AC → init → targ → DS
data-in buffer	cmd ^b	DS → targ → init → AC
data-out buffer	cmd ^b	AC → init → targ → DS
device server buffer	cmd ^b	DS → targ → init
I_T_L_x nexus	this standard	AC → init → targ → DS or AC → init → targ → TM or DS → targ → init
request byte count	SAM-2	DS → targ
service response	this standard ^c	DS → targ → init → AC or targ → DS
autosense request	SAM-2	AC → init → targ → DS
sense data	SPC-2	DS → targ → init → AC
status	SAM-2	DS → targ → init → AC
task attribute	this standard	AC → init → targ → DS
Key: AC=application client, cmd=SCSI command standards, DS=device server, init=SRP initiator port, SAM-2=SAM-2, TM=task manager, targ=SRP target port		
^a The portions not defined in SAM-2 are defined in the SCSI command standards (e.g., SPC-2).		
^b Parameter lists are defined within one of the SCSI command standards (e.g., SPC-2). SCSI standards do not define non-parameter list information.		
^c The SERVICE DELIVERY OR TARGET FAILURE value of the service response is not defined in SCSI.		

A.5 Application client SCSI command services

A.5.1 Application client SCSI command services overview

The SCSI command services shall be requested by the application client using a procedure call defined as:

Execute Command (IN (I_T_L_x nexus, command descriptor block, [task attribute], [data-in buffer size], [data-out buffer], [data-out buffer size], autosense request), OUT ([data-in buffer], [sense data], status, service response))

A.5.2 Send SCSI command service

The send SCSI command service is a four-step confirmed service (see figure A.2) that provides the means to transfer a command data block to a device server.

Processing the execute command procedure call for a send SCSI command service shall be composed of the four-step confirmed service shown in table A.3.

Table A.3 - Processing of execute command procedure call for a send SCSI command service

Step (step number) ^a	Source to Destination	Protocol service name	SCSI Protocol Service Interface procedure calls
request (1)	application client to client consumer	send SCSI command request	Send SCSI command (IN (I_T_L_x nexus, command descriptor block, [task attribute], [data-in buffer size], [data-out buffer], [data-out buffer size], autosense request))
information unit transfer (1a)	client consumer to server consumer	SRP_CMD request or SRP_TSK_MGMT request	See 6.7 and 6.8
indication (2)	server consumer to device server	send SCSI command indication	SCSI command received (IN (I_T_L_x nexus, command descriptor block, [task attribute], autosense request))
If the send SCSI command requires a data transfer see A.6.2 for data-out delivery services and A.6.3 for data-in delivery services			
response (3)	device server to server consumer	send SCSI command response	Send command complete (IN (I_T_L_x nexus, [sense data], status, service response))
information unit transfer (3a)	server consumer to client consumer	SRP_RSP response	See 6.9
confirmation (4)	client consumer to application client	send SCSI command confirmation	Command complete received (IN (I_T_L_x nexus, [data-in buffer], [sense data], status, service response))
^a See figure A.2 for step number			

A.6 Device server SCSI command services

A.6.1 Device server SCSI command services overview

The SCSI data buffer movement services shall be requested from the device server using a procedure call defined as:

Move data buffer (IN (I_T_L_x nexus, device server buffer, application client buffer offset, request byte count)).

Either data-in delivery, data-out delivery, both data-in and data-out delivery, or neither data delivery may be used while processing one command. If both are used, the device server shall combine the data-in and data-out service responses into one service response.

A.6.2 Data-out delivery service

The data-out delivery service is a two-step confirmed service (see figure A.3) that provides the means to transfer a parameter list or data from an SRP initiator port to a device server.

Processing the execute command procedure call for a data-out delivery service shall be composed of the two-step confirmed service shown in table A.4.

