

draft-recio-iwarp-rdmap-v1.0	R. Recio	1
	IBM Corporation	2
	P. Culley	3
	Hewlett-Packard Company	4
	D. Garcia	5
	Hewlett-Packard Company	6
	J. Hilland	7
	Hewlett-Packard Company	8
		9
	October 2002	10
		11
		12
	An RDMA Protocol Specification (Version 1.0)	13
		14
1	Status of this Memo	15
		16
	This document is a Release Specification of the RDMA Consortium.	17
	Copies of this document and associated errata may be found at	18
	http://www.rdmaconsortium.org .	19
		20
2	Abstract	21
		22
	This document defines a Remote Direct Memory Access Protocol (RDMA)	23
	that operates over the Direct Data Placement Protocol (DDP	24
	protocol). RDMA provides read and write services directly to	25
	applications and enables data to be transferred directly into ULP	26
	Buffers without intermediate data copies. It also enables a kernel	27
	bypass implementation.	28
		29
		30
		31
		32
		33
		34
		35
		36
		37
		38
		39
		40
		41
		42
		43
		44
		45
		46
		47
		48
		49
		50
		51

Table of Contents

		1
		2
		3
1	Status of this Memo	4
2	Abstract	5
3	Introduction	6
3.1	Architectural Goals	7
3.2	Protocol Overview	8
3.3	RDMAP Layering	9
4	Glossary	10
4.1	General	11
4.2	LLP	12
4.3	Direct Data Placement (DDP)	13
4.4	Remote Direct Memory Access (RDMA)	14
5	ULP and Transport Attributes	15
5.1	Transport Requirements & Assumptions	16
5.2	RDMAP Interactions with the ULP	17
6	Header Format	18
6.1	RDMAP Control and Invalidate STag Field	19
6.2	RDMA Message Definitions	20
6.3	RDMA Write Header	21
6.4	RDMA Read Request Header	22
6.5	RDMA Read Response Header	23
6.6	Send Header and Send with Solicited Event Header	24
6.7	Send with Invalidate Header and Send with SE and Invalidate Header.....	25 26
6.8	Terminate Header	27
7	Data Transfer	28
7.1	RDMA Write Message	29
7.2	RDMA Read Operation	30
7.2.1	RDMA Read Request Message.....	31
7.2.2	RDMA Read Response Message.....	32
7.3	Send Message Type	33
7.4	Terminate Message	34
7.5	Ordering and Completions	35
8	RDMAP Stream Management	36
8.1	Stream Initialization	37
8.2	Stream Teardown	38
8.2.1	RDMAP Abortive Termination.....	39
9	RDMAP Error Management	40
9.1	RDMAP Error Surfacing	41
9.2	Errors Detected at the Remote Peer on Incoming RDMA Messages.....	42
10	Security Considerations	43
10.1	Protocol-specific Security Considerations.....	44
10.2	Using IPSec with RDMAP.....	45
10.3	Other Security Considerations.....	46
11	References	47
11.1	Normative References.....	48
11.2	Informative References.....	49
12	Appendix	50
12.1	DDP Segment Formats for RDMA Messages.....	51

12.1.1	DDP Segment for RDMA Write	50	1
12.1.2	DDP Segment for RDMA Read Request	50	2
12.1.3	DDP Segment for RDMA Read Response	51	3
12.1.4	DDP Segment for Send and Send with Solicited Event	52	4
12.1.5	DDP Segment for Send with Invalidate and Send with SE and Invalidate.....	52	5
12.1.6	DDP Segment for Terminate	53	6
12.2	Ordering and Completion Table.....	53	7
13	Authors Addresses	56	8
14	Acknowledgments	57	9
15	Full Copyright Statement	60	10

Table of Figures

Figure 1	RDMA Layering.....	7	11
Figure 2	Example of MPA, DDP, and RDMA Header Alignment over TCP	8	12
Figure 3	DDP Control, RDMA Control, and Invalidate STag Fields.	21	13
Figure 4	RDMA Usage of DDP Fields.....	22	14
Figure 5	RDMA Message Definitions.....	23	15
Figure 6	RDMA Read Request Header Format.....	24	16
Figure 7	Terminate Header Format.....	26	17
Figure 8	Terminate Control Field.....	27	18
Figure 9	Terminate Control Field Values.....	28	19
Figure 10	Error Type to RDMA Message Mapping.....	30	20
Figure 11	RDMA Write, DDP Segment format.....	50	21
Figure 12	RDMA Read Request, DDP Segment format.....	51	22
Figure 13	RDMA Read Response, DDP Segment format.....	51	23
Figure 14	Send and Send with Solicited Event, DDP Segment format	52	24
Figure 15	Send with Invalidate and Send with SE and Invalidate, DDP Segment.....	52	25
Figure 16	Terminate, DDP Segment format.....	53	26
Figure 17	Operation Ordering.....	55	27

3 Introduction

Today, communications over TCP/IP typically require copy operations, which add latency and consume significant CPU and memory resources. The Remote Direct Memory Access Protocol (RDMA) enables removal of data copy operations and enables reduction in latencies by allowing a local application to read or write data on a remote computer's memory with minimal demands on memory bus bandwidth and CPU processing overhead, while preserving memory protection semantics.

RDMA is layered on top of Direct Data Placement (DDP) and uses the two Buffer Models available from DDP [DDP].

3.1 Architectural Goals

RDMA has been designed with the following high-level architectural goals:

- * Provide a data transfer operation that allows a Local Peer to transfer up to $2^{32} - 1$ octets directly into a previously advertised buffer (i.e. Tagged buffer) located at a Remote Peer without requiring a copy operation. This is referred to as the RDMA Write data transfer operation.
- * Provide a data transfer operation that allows a Local Peer to retrieve up to $2^{32} - 1$ octets directly from a previously advertised buffer (i.e. Tagged buffer) located at a Remote Peer without requiring a copy operation. This is referred to as the RDMA Read data transfer operation.
- * Provide a data transfer operation that allows a Local Peer to send up to $2^{32} - 1$ octets directly into a buffer located at a Remote Peer that has not been explicitly advertised. This is referred to as the Send (Send with Invalidate, Send with Solicited Event, and Send with Solicited Event and Invalidate) data transfer operation.
- * Enable the local ULP to use the Send Operation Type (includes Send, Send with Invalidate, Send with Solicited Event, and Send with Solicited Event and Invalidate) to signal to the remote ULP the Completion of all previous Messages initiated by the local ULP.
- * Provide for all Operations on a single RDMA Stream to be reliably transmitted in the order that they were submitted.
- * Provide RDMA capabilities independently for each Stream when the LLP supports multiple data Streams within an LLP connection.

3.2 Protocol Overview

RDMA provides seven data transfer operations. Except for the RDMA Read operation, each operation generates exactly one RDMA Message. Following is a brief overview of the RDMA Operations and RDMA Messages:

1. Send - A Send operation uses a Send Message to transfer data from the Data Source into a buffer that has not been explicitly Advertised by the Data Sink. The Send Message uses the DDP Untagged Buffer Model to transfer the ULP Message into the Data Sink's Untagged Buffer.
2. Send with Invalidate - A Send with Invalidate operation uses a Send with Invalidate Message to transfer data from the Data Source into a buffer that has not been explicitly Advertised by the Data Sink. The Send with Invalidate Message includes all functionality of the Send Message, with one addition: an STag field is included in the Send With Invalidate Message and after the message has been Placed and Delivered at the Data Sink the remote peer's buffer identified by the STag can no longer be accessed remotely until the remote peer's ULP re-enables access and Advertises the buffer.
3. Send with Solicited Event (Send with SE) - A Send with Solicited Event operation uses a Send with Solicited Event Message to transfer data from the Data Source into an Untagged Buffer at the Data Sink. The Send with Solicited Event Message is similar to the Send Message, with one addition: when the Send with Solicited Event Message has been Placed and Delivered, an Event may be generated at the recipient, if the recipient is configured to generate such an Event.
4. Send with Solicited Event and Invalidate (Send with SE and Invalidate) - A Send with Solicited Event and Invalidate operation uses a Send with Solicited Event and Invalidate Message to transfer data from the Data Source into a buffer that has not been explicitly Advertised by the Data Sink. The Send with Solicited Event and Invalidate Message is similar to the Send with Invalidate Message, with one addition: when the Send with Solicited Event and Invalidate Message has been Placed and Delivered, an Event may be generated at the recipient, if the recipient is configured to generate such an Event.
5. Remote Direct Memory Access Write - An RDMA Write operation uses an RDMA Write Message to transfer data from the Data Source to a previously advertised buffer at the Data Sink.

The ULP at the Remote Peer, which in this case is the Data Sink, enables the Data Sink Tagged Buffer for access and Advertises

the buffer's size (length), location (Tagged Offset), and Steering Tag (STag) to the Data Source through a ULP specific mechanism. The ULP at the Local Peer, which in this case is the Data Source, initiates the RDMA Write operation. The RDMA Write Message uses the DDP Tagged Buffer Model to transfer the ULP Message into the Data Sink's Tagged Buffer. Note: the STag associated with the Tagged Buffer remains valid until the ULP at the Remote Peer invalidates it or the ULP at the Local Peer invalidates it through a Send with Invalidate or Send with Solicited Event and Invalidate.

6. Remote Direct Memory Access Read - The RDMA Read operation transfers data to a Tagged Buffer at the Local Peer, which in this case is the Data Sink, from a Tagged Buffer at the Remote Peer, which in this case is the Data Source. The ULP at the Data Source enables the Data Source Tagged Buffer for access and Advertises the buffer's size (length), location (Tagged Offset), and Steering Tag (STag) to the Data Sink through a ULP specific mechanism. The ULP at the Data Sink enables the Data Sink Tagged Buffer for access and initiates the RDMA Read operation. The RDMA Read operation consists of a single RDMA Read Request Message and a single RDMA Read Response Message, and the latter may be segmented into multiple DDP Segments.

The RDMA Read Request Message uses the DDP Untagged Buffer Model to Deliver the STag, starting Tagged Offset and length for both the Data Source and Data Sink Tagged Buffers to the remote peer's RDMA Read Request Queue.

The RDMA Read Response Message uses the DDP Tagged Buffer Model to Deliver the Data Source's Tagged Buffer to the Data Sink, without any involvement from the ULP at the Data Source.

Note: the Data Source STag associated with the Tagged Buffer remains valid until the ULP at the Data Source invalidates it or the ULP at the Data Sink invalidates it through a Send with Invalidate or Send with Solicited Event and Invalidate. The Data Sink STag associated with the Tagged Buffer remains valid until the ULP at the Data Sink invalidates it.

7. Terminate - A Terminate operation uses a Terminate Message to transfer to the Remote Peer information associated with an error that occurred at the Local Peer. The Terminate Message uses the DDP Untagged Buffer Model to transfer the Message into the Data Sink's Untagged Buffer.

3.3 RDMAP Layering

RDMAP is dependent on DDP, subject to the requirements defined in section 5 ULP and Transport Attributes

Transport Requirements & Assumptions. Figure 1 RDMAP Layering depicts the relationship between Upper Layer Protocols (ULPs), RDMAP, DDP protocol, the framing layer, and the transport For LLP protocol definitions of each LLP, see [MPA], [TCP], and [SCTP].

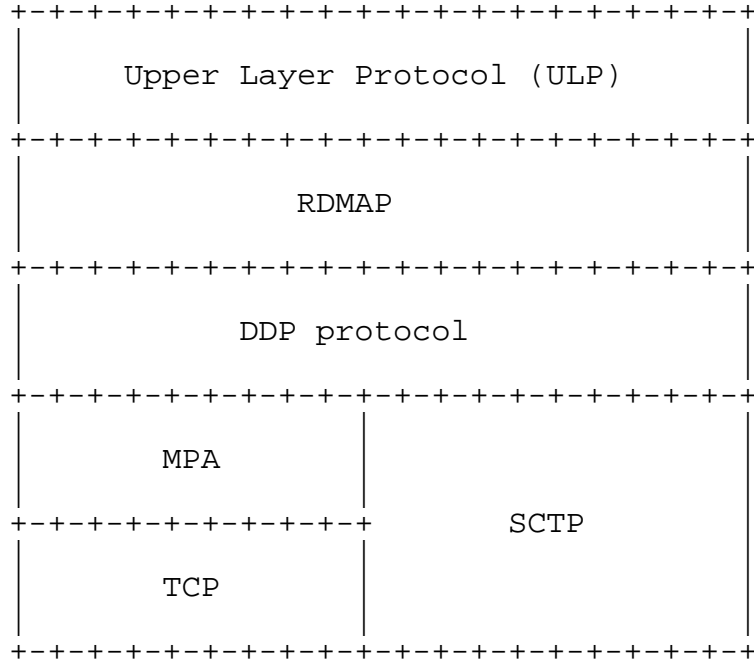


Figure 1 RDMAP Layering

If RDMAP is layered over DDP/MPA/TCP, then the respective headers and ULP Payload are arranged as follows (Note: For clarity, MPA header and CRC fields are included but MPA markers are not shown):

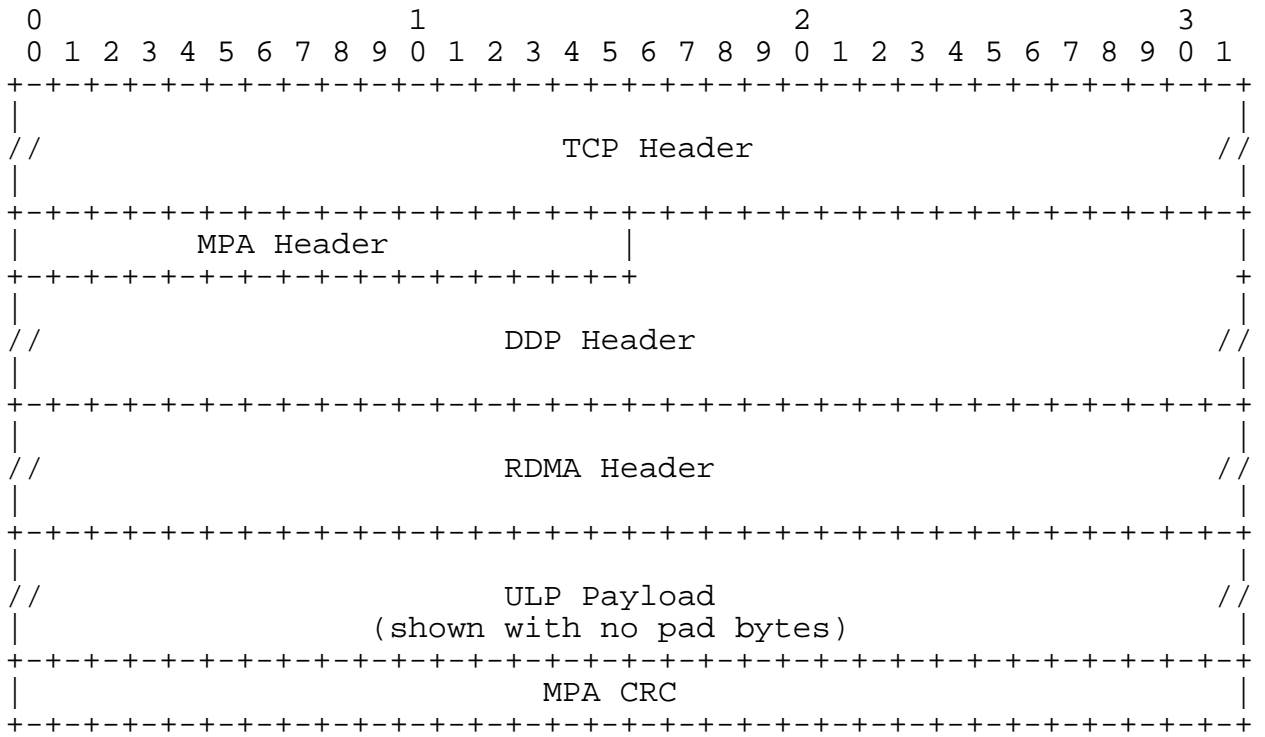


Figure 2 Example of MPA, DDP, and RDMAP Header Alignment over TCP

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

4 Glossary

4.1 General

- Advertisement (Advertised, Advertise, Advertisements, Advertises) - the act of informing a Remote Peer that a local RDMA Buffer is available to it. A Node makes available an RDMA Buffer for incoming RDMA Read or RDMA Write access by informing its RDMA/DDP peer of the Tagged Buffer identifiers (STag, base address, and buffer length). This advertisement of Tagged Buffer information is not defined by RDMA/DDP and is left to the ULP. A typical method would be for the Local Peer to embed the Tagged Buffer's Steering Tag, base address, and length in a Send Message destined for the Remote Peer.
- Data Sink - The peer receiving a data payload. Note that the Data Sink can be required to both send and receive RDMA/DDP Messages to transfer a data payload.
- Data Source - The peer sending a data payload. Note that the Data Source can be required to both send and receive RDMA/DDP Messages to transfer a data payload.
- Data Delivery (Delivery, Delivered, Delivers) - Delivery is defined as the process of informing the ULP or consumer that a particular Message is available for use. This is specifically different from "Placement", which may generally occur in any order, while the order of "Delivery" is strictly defined. See "Data Placement".
- Fabric - The collection of links, switches, and routers that connect a set of Nodes with RDMA/DDP protocol implementations.
- Fence (Fenced, Fences) - To block the current RDMA Operation from executing until prior RDMA Operations have Completed.
- iWARP - A suite of wire protocols comprised of RDMAP, DDP, and MPA. The iWARP protocol suite may be layered above TCP, SCTP, or other transport protocols.
- Local Peer - The RDMA/DDP protocol implementation on the local end of the connection. Used to refer to the local entity when describing a protocol exchange or other interaction between two Nodes.
- Node - A computing device attached to one or more links of a Fabric (network). A Node in this context does not refer to a specific application or protocol instantiation running on the computer. A Node may consist of one or more RNICs installed in a host computer.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

Remote Peer - The RDMA/DDP protocol implementation on the opposite end of the connection. Used to refer to the remote entity when describing protocol exchanges or other interactions between two Nodes.

RNIC - RDMA Network Interface Controller. In this context, this would be a network I/O adapter or embedded controller with iWARP and verbs functionality.

RNIC Interface (RI) - The presentation of the RNIC to the verbs Consumer as implemented through the combination of the RNIC and the RNIC driver.

ULP - Upper Layer Protocol. The protocol layer above the protocol layer currently being referenced. The ULP for RDMA/DDP is expected to be an OS, Application, adaptation layer, or proprietary device. The RDMA/DDP documents do not specify a ULP - they provide a set of semantics that allow a ULP to be designed to utilize RDMA/DDP.

ULP Payload - The ULP data that is contained within a single protocol segment or packet (e.g. a DDP Segment).

Verbs - An abstract description of the functionality of a RNIC Interface. The OS may expose some or all of this functionality via one or more APIs to applications. The OS will also use some of the functionality to manage the RNIC Interface.

4.2 LLP

LLP - Lower Layer Protocol. The protocol layer beneath the protocol layer currently being referenced. For example, for DDP the LLP is SCTP, MPA, or other transport protocols. For RDMA, the LLP is DDP.

LLP Connection - Corresponds to an LLP transport-level connection between the peer LLP layers on two nodes.

LLP Stream - Corresponds to a single LLP transport-level Stream between the peer LLP layers on two Nodes. One or more LLP Streams may map to a single transport-level LLP connection. For transport protocols that support multiple Streams per connection (e.g. SCTP), a LLP Stream corresponds to one transport-level Stream.

MULPDU - Maximum ULPDU. The current maximum size of the record that is acceptable for DDP to pass to the LLP for transmission.

ULPDU - Upper Layer Protocol Data Unit. The data record defined by the layer above MPA.

4.3 Direct Data Placement (DDP)

- Data Placement (Placement, Placed, Places) - For DDP, this term is specifically used to indicate the process of writing to a data buffer by a DDP implementation. DDP Segments carry Placement information, which may be used by the receiving DDP implementation to perform Data Placement of the DDP Segment ULP Payload. See "Data Delivery".
- DDP Abortive Teardown - The act of closing a DDP Stream without attempting to Complete in-progress and pending DDP Messages.
- DDP Graceful Teardown - The act of closing a DDP Stream such that all in-progress and pending DDP Messages are allowed to Complete successfully.
- DDP Control Field - a fixed 16-bit field in the DDP Header. The DDP Control Field contains an 8-bit field whose contents are reserved for use by the ULP.
- DDP Header - The header present in all DDP segments. The DDP Header contains control and Placement fields that are used to define the final Placement location for the ULP payload carried in a DDP Segment.
- DDP Message - A ULP defined unit of data interchange, which is subdivided into one or more DDP segments. This segmentation may occur for a variety of reasons, including segmentation to respect the maximum segment size of the underlying transport protocol.
- DDP Segment - The smallest unit of data transfer for the DDP protocol. It includes a DDP Header and ULP Payload (if present). A DDP Segment should be sized to fit within the underlying transport protocol MULPDU.
- DDP Stream - a sequence of DDP Messages whose ordering is defined by the LLP. For SCTP, a DDP Stream maps directly to an SCTP Stream. For MPA, a DDP Stream maps directly to a TCP connection and a single DDP Stream is supported. Note that DDP has no ordering guarantees between DDP Streams.
- Direct Data Placement - A mechanism whereby ULP data contained within DDP Segments may be Placed directly into its final destination in memory without processing of the ULP. This may occur even when the DDP Segments arrive out of order. Out of order Placement support may require the Data Sink to implement the LLP and DDP as one functional block.

Direct Data Placement Protocol (DDP) - Also, a wire protocol that supports Direct Data Placement by associating explicit memory buffer placement information with the LLP payload units.

Message Offset (MO) - For the DDP Untagged Buffer Model, specifies the offset, in bytes, from the start of a DDP Message.

Message Sequence Number (MSN) - For the DDP Untagged Buffer Model, specifies a sequence number that is increasing with each DDP Message.

Queue Number (QN) - For the DDP Untagged Buffer Model, identifies a destination Data Sink queue for a DDP Segment.

Steering Tag - An identifier of a Tagged Buffer on a Node, valid as defined within a protocol specification.