Table A.4 - Processing of execute command procedure call for a data-out delivery service

Step (step number) ^a	Source/Destination	Protocol service name	SCSI Protocol Service Interface procedure call
request (1)	device server to server consumer	data-out delivery request	Receive data-out (IN (I_T_L_x nexus, application client buffer offset, request byte count, device server buffer))
data-out transfer (1a and 1b)	server consumer to client consumer ^b	RDMA data-out transfer	See 4.4.3.
confirmation (2)	server consumer to device server	data-out delivery confirmation	Data-out received (IN (I_T_L_x nexus))
^a See figure A.3 for step number			
^b RDMA transfers are typically performed by hardware without the intervention of the client consumer.			

A.6.3 Data-in delivery service

The data-in delivery service is a two-step confirmed service (see figure A.3) that provides the means to transfer a parameter list or data from a device server to an SRP initiator port.

Processing the execute command procedure call for a data-in delivery service shall be composed of the two-step confirmed service shown in table A.5.

Table A.5 - Processing of execute command procedure call for a data-in delivery service

Step (step number) ^a	Source to Destination	Protocol service name	SCSI Protocol Service Interface procedure call
request (1)	device server to server consumer	data-in delivery request	Send data-in (IN (I_T_L_x nexus, device server buffer, application client buffer offset, request byte count))
data-in transfer (1a and 1b)	server consumer to client consumer ^b	RDMA data-in transfer	See 4.4.
confirmation (2)	server consumer to device server	data-in delivery confirmation	Data-In delivered (IN (I_T_L_x nexus))
^a See figure A.3 for step number.			
^b RDMA transfers are typically performed by hardware without the intervention of the client consumer.			

A.7 Task management services

A.7.1 Task management functions overview

The task management services shall be requested from the application client using a procedure call defined as:

Function name (IN (nexus), service response)

A.7.2 Task management functions

This standard handles task management functions as a four-step confirmed service that provides the means to transfer task management functions to a task manager.

The task management functions are defined in the SAM-2. This standard defines the actions taken by the SRP services to carry out the requested task management functions.

A.7.3 ABORT TASK

The SRP services request the SRP initiator port issue an SRP_TSK_MGMT request (see 6.7) with a TASK MANAGEMENT FLAGS field set to indicate an ABORT TASK function to be sent to the selected SCSI device.

A.7.4 ABORT TASK SET

The SRP services request the SRP initiator port issue an SRP_TSK_MGMT request (see 6.7) with a TASK MANAGEMENT FLAGS field set to indicate an ABORT TASK SET function to be sent to the selected SCSI device.

A.7.5 CLEAR ACA

The SRP services request the SRP initiator port issue an SRP_TSK_MGMT request (see 6.7) with a TASK MANAGEMENT FLAGS field set to indicate a CLEAR ACA function to be sent to the selected SCSI device.

A.7.6 CLEAR TASK SET

The SRP services request the SRP initiator port issue an SRP_TSK_MGMT request (see 6.7) with a TASK MANAGEMENT FLAGS field set to indicate a CLEAR TASK SET function to be sent to the selected SCSI device.

A.7.7 LOGICAL UNIT RESET

The SRP services request the SRP initiator port issue an SRP_TSK_MGMT request (see 6.7) with a TASK MANAGEMENT FLAGS field set to indicate a LOGICAL UNIT RESET function to be sent to the selected SCSI device.

A.7.8 TARGET RESET

This protocol does not support use of the TARGET RESET task management function.

A.7.9 WAKEUP

This protocol does not support use of the WAKEUP task management function.

Working Draft

Annex B

(normative)

SRP for the InfiniBand™ Architecture

B.1 Overview

This annex specifies requirements for mapping SRP onto the InfiniBand™ Architecture, a transport that implements a superset of the RDMA communication service (see clause 4). See InfiniBand™ Architecture Specification Volume 1 Release 1.0.a (IBAS) for a description of the InfiniBand™ Architecture.

B.2 Normative references

InfiniBand™ Architecture Specification Volume 1 Release 1.0.a, Infiniband Trade Association (www.infinibandta.org).

IETF RFC 2373, IP Version 6 Addressing Architecture. R. Hinden and S. Deering. Internet Engineering Task Force (www.ietf.org).

B.3 Definitions and abbreviations

B.3.1 Introduction to definitions and abbreviations

The definitions in B.3.2 and the abbreviations in B.3.3 are incomplete without reference to IBAS.