STag - Steering Tag

Tagged Buffer - A buffer that is explicitly Advertised to the Remote Peer through exchange of an STag, Tagged Offset, and length.

Tagged Buffer Model - A DDP data transfer model used to transfer Tagged Buffers from the Local Peer to the Remote Peer.

Tagged DDP Message - A DDP Message that targets a Tagged Buffer.

Tagged Offset (TO) - The offset within a Tagged Buffer on a Node.

Untagged Buffer - A buffer that is not explicitly Advertised to the Remote Peer.

Untagged Buffer Model - A DDP data transfer model used to transfer Untagged Buffers from the Local Peer to the Remote Peer.

Untagged DDP Message - A DDP Message that targets an Untagged Buffer.

4.4 Remote Direct Memory Access (RDMA)

Event - An indication provided by the RDMA Layer to the ULP to indicate a Completion or other condition requiring immediate attention.

Invalidate STag - A mechanism used to prevent the Remote Peer from reusing a previous explicitly Advertised STag, until the Local Peer makes it available through a subsequent explicit Advertisement. The STag cannot be accessed remotely until it is explicit Advertised again.

RDMA Completion (Completion, Completed, Complete, Completes) - For RDMA, Completion is defined as the process of informing the ULP that a particular RDMA Operation has performed all functions specified for the RDMA Operations, including Placement and Delivery. The Completion semantic of each RDMA Operation is distinctly defined.

RDMA Message - A data transfer mechanism used to fulfill an RDMA Operation.

RDMA Operation - A sequence of RDMA Messages, including control Messages, to transfer data from a Data Source to a Data Sink. The following RDMA Operations are defined - RDMA Writes, RDMA Read, Send, Send with Invalidate, Send with Solicited Event, Send with Solicited Event and Invalidate, and Terminate.

RDMA Protocol (RDMAP) - A wire protocol that supports RDMA Operations to transfer ULP data between a Local Peer and the Remote Peer.

RDMAP Abortive Termination (Termination, Terminated, Terminate, Terminates) - The act of closing an RDMAP Stream without attempting to Complete in-progress and pending RDMA Operations.

RDMAP Graceful Termination - The act of closing an RDMAP Stream such that all in-progress and pending RDMA Operations are allowed to Complete successfully.

RDMA Read - An RDMA Operation used by the Data Sink to transfer the contents of a source RDMA buffer from the Remote Peer to the Local Peer. An RDMA Read operation consists of a single RDMA Read Request Message and a single RDMA Read Response Message.

RDMA Read Request - An RDMA Message used by the Data Sink to request the Data Source to transfer the contents of an RDMA buffer. The RDMA Read Request Message describes both the Data Source and Data Sink RDMA buffers.

RDMA Read Request Queue - The queue used for processing RDMA Read Requests. The RDMA Read Request Queue has a DDP Queue Number of 1.

RDMA Read Response - An RDMA Message used by the Data Source to transfer the contents of an RDMA buffer to the Data Sink, in response to an RDMA Read Request. The RDMA Read Response Message only describes the data sink RDMA buffer.

RDMAP Stream - An association between a pair of RDMAP implementations, possibly on different Nodes, which transfer ULP data using RDMA Operations. There may be multiple RDMAP Streams

on a single Node. An RDMAP Stream maps directly to a single DDP Stream. 1
2
3
4

RDMA Write - An RDMA Operation that transfers the contents of a 5
source RDMA Buffer from the Local Peer to a destination RDMA 6
Buffer at the Remote Peer using RDMA. The RDMA Write Message 7
only describes the Data Sink RDMA buffer. 8
9

Remote Direct Memory Access (RDMA) - A method of accessing memory on 10
a remote system in which the local system specifies the remote 11
location of the data to be transferred. Employing a RNIC in the 12
remote system allows the access to take place without 13
interrupting the processing of the CPU(s) on the system. 14
15

Send - An RDMA Operation that transfers the contents of a ULP Buffer 16
from the Local Peer to an Untagged Buffer at the Remote Peer. 17
18

Send Message Type - A Send Message, Send with Invalidate Message, 19
Send with Solicited Event Message, or Send with Solicited Event 20
and Invalidate Message. 21
22

Send Operation Type - A Send Operation, Send with Invalidate 23
Operation, Send with Solicited Event Operation, or Send with 24
Solicited Event and Invalidate Operation. 25
26

Solicited Event (SE) - A facility by which an RDMA Operation sender 27
may cause an Event to be generated at the recipient, if the 28
recipient is configured to generate such an Event, when a Send 29
with Solicited Event or Send with Solicited Event and Invalidate 30
Message is received. Note: The Local Peer's ULP can use the 31
Solicited Event mechanism to ensure that Messages designated as 32
important to the ULP are handled in an expeditious manner by the 33
Remote Peer's ULP. The ULP at the Local Peer can indicate a given 34
Send Message Type is important by using the Send with Solicited 35
Event Message or Send with Solicited Event and Invalidate 36
Message. The ULP at the Remote Peer can choose to only be 37
notified when valid Send with Solicited Event Messages and/or 38
Send with Solicited Event and Invalidate Messages arrive and 39
handle other valid incoming Send Messages or Send with Invalidate 40
Messages at its leisure. 41
42

Terminate - An RDMA Message used by a Node to pass an error 43
indication to the peer Node on an RDMAP Stream. This operation 44
is for RDMAP use only. 45
46

ULP Buffer - A buffer owned above the RDMAP Layer and advertised to 47
the RDMAP Layer either as a Tagged Buffer or an Untagged ULP 48
Buffer. 49
50
51

ULP Message - The ULP data that is handed to a specific protocol layer for transmission. Data boundaries are preserved as they are transmitted through iWARP.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

5 ULP and Transport Attributes

5.1 Transport Requirements & Assumptions

RDMAP MUST be layered on top of the Direct Data Placement Protocol [DDP].

RDMAP requires the following DDP support:

- * RDMAP uses three queues for Untagged Buffers:
 - * Queue Number 0 (used by RDMAP for Send, Send with Invalidate, Send with Solicited Event, and Send with Solicited Event and Invalidate operations).
 - * Queue Number 1 (used by RDMAP for RDMA Read operations).
 - * Queue Number 2 (used by RDMAP for Terminate operations).
- * DDP maps a single RDMA Message to a single DDP Message.
- * DDP uses the STag and Tagged Offset provided by the RDMAP for Tagged Buffer Messages (i.e. RDMA Write and RDMA Read Response).
- * When the DDP layer Delivers an Untagged DDP Message to the RDMAP layer, DDP provides the length of the DDP Message. This ensures that RDMAP does not have to carry a length field in its header.
- * When the RDMAP layer provides an RDMA Message to the DDP Layer, DDP must insert the RsvdULP field value provided by the RDMAP Layer into the associated DDP Message.
- * When the DDP layer Delivers a DDP Message to the RDMAP layer, DDP provides the RsvdULP field.
- * The RsvdULP field must be 1 octet for DDP Tagged Messages and 5 octets for DDP Untagged Messages.
- * DDP propagates to RDMAP all operation or protection errors (used by RDMAP Terminate) and, when appropriate, the DDP Header fields of the DDP Segment that encountered the error.
- * If an RDMA Operation is aborted by DDP or a lower layer, the contents of the Data Sink buffers associated with the operation are considered indeterminate.
- * DDP in conjunction with the lower layers provide reliable, in-order Delivery.

5.2 RDMAP Interactions with the ULP

RDMAP provides the ULP with access to the following RDMA Operations as defined in this specification:

- * Send
- * Send with Solicited Event
- * Send with Invalidate
- * Send with Solicited Event and Invalidate
- * RDMA Write
- * RDMA Read

For Send Operation Types, the following are the interactions between the RDMAP Layer and the ULP:

- * At the Data Source:
 - * The ULP passes to the RDMAP Layer the following:
 - * ULP Message Length
 - * ULP Message
 - * An indication of the Send Operation Type, where the valid types are: Send, Send with Solicited Event, Send with Invalidate, or Send with Solicited Event and Invalidate.
 - * An Invalidate STag, if the Send Operation Type was Send with Invalidate or Send with Solicited Event and Invalidate.
 - * When the Send Operation Type Completes, an indication of the Completion results.
- * At the Data Sink:
 - * If the Send Operation Type Completed successfully, the RDMAP Layer passes the following information to the ULP Layer:
 - * ULP Message Length
 - * ULP Message

- * An Event, if the Data Sink is configured to generate an Event. 1
2
3
4
- * An Invalidated STag, if the Send Operation Type was Send with Invalidate or Send with Solicited Event and Invalidate. 5
6
7
8
- * If the Send Operation Type Completed in error, the Data Sink RDMAP Layer will pass up the corresponding error information to the Data Sink ULP and send a Terminate Message to the Data Source RDMAP Layer. The Data Source RDMAP Layer will then pass up the Terminate Message to the ULP. 9
10
11
12
13
14

For RDMA Write Operations, the following are the interactions between the RDMAP Layer and the ULP: 15
16
17

* At the Data Source: 18
19

- * The ULP passes to the RDMAP Layer the following: 20
21
 - * ULP Message Length 22
23
 - * ULP Message 24
25
 - * Data Sink STag 26
27
 - * Data Sink Tagged Offset 28
29
- * When the RDMA Write Operation Completes, an indication of the Completion results. 30
31
32

* At the Data Sink: 33
34

- * If the RDMA Write completed successfully, the RDMAP Layer does not Deliver the RDMA Write to the ULP. It does Place the ULP Message transferred through the RDMA Write Message into the ULP Buffer. 35
36
37
38
39
- * If the RDMA Write completed in error, the Data Sink RDMAP Layer will pass up the corresponding error information to the Data Sink ULP and send a Terminate Message to the Data Source RDMAP Layer. The Data Source RDMAP Layer will then pass up the Terminate Message to the ULP. 40
41
42
43
44
45

For RDMA Read Operations, the following are the interactions between the RDMAP Layer and the ULP: 46
47
48

* At the Data Sink: 49
50

- * The ULP passes to the RDMAP Layer the following: 51

- * ULP Message Length
- * Data Source STag
- * Data Sink STag
- * Data Source Tagged Offset
- * Data Sink Tagged Offset
- * When the RDMA Read Operation Completes, an indication of the Completion results.
- * At the Data Source:
 - * If no error occurred while processing the RDMA Read Request, the Data Source will not pass up any information to the ULP.
 - * If an error occurred while processing the RDMA Read Request, the Data Source RDMAP Layer will pass up the corresponding error information to the Data Source ULP and send a Terminate Message to the Data Sink RDMAP Layer. The Data Sink RDMAP Layer will then pass up the Terminate Message to the ULP.

For STags made available to the RDMAP Layer, following are the interactions between the RDMAP Layer and the ULP:

- * If the ULP enables an STag, the ULP passes to the RDMAP Layer the:
 - * STag;
 - * range of Tagged Offsets that are associated with a given STag;
 - * remote access rights (read, write, or read and write) associated with a given, valid STag; and
 - * association between a given STag and a given RDMAP Stream.
- * If the ULP disables an STag, the ULP passes to the RDMAP Layer the STag.

If an error occurs at the RDMAP Layer, the RDMAP Layer may pass back error information (e.g. the content of a Terminate Message) to the ULP.