B.3.2 Definitions

B.3.2.1 IB channel adapter: A device that terminates an InfiniBand™ Architecture link and processes transport-level functions.

B.3.2.2 IB channel adapter GUID: An IB GUID that uniquely identifies an IB channel adapter.

B.3.2.3 IB communication manager: The software, hardware, or combination of the two that supports the InfiniBand™ Architecture communication management mechanisms and protocols.

B.3.2.4 IB consumer: An object that communicates with other IB consumers using the InfiniBand™ Architecture.

B.3.2.5 IB GID: A 128-bit value that conforms to the IPv6 address format.

B.3.2.6 IB GUID: A globally unique value that identifies an InfiniBand™ Architecture device or component.

B.3.2.7 IB General Service Interface: An interface providing management services other than IB subnet management.

B.3.2.8 IB I/O controller: The part of an IB I/O unit that provides I/O services.

B.3.2.9 IB I/O controller GUID: An IB GUID that uniquely identifies an IB I/O controller. This value is present as the GUID field of the IOControllerProfile attribute. (See Table B.7)

B.3.2.10 IB I/O unit: One or more IB I/O controllers attached to the IB fabric through a single IB channel adapter.

B.3.2.11 IB LID: A port address used for directing IB packets within an IB subnet.

B.3.2.12 IB MAD: An IB packet used to manage an InfiniBand™ Architecture network.

B.3.2.13 IB packet: The indivisible unit of InfiniBand™ Architecture data transfer and routing, consisting of one or more headers, a packet payload, and one or two CRCs.

B.3.2.14 IB port: A location on an IB channel adapter, switch, or router to which a link connects.

B.3.2.15 IB port GUID: An IB GUID that uniquely identifies an IB port.

Working Draft

B.3.2.16 IB Queue Pair: An interface used for communication, consisting of a Send work queue and a Receive work queue.

B.3.2.17 IB service ID: A value that allows an IB communication manager to associate an incoming connection request with the entity providing the service.

B.3.2.18 IB subnet: A set of IB ports connected via IB switches that have a common IB subnet ID and are managed by a common IB subnet manager.

B.3.2.19 IB subnet manager: Entity that configures and controls an IB subnet. See InfiniBand™ Architecture Specification Volume 1 Release 1.0.a.

B.3.2.20 IPv6 address: A 128-bit address constructed in accordance with IETF RFC 2373 for Internet Protocol version 6. See IETF RFC 2373.

B.3.3 Abbreviations

CM:Ready To Use IB communication manager Ready to Use message

CM:Reject IB communication manager Reject message

CM:Request IB communication manager Request message

CM:Response IB communication manager Response message

CRC Cyclic Redundancy Check

GID IB Global ID

GUID Globally unique identifier

IB InfiniBand™ Architecture

IBAS *InfiniBand™* Architecture Specification Volume 1 Release 1.0.a, Infiniband Trade Association (www.infinibandta.org)

IOC IB IO Controller

IPv6 Internet Protocol version 6

LID IB Local ID

MAD IB Management datagram

QP IB Queue pair

B.4 InfiniBand™ Architecture overview

InfiniBand™ Architecture devices contain IB consumers and one or more IB channel adapters. Each IB channel adapter contains one or more IB ports. Associated with each IB channel adapter are IB QPs that interface between IB consumers and the IB channel adapter. Figure B.1 shows an example InfiniBand™ Architecture device.