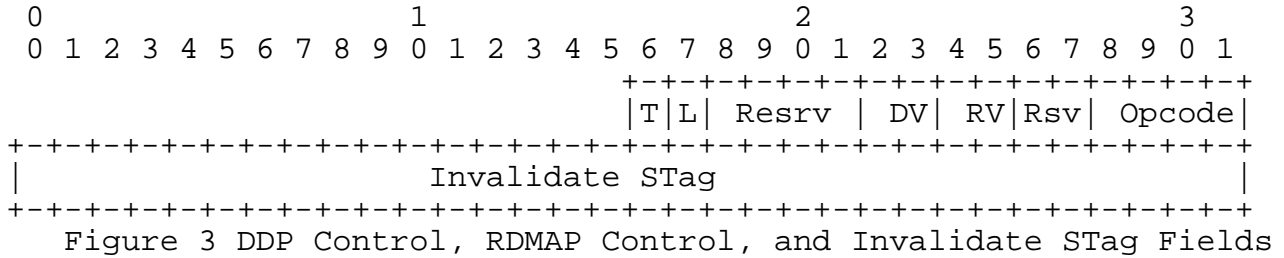
6 Header Format

The control information of RDMA Messages is included in DDP protocol defined header fields, with the following exceptions:

- * The first octet reserved for ULP usage on all DDP Messages in the DDP Protocol (i.e. the RsvdULP Field) is used by RDMAP to carry the RDMA Message Opcode and the RDMAP version. This octet is known as the RDMAP Control Field in this specification. For Send with Invalidate and Send with Solicited Event and Invalidate, RDMAP uses the second through fifth octets provided by DDP on Untagged DDP Messages to carry the STag that will be Invalidated.
- * The RDMA Message length is passed by the RDMAP layer to the DDP layer on all outbound transfers.
- * For RDMA Read Request Messages, the RDMA Read Message Size is included in the RDMA Read Request Header.
- * The RDMA Message length is passed to the RDMAP Layer by the DDP layer on inbound Untagged Buffer transfers.
- * Two RDMA Messages carry additional RDMAP headers. The RDMA Read Request carries the Data Sink and Data Source buffer descriptions, including buffer length. The Terminate carries additional information associated with the error that caused the Terminate.

6.1 RDMAP Control and Invalidate STag Field

The version of RDMAP defined by this specification uses all 8 bits of the RDMAP Control Field. The first octet reserved for ULP use in the DDP Protocol MUST be used by the RDMAP to carry the RDMAP Control Field. The ordering of the bits in the first octet MUST be as defined in Figure 3 DDP Control, RDMAP Control, and Invalidate STag Field. For Send with Invalidate and Send with Solicited Event and Invalidate, the second through fifth octets of the DDP RsvdULP field MUST be used by RDMAP to carry the Invalidate STag. Figure 3 DDP Control, RDMAP Control, and Invalidate STag Field depicts the format of the DDP Control and RDMAP Control fields. (Note: In Figure 3 DDP Control, RDMAP Control, and Invalidate STag Field, the DDP Header is offset by 16 bits to accommodate the MPA header defined in [MPA]. The MPA header is only present if DDP is layered on top of MPA.)



All RDMA Messages handed by the RDMAP Layer to the DDP layer MUST define the value of the Tagged flag in the DDP Header. Figure 4 RDMA Usage of DDP Fields MUST be used to define the value of the Tagged flag that is handed to the DDP Layer for each RDMA Message.

Figure 4 RDMA Usage of DDP Fields defines the value of the RDMA Opcode field that MUST be used for each RDMA Message.

Figure 4 RDMA Usage of DDP Fields defines when the STag, Queue Number, and Tagged Offset fields MUST be provided for each RDMA Message.

For this version of the RDMAP, all RDMA Messages MUST have:

- * Bits 24-25; RDMA Version field: 00b .
- * Bits 26-27; Reserved. MUST be set to zero by sender, ignored by the receiver.
- * Bits 28-31; OpCode field: see Figure 4 RDMA Usage of DDP Fields.
- * Bits 32-63; Invalidate STag. However, this field is only valid for Send with Invalidate and Send with Solicited Event and Invalidate Messages (see Figure 4 RDMA Usage of DDP Fields). For Send, Send with Solicited Event, RDMA Read Request, and Terminate, the Invalidate STag field MUST be set to zero on transmit and ignored by the receiver.

RDMA Message OpCode	Message Type	Tagged Flag	S Tag and TO	Queue Number	Invalidate S Tag	Message Length Communicated between DDP and RDMAP
0000b	RDMA Write	1	Valid	N/A	N/A	Yes
0001b	RDMA Read Request	0	N/A	1	N/A	Yes
0010b	RDMA Read Response	1	Valid	N/A	N/A	Yes
0011b	Send	0	N/A	0	N/A	Yes
0100b	Send with Invalidate	0	N/A	0	Valid	Yes
0101b	Send with SE	0	N/A	0	N/A	Yes
0110b	Send with SE and Invalidate	0	N/A	0	Valid	Yes
0111b	Terminate	0	N/A	2	N/A	Yes
1000b to 1111b	Reserved	Not Specified				

Figure 4 RDMA Usage of DDP Fields

Note: N/A means Not Applicable.

6.2 RDMA Message Definitions

The following figure defines which RDMA Headers MUST be used on each RDMA Message and which RDMA Messages are allowed to carry ULP payload:

RDMA Message OpCode	Message Type	RDMA Header Used	ULP Message allowed in the RDMA Message
0000b	RDMA Write	None	Yes
0001b	RDMA Read Request	RDMA Read Request Header	No
0010b	RDMA Read Response	None	Yes
0011b	Send	None	Yes
0100b	Send with Invalidate	None	Yes
0101b	Send with SE	None	Yes
0110b	Send with SE and Invalidate	None	Yes
0111b	Terminate	Terminate Header	No
1000b to 1111b	Reserved	Not Specified	

Figure 5 RDMA Message Definitions

6.3 RDMA Write Header

The RDMA Write Message does not include an RDMAP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field. The RDMA Write Message is fully described by the DDP Headers of the DDP Segments associated with the Message.

See section 12 Appendix for a description of the DDP Segment format associated with RDMA Write Messages.

6.4 RDMA Read Request Header

The RDMA Read Request Message carries an RDMA Read Request Header that describes the Data Sink and Data Source Buffers used by the RDMA Read operation. The RDMA Read Request Header immediately follows the DDP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field. The following figure depicts the RDMA Read Request Header that MUST be used for all RDMA Read Request Messages:

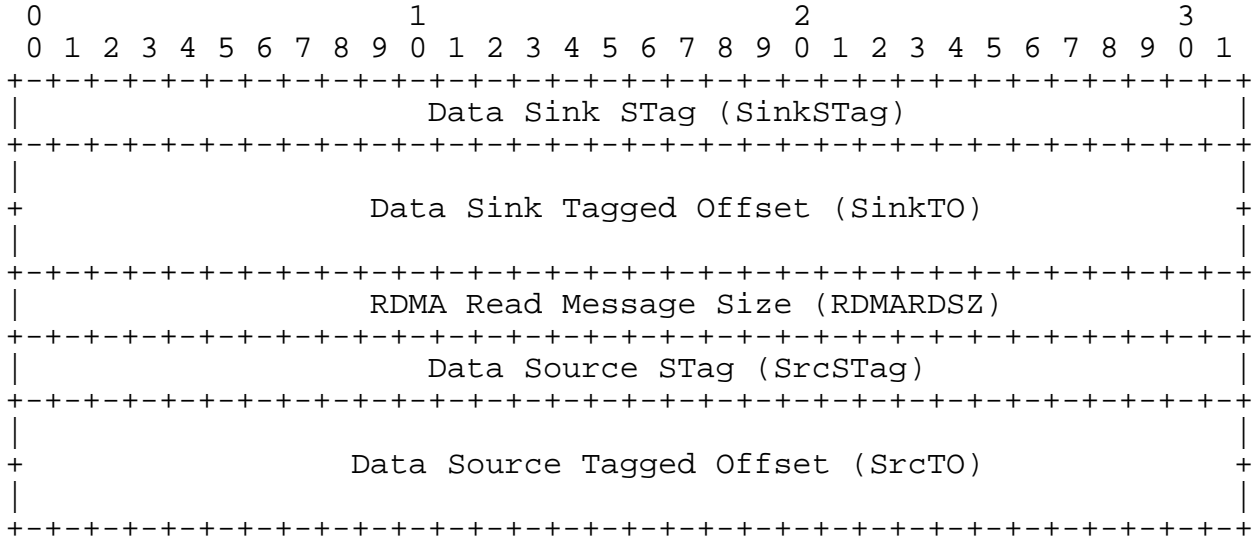


Figure 6 RDMA Read Request Header Format

Data Sink Steering Tag: 32 bits.

The Data Sink Steering Tag identifies the Data Sink's Tagged Buffer. This field MUST be copied, without interpretation, from the RDMA Read Request into the corresponding RDMA Read Response and allows the Data Sink to place the returning data. The STag is associated with the RDMAP Stream through a mechanism that is outside the scope of the RDMAP specification (see Section 10.3 Other Security Considerations).

Data Sink Tagged Offset: 64 bits.

The Data Sink Tagged Offset specifies the starting offset, in octets, from the base of the Data Sink's Tagged Buffer, where the data is to be written by the Data Source. This field is copied from the RDMA Read Request into the corresponding RDMA Read Response and allows the Data Sink to place the returning data. The Data Sink Tagged Offset MAY start at an arbitrary offset.

The Data Sink STag and Data Sink Tagged Offset fields describe the buffer to which the RDMA Read data is written.

Note: the DDP Layer protects against a wrap of the Data Sink Tagged Offset.

RDMA Read Message Size: 32 bits.

The RDMA Read Message Size is the amount of data, in octets, read from the Data Source. A single RDMA Read Request Message can retrieve from 0 to $2^{32}-1$ data octets from the Data Source.

Data Source Steering Tag: 32 bits.

The Data Source Steering Tag identifies the Data Source's Tagged Buffer. The STag is associated with the RDMAP Stream through a mechanism that is outside the scope of the RDMAP specification (see Section 10.3 Other Security Considerations).

Data Source Tagged Offset: 64 bits.

The Tagged Offset specifies the starting offset, in octets, that is to be read from the Data Source's Tagged Buffer. The Data Source Tagged Offset MAY start at an arbitrary offset.

The Data Source STag and Data Source Tagged Offset fields describe the buffer from which the RDMA Read data is read.

See Section 9.2 Errors Detected at the Remote Peer on Incoming RDMA Messages for a description of error checking required upon processing of an RDMA Read Request at the Data Source.

6.5 RDMA Read Response Header

The RDMA Read Response Message does not include an RDMAP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field. The RDMA Read Response Message is fully described by the DDP Headers of the DDP Segments associated with the Message.

See Section 12 Appendix for a description of the DDP Segment format associated with RDMA Read Response Messages.

6.6 Send Header and Send with Solicited Event Header

The Send and Send with Solicited Event Message do not include an RDMAP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field. The Send and Send with Solicited Event Message are fully described by the DDP Headers of the DDP Segments associated with the Message.

See Section 12 Appendix for a description of the DDP Segment format associated with Send and Send with Solicited Event Messages.