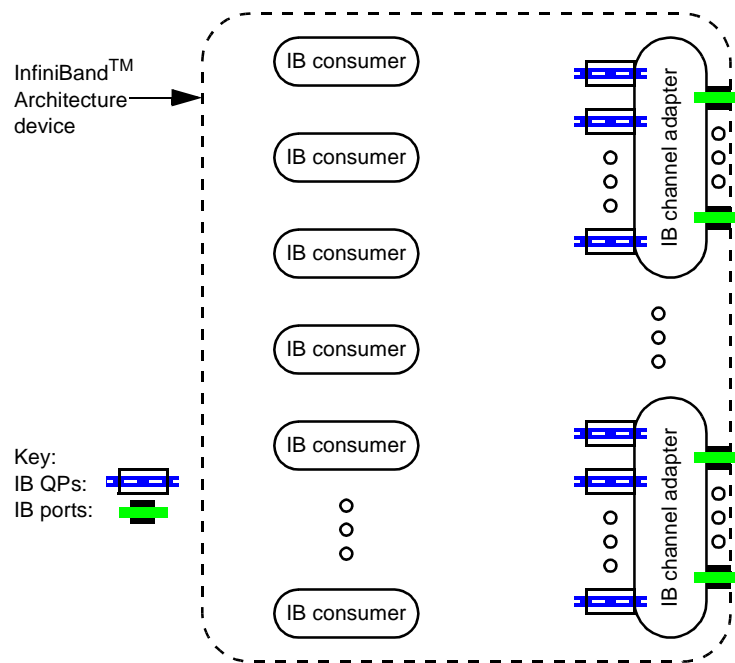


Figure B.1 - InfiniBand™ Architecture device example

An IB I/O unit is an InfiniBand™ Architecture device that contains an IB channel adapter with one or more IB ports, IB QPs, and one or more IB I/O controllers. Figure B.2 shows an example IB I/O unit.

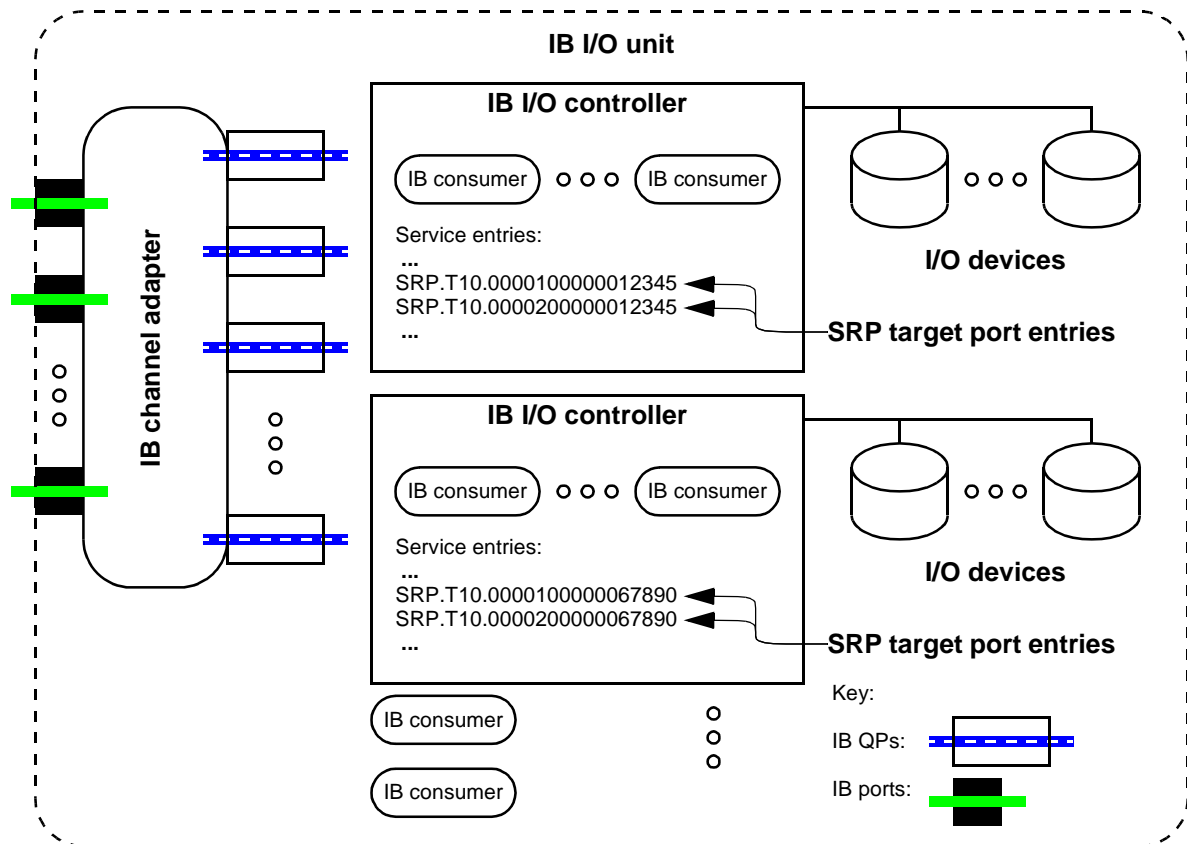


Figure B.2 - IB I/O unit example

Working Draft

Each IB port has a 64-bit globally unique identifier called an IB port GUID. Each IB channel adapter has a IB channel adapter GUID (which is shared by all IB ports on the IB channel adapter). Each IB I/O controller has an IB I/O controller GUID.

The IB subnet manager assigns one or more IB LIDs and one or more IB GIDs to each IB port.

Table B.1 summarizes the InfiniBand™ Architecture names (IB GUIDs) and addresses (IDs) relevant to this protocol.

Table B.1 - InfiniBand™ Architecture names and addresses

Name	Scope of uniqueness	Size	Description
IB port GUID	worldwide	64 bits	Identifies an IB port
IB channel adapter GUID	worldwide	64 bits	Identifies a IB channel adapter
IB I/O controller GUID	worldwide	64 bits	Identifies an IB I/O controller
IB LID	IB subnet	16 bits	Local routing address assigned to each IB port by the IB subnet manager
IB GID	varies ^a	128 bits	Address assigned by the IB subnet manager; (e.g., IB subnet prefix plus the IB port GUID)
^a)Refer to Infiniband™ Architecture Specification Volume 1 Release 1.0.a			

B.5 SCSI architecture mapping

Figure B.3 illustrates how SCSI initiator devices, SRP initiator ports, SRP target ports, and SCSI target devices map to InfiniBand™ Architecture objects.