6.7 Send with Invalidate Header and Send with SE and Invalidate Header

The Send with Invalidate and Send with Solicited Event and Invalidate Message do not include an RDMAP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field and the Invalidate STag field (see section 6.1 RDMAP Control and Invalidate STag Field). The Send with Invalidate and Send with Solicited Event and Invalidate Message are fully described by the DDP Headers of the DDP Segments associated with the Message.

See Section 12 Appendix for a description of the DDP Segment format associated with Send and Send with Solicited Event Messages.

6.8 Terminate Header

The Terminate Message carries a Terminate Header that contains additional information associated with the cause of the Terminate. The Terminate Header immediately follows the DDP header. The RDMAP layer passes to the DDP layer an RDMAP Control Field. The following figure depicts a Terminate Header that MUST be used for the Terminate Message:

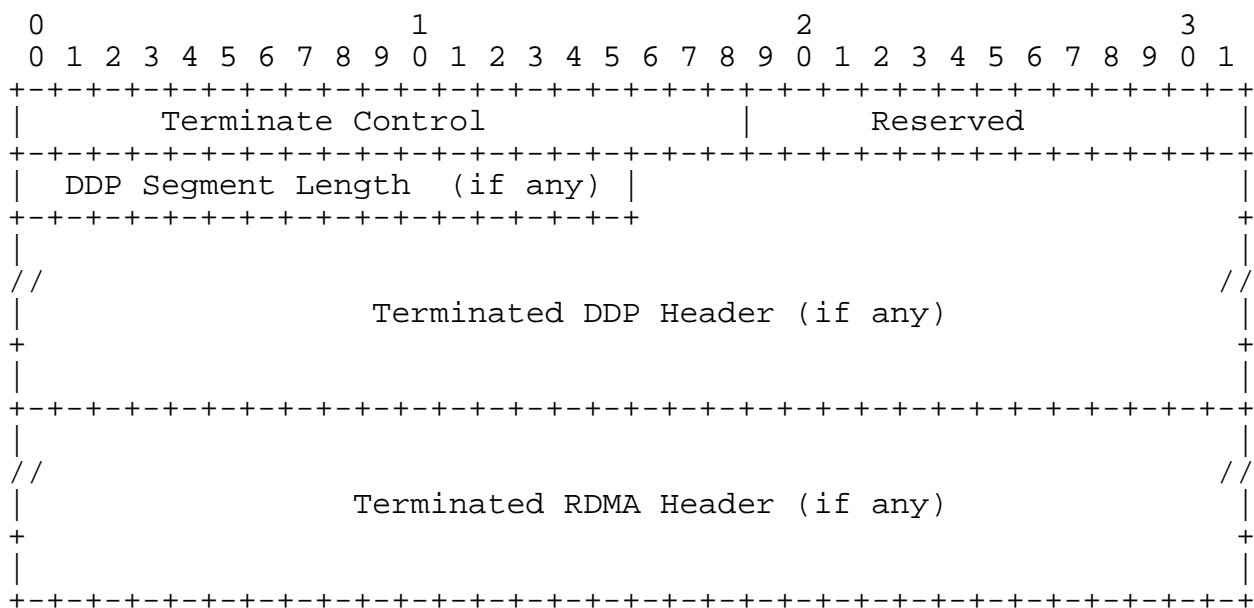


Figure 7 Terminate Header Format

Terminate Control: 19 bits.

The Terminate Control field MUST have the format defined in Figure 8 Terminate Control Field.

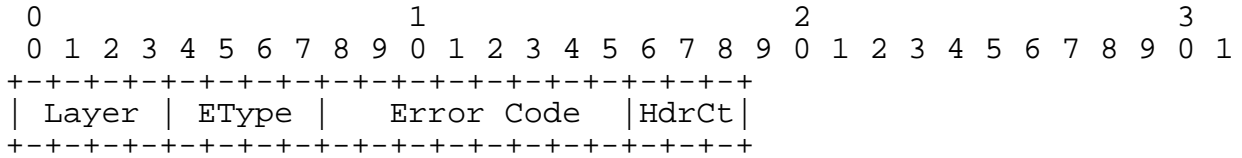


Figure 8 Terminate Control Field

- * Figure 9 Terminate Control Field Values defines the valid values that MUST be used for this field.
- * Layer: 4 bits.
Identifies the layer that encountered the error.
- * EType (RDMA Error Type): 4 bits.
Identifies the type of error that caused the Terminate. When the error is detected at the RDMAP Layer, the RDMAP Layer inserts the Error Type into this field. When the error is detected at a LLP layer, a LLP layer creates the Error Type and the DDP layer passes it up to the RDMAP Layer, and the RDMAP Layer inserts it into this field.
- * Error Code: 8 bits.
This field identifies the specific error that caused the Terminate. When the error is detected at the RDMAP Layer, the RDMAP Layer creates the Error Code. When the error is detected at a LLP layer, a LLP layer creates the Error Code and the DDP layer passes it up to the RDMAP Layer, and the RDMAP Layer inserts it into this field.
- * HdrCt: 3 bits.
Header control bits:
 - * M: 1 bit. DDP Segment Length valid. See Figure 10 for when this bit SHOULD be set.
 - * D: 1 bit. DDP Header Included. See Figure 10 for when this bit SHOULD be set.
 - * R: 1 bit. RDMAP Header Included. See Figure 10 for when this bit SHOULD be set.

Layer	Layer Name	Error Type	Error Type Name	Error Code	Error Code Name
0000b	RDMA	0000b	Local Catastrophic Error	None	None
		0001b	Remote Protection Error	00X	Invalid STag
				01X	Base or bounds violation
				02X	Access rights violation
				03X	STag not associated with RDMAP Stream
				04X	TO wrap
		0010b	Remote Operation Error	05X	Invalid RDMAP version
				06X	Unexpected OpCode
				07X	Catastrophic error, localized to RDMAP Stream
				08X	Catastrophic error, global
0001b	DDP	See DDP Specification [DDP] for a description of the values and names.			
0010b	MPA	See MPA Specification [MPA] for a description of the values and names.			

Figure 9 Terminate Control Field Values

Reserved: 8 bits. This field MUST be set to zero on transmit, ignored on receive.

DDP Segment Length: 16 bits

The length handed up by the DDP Layer when the error was detected. It MUST be valid if the M bit is set. It MUST be present when the D bit is set.

Terminated DDP Header: 112 bits for Tagged Messages and 144 bits for Untagged Messages.

The DDP Header of the incoming Message that is associated with the Terminate. The DDP Header is not present if the Terminate Error Type is a Local Catastrophic Error. It MUST be present if the D bit is set.

Terminated RDMA Header: 224 bits.

The Terminated RDMA Header is only sent back if the terminate is associated with an RDMA Read Request Message. It MUST be present if the R bit is set.

If the terminate occurs before the first RDMA Read Request byte is processed, the original RDMA Read Request Header is sent back.

If the terminate occurs after the first RDMA Read Request byte is processed, the RDMA Read Request Header is updated to reflect the current location of the RDMA Read operation that is in process:

- * Data Sink STag = Data Sink STag originally sent in the RDMA Read Request.
- * Data Sink Tagged Offset = Current offset into the Data Sink Tagged Buffer. For example if the RDMA Read Request was terminated after 2048 octets were sent, then the Data Sink Tagged Offset = the original Data Sink Tagged Offset + 2048.
- * Data Message size = Number of bytes left to transfer.
- * Data Source STag = Data Source STag in the RDMA Read Request.
- * Data Source Tagged Offset = Current offset into the Data Source Tagged Buffer. For example if the RDMA Read Request was terminated after 2048 octets were sent, then the Data Source Tagged Offset = the original Data Source Tagged Offset + 2048.

Figure 10 Error Type to RDMA Message Mapping maps layer name and error types to each RDMA Message type:

Layer Name	Error Type Name	Terminate Includes DDP Header and DDP Segment Length	Terminate Includes RDMA Header	What type of RDMA Message can cause the error
	Local Catastrophic Error	No	No	Any
RDMA	Remote Protection Error	Yes, if possible	Yes	Only RDMA Read Request, Send with Invalidate, and Send with SE and Invalidate
	Remote Operation Error	Yes, if possible	No	Any
DDP	See DDP Spec [DDP]	Yes	No	Any
MPA	See MPA Spec [MPA]	No	No	Any

Figure 10 Error Type to RDMA Message Mapping

7 Data Transfer

7.1 RDMA Write Message

An RDMA Write is used by the Data Source to transfer data to a previously Advertised Tagged Buffer at the Data Sink. The RDMA Write Message has the following semantics:

- * AN RDMA Write Message MUST reference a Tagged Buffer. That is, the Data Source RDMAP Layer MUST request that the DDP layer mark the Message as Tagged.
- * A valid RDMA Write Message MUST NOT be delivered to the Data Sink's ULP (i.e. it is placed by the DDP layer).
- * At the Remote Peer, when an invalid RDMA Write Message is delivered to the Remote Peer's RDMAP Layer, an error is surfaced (see section 9.1 RDMAP Error Surfacing).
- * The Tagged Offset of a Tagged Buffer MAY start at a non-zero value.
- * AN RDMA Write Message MAY target all or part of a previously Advertised buffer.
- * The RDMAP does not define how the buffer(s) used by an outbound RDMA Write is defined and how it is addressed. For example, an implementation of RDMA may choose to allow a gather-list of non-contiguous data blocks to be the source of an RDMA Write. In this case, the data blocks would be combined by the Data Source and sent as a single RDMA Write Message to the Data Sink.
- * The Data Source RDMAP Layer MUST issue RDMA Write Messages to the DDP layer in the order they were submitted by the ULP.
- * At the Data Source, a subsequent Send (Send with Invalidate, Send with Solicited Event, or Send with Solicited Event and Invalidate) Message MAY be used to signal Delivery of previous RDMA Write Messages to the Data Sink, if desired by the ULP.
- * If the Local Peer wishes to write to multiple Tagged Buffers on the Remote Peer, the Local Peer MUST use multiple RDMA Write Messages. That is, a single RDMA Write Message can only write to one remote Tagged Buffer.
- * The Data Source MAY issue a zero length RDMA Write Message.

7.2 RDMA Read Operation

The RDMA Read operation MUST consist of a single RDMA Read Request Message and a single RDMA Read Response Message.