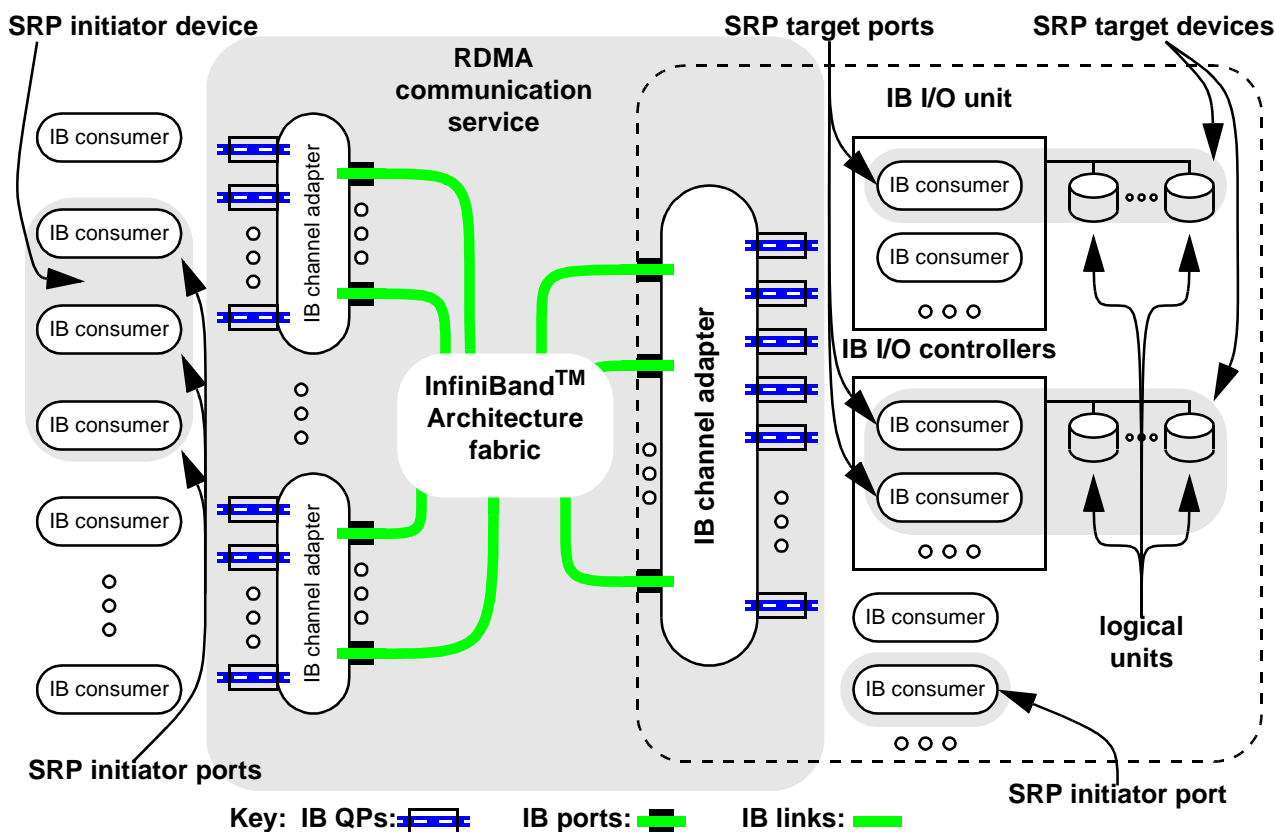


Figure B.3 - SCSI architecture mapping

The RDMA communication service (see clause 4) includes IB queue pairs, IB channel adapters, IB ports, software, and the InfiniBand™ Architecture fabric.

An IB consumer in any InfiniBand™ Architecture device may be an SRP initiator port. An SRP initiator device may consist of one or more IB consumers. The SRP initiator port identifier shall be constructed as shown in table B.2.

Table B.2 - InfiniBand™ Architecture SRP initiator port identifier

Bit Byte	7	6	5	4	3	2	1	0
0	(MSB)							
...	IDENTIFIER EXTENSION							
7								(LSB)
8	(MSB)							
...	GUID (e.g., IB channel adapter GUID)							
15								(LSB)

The IDENTIFIER EXTENSION field shall be chosen by the SRP initiator port to ensure that all SRP initiator port identifiers are unique.

Working Draft

The GUID field should be an IB GUID available to the SRP initiator port.

SRP target ports shall be implemented by IB I/O Controllers in IB I/O units. The IB I/O unit shall include an IB device management agent to provide the IOUnitInfo, IOControllerProfile, and ServiceEntries attributes.

An SRP target port is identified by a ServiceEntries attribute of an IB I/O controller. The SRP target port identifier shall be constructed as shown in table B.3.

Table B.3 - InfiniBand™ Architecture SRP target port identifier

Bit Byte	7	6	5	4	3	2	1	0
0	(MSB)							
...	IDENTIFIER EXTENSION							
7								(LSB)
8	(MSB)							
...	IO CONTROLLER GUID							
15								(LSB)

The IDENTIFIER EXTENSION field shall be the value from the ServiceEntries attribute that identifies the SRP target port (see table B.8).

The IO CONTROLLER GUID field shall be the IB I/O controller GUID of the IB I/O controller providing the SRP target port.

B.6 Communication management

B.6.1 Communication management overview

IB communications managers on each InfiniBand™ Architecture device manage InfiniBand™ Architecture connections using IB MADs transported over the IB General Service Interface. SRP initiator ports and SRP target ports shall use the active/passive (client/server) connection establishment protocol. The processor unit or IB I/O controller containing the SRP target port shall act as the server and the processor unit or IB I/O controller containing the SRP initiator port shall act as the client.

B.6.2 Discovering SRP target ports

To discover the IB service ID of an SRP target port in an IB I/O unit, an SRP initiator port may use this sequence:

1. Retrieve the IOUnitInfo attribute from an IB I/O unit using a DevMgtGet IB MAD to determine the presence and slot number of each IB I/O controller attached to the IB I/O unit.
2. Retrieve the IOControllerProfile attributes from each IB I/O controller, each of which includes a ServiceEntries table.
3. Search the ServiceEntries table for service names matching the rules described in table B.8.

The IB service ID associated with each matching service name may be used in the communication management process to establish InfiniBand™ Architecture connections to IB I/O controllers providing SRP target ports. The SRP target port identifier for each SRP target port is constructed as described in table B.3.

B.6.3 Establishing a connection

To establish an InfiniBand™ Architecture connection, the client places the IB service ID in an IB communication management CM:Request message. The server associates the request with the appropriate SRP target port. The PrivateData field of the CM:Request message shall include an SRP_LOGIN_REQ request (see 6.2).

The SRP target port may choose to refuse the connection based on the SRP_LOGIN_REQ request content by returning a CM:Reject message with the reason code set to Consumer Reject. The PrivateData field of the CM:Reject message shall include an SRP_LOGIN_REJ response (see 6.4).

The SRP target port may choose to redirect the connection to a different endpoint (e.g. another IB port) by returning a CM:Reject message with the reason code set to either PORT AND CM REDIRECTION or PORT REDIRECTION. The SRP initiator port should retry the connection establishment using the new endpoint. See InfiniBand™ Architecture Specification Volume 1 Release 1.0.a.

If the server accepts the connection request and SRP login, the server returns a CM:Response message. The PrivateData field of the CM:Response message shall include an SRP_LOGIN_RSP response (see 6.3). The SRP initiator port may choose to refuse the connection based on the SRP_LOGIN_RSP response content by returning a CM:Reject message with a Reason code set to Consumer Reject. In this case, the PrivateData field of the CM:Reject message is reserved.

If the client accepts the connection reply and the SRP login response, it replies with a CM:Ready To Use message indicating both an InfiniBand™ Architecture and an SRP connection are open. The client may then start using the connection for communication.

B.6.4 Releasing a connection

The SRP initiator port may send an SRP_I_LOGOUT request or the SRP target port may send an SRP_T_LOGOUT request with a SEND operation. The sender shall send a CM Disconnect Request as described in IBAS upon receipt of an InfiniBand™ Architecture transport level acknowledgement to the SRP_I_LOGOUT request or SRP_T_LOGOUT request information unit. The receiver of an SRP_I_LOGOUT request or SRP_T_LOGOUT request information unit shall respond with an InfiniBand™ Architecture transport acknowledgement and may send a CM Disconnect Request as described in IBAS, or may wait to receive a CM Disconnect Request.

B.6.5 Errors

Some errors cause an IB queue pair to enter the Error state, which destroys the connection. The IB communication manager for the queue pair consumer should send a CM Disconnect Request as described in IBAS.

B.6.6 Data-out and data-in operations

An SRP target port shall map a Receive Data-out SCSI protocol service interface procedure call to one or more InfiniBand™ Architecture RDMA READ requests. An SRP target port shall map a Send Data-in SCSI protocol

service interface procedure call to one or more InfiniBand™ Architecture RDMA WRITE requests. Table B.4 specifies the value of the InfiniBand™ Architecture RDMA header fields.

Table B.4 - InfiniBand™ Architecture RDMA header fields

InfiniBand™ Architecture RDMA Extended Transport Header field	Value
Virtual Address	VIRTUAL ADDRESS ^a + application client buffer offset ^b
Remote Key	MEMORY HANDLE ^c
DMA Length	request byte count ^d
^a The contents of the VIRTUAL ADDRESS field in the memory descriptor (see table 1). ^b The application client buffer offset parameter to the receive data-out (see table A.4) or send data-in (see table A.5) SCSI protocol service interface procedure call. ^c The contents of the MEMORY HANDLE field in the memory descriptor (see table 1). ^d The request byte count parameter to the receive data-out (see table A.4) or send data-in (see table A.5) SCSI protocol service interface procedure call.	

B.7 InfiniBand™ Architecture protocol requirements

SRP target ports and SRP initiator ports shall support the Reliable Connection transport service type.

SRP target ports shall implement the device management class of general management services.

SRP initiator ports and SRP target ports shall support the transport functions described in table B.5.