7.2.1 RDMA Read Request Message

An RDMA Read Request is used by the Data Sink to transfer data from a previously Advertised Tagged Buffer at the Data Source to a Tagged Buffer at the Data Sink. The RDMA Read Request Message has the following semantics:

- * AN RDMA Read Request Message MUST reference an Untagged Buffer. That is, the Local Peer's RDMAP Layer MUST request that the DDP mark the Message as Untagged.
- * One RDMA Read Request Message MUST consume one Untagged Buffer.
- * The Remote Peer's RDMAP Layer MUST process an RDMA Read Request Message. A valid RDMA Read Request Message MUST NOT be delivered to the Data Sink's ULP (i.e. it is processed by the RDMAP layer).
- * At the Remote Peer, when an invalid RDMA Read Request Message is delivered to the Remote Peer's RDMAP Layer, an error is surfaced (see section 9.1 RDMAP Error Surfacing).
- * AN RDMA Read Request Message MUST reference the RDMA Read Request Queue. That is, the Local Peer's RDMAP Layer MUST request that the DDP layer set the Queue Number field to one.
- * The Local Peer MUST pass to the DDP Layer RDMA Read Request Messages in the order they were submitted by the ULP.
- * The Remote Peer MUST process the RDMA Read Request Messages in the order they were sent.
- * If the Local Peer wishes to read from multiple Tagged Buffers on the Remote Peer, the Local Peer MUST use multiple RDMA Read Request Messages. That is, a single RDMA Read Request Message MUST only read from one remote Tagged Buffer.
- * AN RDMA Read Request Message MAY target all or part of a previously Advertised buffer.
- * If the Data Source receives a valid RDMA Read Request Message it MUST respond with a valid RDMA Read Response Message.
- * The Data Sink MAY issue a zero length RDMA Read Request Message, by setting the RDMA Read Message Size field to zero in the RDMA Read Request Header.

- * If the Data Source receives a non-zero length RDMA Read Message Size, the Data Source RDMAP MUST validate the Data Source STag and Data Source Tagged Offset contained in the RDMA Read Request Header.
- * If the Data Source receives an RDMA Read Request Header with the RDMA Read Message Size set to zero, the Data Source RDMAP:
 - * MUST NOT validate the Data Source STag and Data Source Tagged Offset contained in the RDMA Read Request Header, and
 - * MUST respond with a zero length RDMA Read Response Message.

7.2.2 RDMA Read Response Message

The RDMA Read Response Message uses the DDP Tagged Buffer Model to Deliver the contents of a previously requested Data Source Tagged Buffer to the Data Sink, without any involvement from the ULP at the Remote Peer. The RDMA Read Response Message has the following semantics:

- * The RDMA Read Response Message for the associated RDMA Read Request Message travels in the opposite direction.
- * An RDMA Read Response Message MUST reference a Tagged Buffer. That is, the Data Source RDMAP Layer MUST request that the DDP mark the Message as Tagged.
- * The Data Source MUST ensure that a sufficient number of Untagged Buffers are available on the RDMA Read Request Queue (Queue with DDP Queue Number 1) to support the maximum number of RDMA Read Requests negotiated by the ULP.
- * The RDMAP Layer MUST Deliver the RDMA Read Response Message to the ULP.
- * At the Remote Peer, when an invalid RDMA Read Response Message is delivered to the Remote Peer's RDMAP Layer, an error is surfaced (see section 9.1 RDMAP Error Surfacing).
- * The Tagged Offset of a Tagged Buffer MAY start at a non-zero value.
- * The Data Source RDMAP Layer MUST pass RDMA Read Response Messages to the DDP layer in the order that the RDMA Read Request Messages were received by the RDMAP Layer at the Data Source.
- * The Data Sink MAY validate that the STag, Tagged Offset, and length of the RDMA Read Response Message are the same as the

STag, Tagged Offset, and length included in the corresponding RDMA Read Request Message.

- * A single RDMA Read Response Message MUST write to one remote Tagged Buffer. If the Data Sink wishes to Read multiple Tagged Buffers, the Data Sink can use multiple RDMA Read Request Messages.

7.3 Send Message Type

The Send Message Type uses the DDP Untagged Buffer Model to transfer data from the Data Source into an Untagged Buffer at the Data Sink.

- * A Send Message Type MUST reference an Untagged Buffer. That is, the Local Peer's RDMAP Layer MUST request that the DDP layer mark the Message as Untagged.
- * One Send Message Type MUST consume one Untagged Buffer.
 - * The ULP Message sent using a Send Message Type MAY be less than or equal to the size of the consumed Untagged Buffer. The RDMAP Layer communicates to the ULP the size of the data written into the Untagged Buffer.
 - * If the ULP Message sent via Send Message Type is larger than the Data Sink's Untagged Buffer, it is an error (see section 9.1 RDMAP Error Surfacing).
- * At the Remote Peer, the Send Message Type MUST be Delivered to the Remote Peer's ULP in the order they were sent.
- * After the Send with Solicited Event or Send with Solicited Event and Invalidate Message is Delivered to the ULP, the RDMAP MAY generate an Event, if the Data Sink is configured to generate such an Event.
- * At the Remote Peer, when an invalid Send Message Type is Delivered to the Remote Peer's RDMAP Layer, an error is surfaced (see section 9.1 RDMAP Error Surfacing).
- * The RDMAP does not define how the buffer(s) used by an outbound Send Message Type is defined and how it is addressed. For example, an implementation of RDMA may choose to allow a gather-list of non-contiguous data blocks to be the source of a Send Message Type. In this case, the data blocks would be combined by the Data Source and sent as a single Send Message Type to the Data Sink.
- * For a Send Message Type, the Local Peer's RDMAP Layer MUST request that the DDP layer set the Queue Number field to zero.

- * The Local Peer MUST issue Send Message Type Messages in the order they were submitted by the ULP.
- * The Data Source MAY pass a zero length Send Message Type. A zero length Send Message Type MUST consume an Untagged Buffer at the Data Sink. A Send with Invalidate or Send with Solicited Event and Invalidate Message MUST reference an STag. That is, the Local Peer's RDMAP Layer MUST pass the RDMA control field and the STag that will be Invalidated to the DDP layer.
- * When the Send with Invalidate and Send with Solicited Event and Invalidate Message are Delivered to the Remote Peer's RDMAP Layer, the RDMAP Layer MUST:
 - * Verify the STag that is associated with the RDMAP Stream; and
 - * Invalidate the STag if it is associated with the RDMAP Stream; or Issue a Terminate Message if the STag is not associated with the RDMAP Stream.

7.4 Terminate Message

The Terminate Message uses the DDP Untagged Buffer Model to transfer error related information from the Data Source into an Untagged Buffer at the Data Sink and then ceases all further communications on the underlying DDP Stream. The Terminate Message has the following semantics:

- * A Terminate Message MUST reference an Untagged Buffer. That is, the Local Peer's RDMAP Layer MUST request that the DDP layer mark the Message as Untagged.
- * A Terminate Message references the Terminate Queue. That is, the Local Peer's RDMAP Layer MUST request that the DDP layer set the Queue Number field to two.
- * One Terminate Message MUST consume one Untagged Buffer.
- * On a single RDMAP Stream, the RDMAP layer MUST guarantee placement of a single Terminate Message.
- * A Terminate Message MUST be Delivered to the Remote Peer's RDMAP Layer. The RDMAP Layer MUST Deliver the Terminate Message to the ULP.
- * At the Remote Peer, when an invalid Terminate Message is delivered to the Remote Peer's RDMAP Layer, an error is surfaced (see section 9.1 RDMAP Error Surfacing).

- * The RDMAP Layer Completes in error all ULP Operations that have not been provided to the DDP layer.
- * After sending a Terminate Message on an RDMAP Stream, the Local Peer MUST NOT send any more Messages on that specific RDMAP Stream.
- * After receiving a Terminate Message on an RDMAP Stream, the Remote Peer MAY stop sending Messages on that specific RDMAP Stream.

7.5 Ordering and Completions

It is important to understand the difference between Placement and Delivery ordering since RDMAP provides quite different semantics for the two.

Note that many current protocols, both as used in the Internet and elsewhere, assume that data is both Placed and Delivered in order. This allowed applications to take a variety of shortcuts by taking advantage of this fact. For RDMAP, many of these shortcuts are no longer safe to use, and could cause application failure.

The following rules apply to implementations of the RDMAP protocol. Note, in these rules Send includes Send, Send with Invalidate, Send with Solicited Event, and Send with Solicited Event and Invalidate:

1. RDMAP does not provide ordering among Messages on different RDMAP Streams.
2. RDMAP does not provide ordering between operations that are generated from the two ends of an RDMAP Stream.
3. RDMA Messages that use Tagged and Untagged Buffers MAY be Placed in any order. If an application uses overlapping buffers (points different Messages or portions of a single Message at the same buffer), then it is possible that the last incoming write to the Data Sink buffer will not be the last outgoing data sent from the Data Source.
4. For a Send operation, the contents of an Untagged Buffer at the Data Sink MAY be indeterminate until the Send is Delivered to the ULP at the Data Sink.
5. For an RDMA Write operation, the contents of the Tagged Buffer at the Data Sink MAY be indeterminate until a subsequent Send is Delivered to the ULP at the Data Sink.

6. For an RDMA Read operation, the contents of the Tagged Buffer at the Data Sink MAY be indeterminate until the RDMA Read Response Message has been Delivered at the Local Peer.

Statements 4, 5, and 6 imply "no peeking" at the data to see if it is done. It is possible for some data to arrive before logically earlier data does, and peeking may cause unpredictable application failure

7. If the ULP or Application modifies the contents of Tagged or Untagged Buffers being modified by an RDMA Operation while the RDMAP is processing the RDMA Operation, the state of the Buffers is indeterminate.
8. If the ULP or Application modifies the contents of Tagged or Untagged Buffers read by an RDMA Operation while the RDMAP is processing the RDMA Operation, the results of the read are indeterminate.
9. The Completion of an RDMA Write or Send Operation at the Local Peer does not guarantee that the ULP Message has yet reached the Remote Peer ULP Buffer or been examined by the Remote ULP.
10. Send Messages MUST be Delivered to the ULP at the Remote Peer after they are Delivered to RDMAP by DDP and in the order that they were Delivered to RDMAP.

Note that DDP ordering rules ensure that this will be the same order that they were submitted at the Local Peer and that any prior RDMA Writes have been submitted for ordered Placement at the Remote Peer. This means that when the ULP sees the Delivery of the Send, the memory buffers targeted by any preceding RDMA Writes and Sends are available to be accessed locally or remotely as authorized. If the ULP overlaps its buffers for different operations, the data from the RDMA Write or Send may be overwritten by subsequent RDMA Operations before the ULP receives and processes the Delivery.

11. RDMA Read Response Messages MUST be Delivered to the ULP at the Remote Peer after they are Delivered to RDMAP by DDP and in the order that they were Delivered to RDMAP.

DDP ordering rules ensure that this will be the same order that they were submitted at the Local Peer. This means that when the ULP sees the Delivery of the RDMA Read Response, the memory buffers targeted by the RDMA Read Response are available to be accessed locally or remotely as authorized. If the ULP overlaps its buffers for different operations, the data from the RDMA Read Response may be overwritten by subsequent RDMA Operations before the ULP receives and processes the Delivery.

12. RDMA Read Request Messages, including zero-length RDMA Read Requests, MUST NOT start processing at the Remote Peer until they have been Delivered to RDMAP by DDP.

Note the ULP is assured that data written can be read back. If an RDMA Read Request is issued by the local peer, targeting the same ULP Buffer as a preceding Send or RDMA Write (in the same direction as the RDMA Read Request), the remote peer will send back the data written by the Send or RDMA Write. Thus assuring that subsequent local or remote accesses to the ULP Buffer contain the data written by the Send or RDMA Write. Note: this statement requires that the buffers are not arranged such that the data is over-written either by the incoming data, or other applications.

RDMA Read Response Messages MAY be generated at the Remote Peer after subsequent RDMA Write Messages or Send Messages have been Placed or Delivered. Therefore, when an application does an RDMA Read Request followed by an RDMA Write (or Send) to the same buffer, it may get the data from the later RDMA Write (or Send) in the RDMA Read Response Message, even though the operations completed in order at the Local Peer. If this behavior is not desired, the Local Peer ULP must Fence the later RDMA write (or Send) by withholding the RDMA Write Message until all outstanding RDMA Read Responses have been Delivered.

13. The RDMAP Layer MUST submit RDMA Messages to the DDP layer in the order the RDMA Operations are submitted to the RDMAP Layer by the ULP.
14. A Send or RDMA Write Message MUST NOT be considered Complete at the Local Peer (Data Source) until it has been successfully completed at the DDP layer.
15. RDMA Operations MUST be Completed at the Local Peer in the order that they were submitted by the ULP.
16. At the Data Sink, an incoming Send Message MUST be Delivered to the ULP only after the DDP Message has been Delivered to the RDMAP Layer by the DDP layer.
17. RDMA Read Response Message processing at the Remote Peer (reading the specified Tagged Buffer) MUST be started only after the RDMA Read Request Message has been Delivered by the DDP layer (thus all previous RDMA Messages have been properly submitted for ordered Placement).

18. Send Messages MAY be Completed at the Remote Peer (Data Sink) before prior incoming RDMA Read Request Messages have completed their response processing.
19. An RDMA Read operation MUST NOT be Completed at the Local Peer until the DDP layer Delivers the associated incoming RDMA Read Response Message.
20. If more than one outstanding RDMA Read Request Message is supported by both peers, the RDMA Read Response Messages MUST be submitted to the DDP layer on the Remote Peer in the order the RDMA Read Request Messages were Delivered by DDP, but the actual read of the buffer contents MAY take place in any order at the Remote Peer.

This simplifies Local Peer Completion processing for RDMA Reads in that a Delivered RDMA Read Response MUST be sufficient to Complete the RDMA Read Operation.

8 RDMA Stream Management

RDMA Stream management consists of RDMA Stream Initialization and RDMA Stream Termination.

8.1 Stream Initialization

RDMA Stream initialization occurs after the LLP Stream has been created (e.g. for DDP/MPA over TCP the first TCP Segment after the SYN, SYN/ACK exchange). The ULP is responsible for transitioning the LLP Stream into RDMA enabled mode. The switch to RDMA mode can happen immediately at LLP Stream initialization or at any time thereafter. Once in RDMA enabled mode, an implementation MUST send only RDMA Messages across the transport Stream until the RDMA Stream is torn down.

For each direction of an RDMA Stream:

- * For a given RDMA Stream, the number of outstanding RDMA Read Requests is limited per RDMA Stream direction.
- * It is the ULP's responsibility to set the maximum number of outstanding, inbound RDMA Read Requests per RDMA Stream direction.
- * The RDMA Layer MUST provide the maximum number of outstanding, inbound RDMA Read Requests per RDMA Stream direction that were negotiated between the ULP and the Local Peer's RDMA Layer. The negotiation mechanism is outside the scope of this specification.
- * It is the ULP's responsibility to set the maximum number of outstanding, outbound RDMA Read Requests per RDMA Stream direction.
- * The RDMA Layer MUST provide the maximum number of outstanding, outbound RDMA Read Requests for the RDMA Stream direction that were negotiated between the ULP and the Local Peer's RDMA Layer. The negotiation mechanism is outside the scope of this specification.
- * The Local Peer's ULP is responsible for negotiating with the Remote Peer's ULP the maximum number of outstanding RDMA Read Requests for the RDMA Stream direction. It is recommended that the ULP set the maximum number of outstanding, inbound RDMA Read Requests equal to the maximum number of outstanding, outbound RDMA Read Requests for a given RDMA Stream direction.
- * For outbound RDMA Read Requests, the RDMA Layer MUST NOT exceed the maximum number of outstanding, outbound RDMA Read Requests

that were negotiated between the ULP and the Local Peer's RDMAP Layer.

- * For inbound RDMA Read Requests, the RDMAP Layer MUST NOT exceed the maximum number of outstanding, inbound RDMA Read Requests that were negotiated between the ULP and the Local Peer's RDMAP Layer.

8.2 Stream Teardown

There are three methods for terminating an RDMAP Stream: ULP Graceful Termination, RDMAP Abortive Termination, and LLP Abortive Termination.

The ULP is responsible for performing ULP Graceful Termination. After a ULP Graceful Termination, either side of the Stream can initiate LLP Graceful Termination, using the graceful termination mechanism provided by the LLP.

RDMAP Abortive Termination allows the RDMAP to issue a Terminate Message describing the reason the RDMAP Stream was terminated. The next section (8.2.1 RDMAP Abortive Termination) describes the RDMAP Abortive Termination in detail.

LLP Abortive Termination results due to a LLP error and causes the RDMAP Stream to be torn down midstream, without an RDMAP Terminate Message. While this last method is highly undesirable, it is possible and the ULP should take this into consideration.

8.2.1 RDMAP Abortive Termination

RDMAP defines a Terminate operation that SHOULD be invoked when either an RDMAP error is encountered or a LLP error is surfaced to the RDMAP layer by the LLP.

It is not always possible to send the Terminate Message. For example, certain LLP errors may occur that cause the LLP Stream to be torn down before a) RDMAP is aware of the error, b) before RDMAP is able to send the Terminate Message, or c) after RDMAP has posted the Terminate Message to the LLP, but it has not yet been transmitted by the LLP.

Note that an RDMAP Abortive Termination may entail loss of data. In general, when a Terminate Message is received it is impossible to tell for sure what unacknowledged RDMA Messages were Completed successfully at the Remote Peer. Thus the state of all outstanding RDMA Messages is indeterminate and the Messages SHOULD be considered Completed in error.

When a peer sends or receives a Terminate Message, it MAY immediately teardown the LLP Stream. The peer SHOULD perform a graceful LLP teardown to ensure the Terminate Message is successfully Delivered.

See section 6.8 Terminate Header for a description of the Terminate Message and its contents. See section 7.4 Terminate Message for a description of the Terminate Message semantics.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

9 RDMAP Error Management

The RDMAP protocol does not have RDMAP or DDP layer error recovery operations built in. If everything is working, the LLP guarantees will ensure that the Messages are arriving at the destination.

If errors are detected at the RDMAP or DDP layer, then the RDMAP, DDP and LLP Streams are Abortively Terminated (see section 6.8 Terminate Header on page 26).

In general poor implementations or improper ULP programming causes the errors detected at the RDMAP and DDP layers. In these cases, returning a diagnostic termination error Message and closing the RDMAP Stream is far simpler than attempting to maintain the RDMAP Stream, particularly when the cause of the error is not known.

If an LLP does not support teardown of a Stream independent of other Streams and an RDMAP error results in the Termination of a specific Stream, then the LLP MUST label the Stream as an erroneous Stream and MUST NOT allow any further data transfer on that Stream after RDMAP requests the Stream to be torn down.

For a specific LLP connection, when all Streams are either gracefully torn down or are labeled as erroneous Streams, the LLP connection MUST be torn down.

Since errors are detected at the Remote Peer (possibly long) after RDMA Messages are passed to DDP and the LLP at the Local Peer and Completed, the sender cannot easily determine which of its Messages have been received. (RDMA Reads are an exception to this rule).

For a list of errors returned to the Remote Peer as a result of an Abortive Termination, see section 6.8 Terminate Header on page 26.

9.1 RDMAP Error Surfacing

If an error occurs at the Local Peer, the RDMAP layer MUST attempt to inform the local ULP that the error has occurred.

The Local Peer MUST send a Terminate Message for each of the following cases:

1. For Errors detected while creating RDMA Write, Send, Send with Invalidate, Send with Solicited Event, Send with Solicited Event and Invalidate, or RDMA Read Requests, or other reasons not directly associated with an incoming Message, the Terminate Message and Error code are sent instead of the request. In this case, the Error Type and Error Code fields are included in the Terminate Message, but the Terminated DDP Header and Terminated RDMA Header fields are set to zero.

2. For errors detected on an incoming RDMA Write, Send, Send with Invalidate, Send with Solicited Event, Send with Solicited Event and Invalidate, or Read Response Message (after the Message has been Delivered by DDP), the Terminate Message is sent at the earliest possible opportunity, preferably in the next outgoing RDMA Message. In this case, the Error Type, Error Code, ULP PDU Length, and Terminated DDP Header fields are included in the Terminate Message, but the Terminated RDMA Header field is set to zero.
3. For errors detected on an incoming RDMA Read Request Message (after the Message has been Delivered by DDP), the Terminate Message is sent at the earliest possible opportunity, preferably in the next outgoing RDMA Message. In this case, the Error Type, Error Code, ULP PDU Length, Terminated DDP Header, and Terminated RDMA Header fields are included in the Terminate Message.
4. If more than one error is detected on incoming RDMA Messages, before the Terminate Message can be sent, then the first RDMA Message (and its associated DDP Segment) that experienced an error MUST be captured by the Terminate Message in accordance with rules 2 and 3 above.

9.2 Errors Detected at the Remote Peer on Incoming RDMA Messages

On incoming RDMA Writes, RDMA Read Response, Sends, Send with Invalidate, Send with Solicited Event, Send with Solicited Event and Invalidate, and Terminate Messages, the following must be validated:

1. The DDP Layer MUST validate all DDP Segment fields.
2. The RDMA OpCode MUST be valid.
3. The RDMA Version MUST be valid.

Additionally, on incoming Send with Invalidate and Send with Solicited Event and Invalidate Messages, the following must also be validated:

4. The Invalidate STag MUST be valid.
5. The STag MUST be associated to this RDMAP Stream.

On incoming RDMA Request Messages, the following must be validated:

1. The DDP Layer MUST validate all Untagged DDP Segment fields.
2. The RDMA OpCode MUST be valid.

3. The RDMA Version MUST be valid.
4. For non-zero length RDMA Read Request Messages:
 - a. The Data Source STag MUST be valid.
 - b. The Data Source STag MUST be associated to this RDMAP Stream.
 - c. The Data Source Tagged Offset MUST fall in the range of legal offsets associated with the Data Source STag.
 - d. The sum of the Data Source Tagged Offset and the RDMA Read Message Size MUST fall in the range of legal offsets associated with the Data Source STag.
 - e. The sum of the Data Source Tagged Offset and the RDMA Read Message Size MUST NOT cause the Data Source Tagged Offset to wrap.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

10 Security Considerations

This section discusses both protocol-specific considerations and the implications of using RDMAP with existing security mechanisms.

10.1 Protocol-specific Security Considerations

The vulnerabilities of RDMAP to active third-party interference are no greater than any other protocol running over TCP. A third party, by injecting spoofed packets into the network that are Delivered to an RDMAP Data Sink, could launch a variety of attacks that exploit RDMAP-specific behavior. Since RDMAP directly or indirectly exposes memory addresses on the wire, the Placement information carried in each RDMA Message must be validated, including access rights and octet level granularity base and bounds check, before any data is Placed. For example, a third-party adversary could inject random packets that appear to be valid RDMA Messages and corrupt the memory on an RDMAP Data Sink. Since RDMAP is IP transport protocol independent, communication security mechanisms such as IPsec [IPSEC] or TLS [TLS] may be used to prevent such attacks.