Table B.5 - Transport operation support requirements

Transport functions	SRP initiator port	SRP target port
Send to	Mandatory	Mandatory
Send from	Mandatory	Mandatory
RDMA write to	Mandatory	Not used
RDMA write from	Not used	Mandatory
RDMA read to	Mandatory for data-out commands	Not used
RDMA read from	Not used	Mandatory for data-out commands
RDMA Write with immediate data (to or from)	Not used	Not used
ATOMIC (to or from)	Not used	Not used

IB I/O units containing an IB I/O controller acting as an SRP target port shall report the device management IOUnit attributes defined in Infiniband™ Architecture Specification Volume 1 Release 1.0.a as described in table B.6.

Table B.6 - IOUnit attributes for SRP target ports

Field	SRP requirement
Max Controllers	At least one
Controller List	At least one IB I/O controller acting as an SRP target port shall be present
^a This protocol does not change or override InfiniBand Architecture requirements on the values of fields not listed.	

IB I/O controllers acting as SRP target ports shall report the device management IOControllerProfile attributes defined in Infiniband™ Architecture Specification Volume 1 Release 1.0.a as described in table B.7.

Table B.7 - IOControllerProfile attributes for SRP target ports

Field	SRP requirement
I/O Class	0100h
I/O Subclass	609Eh
Protocol	0108h
Protocol Version	0001h
Service Connections	At least one
Initiators Supported	At least one
Send Message Depth	Reserved
RDMA Read Depth	Maximum IOC-issued RDMA depth ^a
Send Message Size	MAXIMUM INITIATOR TO TARGET IU SIZE ^b
RDMA Transfer Size	Reserved
Controller Operations Capability Mask: 0: ST; Send Messages To IOCs 1: SF; Send Messages From IOCs 2: RT; RDMA Read Requests To IOCs 3: RF; RDMA Read Requests From IOCs 4: WT; RDMA Write Requests To IOCs 5: WF; RDMA Write Requests From IOCs 6: AT; Atomic Operations To IOCs 7: AF; Atomic Operations From IOCs	Shall be set to one. Shall be set to one. No requirement Shall be set to one if an SRP target port supports data-out commands. No requirement otherwise. No requirement Shall be set to one. No requirement No requirement
Controller Services Capability Mask	Reserved ^a
Service Entries	At least one
^a This protocol does not change or override InfiniBand Architecture requirements on the values of fields not listed, or for those marked as having 'no requirement'. ^b The largest number of RDMA Read requests that this IO Controller may have outstanding on one channel. ^c This value shall be no less than the largest value, in bytes, of MAXIMUM INITIATOR TO TARGET IU SIZE that this IO Controller shall return in the SRP_LOGIN_RSP information unit.	

IB I/O controllers providing SRP target ports shall include at least one ServiceName/ServiceID pair in the device management ServiceEntries attribute pair (see IBAS) as described in table B.8.

Table B.8 - ServiceEntries attribute pair for SRP target ports

Field	Length (bits)	SRP requirement
ServiceName_n	320	'SRP.T10:xxxxxxxxxxxxxxxxxx' or 'SRP.T10:xxxxxxxxxxxxxxxxxx:reserved'
ServiceID_n	64	Assigned by the IB I/O controller
<p>^a A service name that identifies an SRP target port shall meet the rules described in this table.</p> <p>^b The string 'SRP.T10' and the colons shall appear exactly as shown (e.g. capital letters only).</p> <p>^c The string 'xxxxxxxxxxxxxxxxxx' in the service name shall be sixteen hexadecimal digits. Only the characters 0 to 9 and A to F (capital letters only) are permitted. If any other character appears the service name shall not be recognized as identifying an SRP target port.</p> <p>^d The string 'xxxxxxxxxxxxxxxxxx' in the service name identifies the 64-bit extension identifier value used to construct the SRP target port identifier (see table B.3)</p> <p>^e The literal string 'reserved' shall either be ignored by SRP initiator ports or treated in accordance with a future revision of this standard.</p> <p>^f If the service name does not completely fill ServiceName_n field (i.e. it is less than 40 bytes), it shall be extended with null characters (i.e., binary zeros).</p>		