10.2 Using IPsec with RDMAP

IPsec can be used to protect against the packet injection attacks outlined above. Because IPsec is designed to secure arbitrary IP packet streams, including streams where packets are lost, RDMAP can run on top of IPsec without any change. IPsec packets are processed (e.g., integrity checked and possibly decrypted) in the order they are received, and an RDMAP Data Sink will process the decrypted RDMA Messages contained in these packets in the same manner as RDMA Messages contained in unsecured IP packets.

10.3 Other Security Considerations

RDMAP has several mechanisms that deal with a number of attacks. These include, but are not limited to:

1. Connection to/from an unauthorized or unauthenticated endpoint.
2. Highjacking of an RDMAP Stream.
3. Attempts to read or write from unauthorized memory regions.
4. Injection of RDMA Messages within a Stream on a multi-user operating system by another application.

RDMAP relies on the LLP to establish the LLP Stream over which RDMA Messages will be carried. RDMAP itself does nothing to authenticate the validity of the LLP Stream of either of the endpoints. It is

the responsibility of the ULP to validate the LLP Stream. This is highly desirable due to the nature of RDMA and DDP.

Hijacking of an RDMAP channel would require that the underlying LLP connection be hijacked. This would require knowledge of Advertised buffers in order to directly Place data into a user buffer and is therefore constrained by the same techniques mentioned to guard against attempts to read or write from unauthorized memory regions.

RDMAP does not require a host to open its buffers to arbitrary attacks over the RDMAP Stream. It may access ULP memory only to the extent that the ULP has enabled and authorized it to do so. The STag access control model is defined by the RNIC Verbs Specification [VERBS]. Specific security operations include:

1. STags are only valid over the exact byte range established by the ULP. RDMAP MUST provide a mechanism for the ULP to establish and revoke the TO range associated with the ULP Buffer referenced by an STag.
2. STags are only valid for the duration established by the ULP. The ULP may revoke them at any time, in accordance with its own upper layer protocol requirements. RDMAP MUST provide a mechanism for the ULP to establish and revoke STag validity.
3. STags are only enabled for read and/or write access by explicit ULP action. RDMAP MUST provide a mechanism for the ULP to establish and revoke read, write, or read and write access to the ULP Buffer referenced by an STag.
4. The implementation is free to choose the value of STags and is encouraged to sparsely populate them over the full range available. This is admittedly weak security protection against a deliberate attack, but does minimize the risk of accidental matches when an incorrect STag is used due to a ULP software error.
5. RDMAP allows and encourages local interactions to restrict the usage of STags to specific Streams and/or user processes. RDMAP MUST provide a mechanism for associating a RDMAP Stream with a STag.
6. A ULP may only expose memory to remote access to the extent that it already had access to that memory itself.
7. RDMAP provides operations to allow the holder of an STag to indicate when it has made its last usage of that STag. This enables automatic deregistration and/or scope reduction of STags as the implementation and ULP may see fit.

8. If an STag is not valid on a connection, RDMAP provides a mechanism for terminating the RDMAP Stream (see section 6.8 Terminate Header).
9. An STag that is associated with an RDMAP Stream becomes invalid upon reception of a valid Send with Invalidate or Send with Solicited Event Message. RDMAP MUST invalidate the STag sent in a valid Send with Invalidate or Send with Solicited Event and Invalidate Message, before Completing the Send with Invalidate or Send with Solicited Event and Invalidate Message.

Further, RDMAP encourages direct Placement of incoming payloads in user-mode ULP Buffers. This avoids the risks of prior solutions that relied upon exposing system buffers for incoming payloads.

There is also a clean data Delivery hand-off between RDMAP and the ULP. This allows the ULP to implement additional security operations without restrictions or interference from RDMAP.

In summary, RDMAP enables both ULP and LLP security. It requires that all of its data access be enabled and authorized by the ULP. It provides no operations for the ULP to gain permissions not already granted by the host operating system. It allows and encourages local interactions to specify even more precise security checks on STag binding and data transfer operations.

By remaining independent of ULP and LLP security protocols, RDMAP will benefit from continuing improvements at those layers. Users are provided flexibility to adapt to their specific security requirements and the ability to adapt to future security challenges.

11 References

11.1 Normative References

[VERBS] Specification under development in the RDMA Consortium (see <http://www.rdmaconsortium.org/>).

[DDP] H. Shah et al., "Direct Data Placement over Reliable Transports", RDMA Consortium Draft Specification draft-shah-iwarp-ddp-v1.0, October 2002 (see <http://www.rdmaconsortium.org/>)

[MPA] P. Culley et al., "Markers with PDU Alignment", RDMA Consortium Draft Specification draft-culley-iwarp-mpa-v1.0, October 2002 (see <http://www.rdmaconsortium.org/>)

[SCTP] R. Stewart et al., "Stream Control Transmission Protocol", RFC 2960, October 2000.

[TCP] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, September 1981.

11.2 Informative References

[RFC2401] Atkinson, R., Kent, S., "Security Architecture for the Internet Protocol", RFC 2401, November 1998.

[TLS] Dierks, T. and C. Allen, "The TLS Protocol Version 1.0", RFC 2246, November 1998.

12 Appendix

12.1 DDP Segment Formats for RDMA Messages

This appendix is for information only and is NOT part of the standard. It simply depicts the DDP Segment format for the various RDMA Messages.

12.1.1 DDP Segment for RDMA Write

The following figure depicts an RDMA Write, DDP Segment:

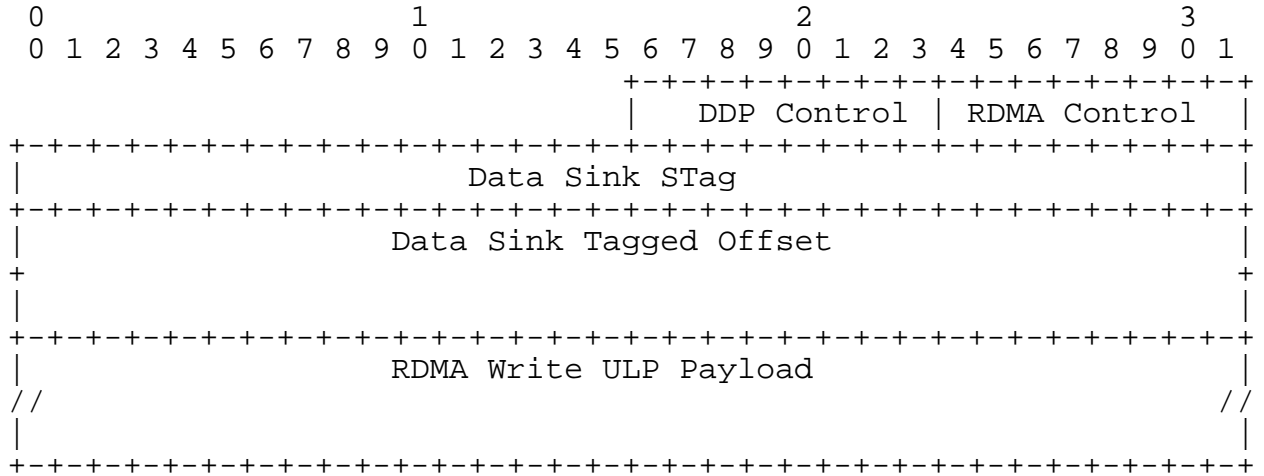


Figure 11 RDMA Write, DDP Segment format

12.1.2 DDP Segment for RDMA Read Request

The following figure depicts an RDMA Read Request, DDP Segment:

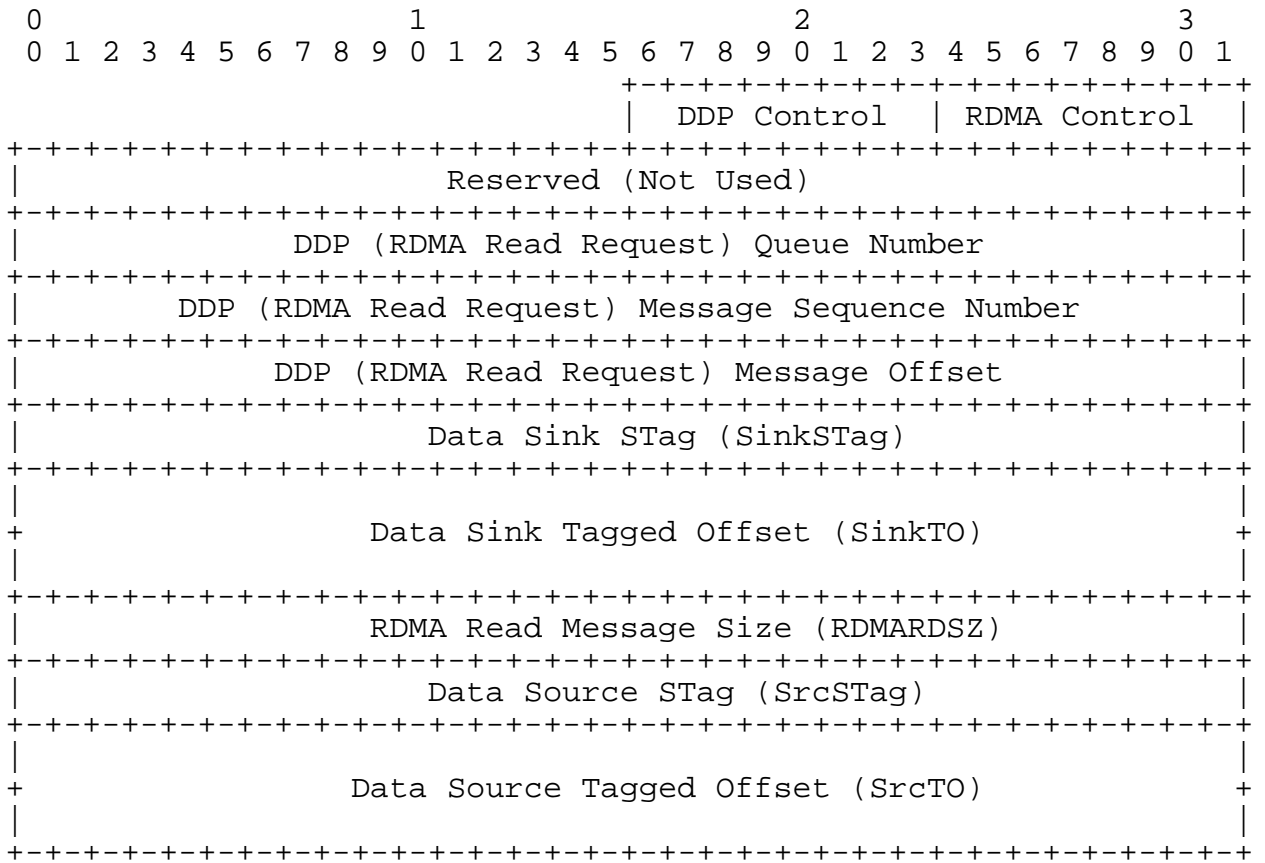


Figure 12 RDMA Read Request, DDP Segment format

12.1.3 DDP Segment for RDMA Read Response

The following figure depicts an RDMA Read Response, DDP Segment:

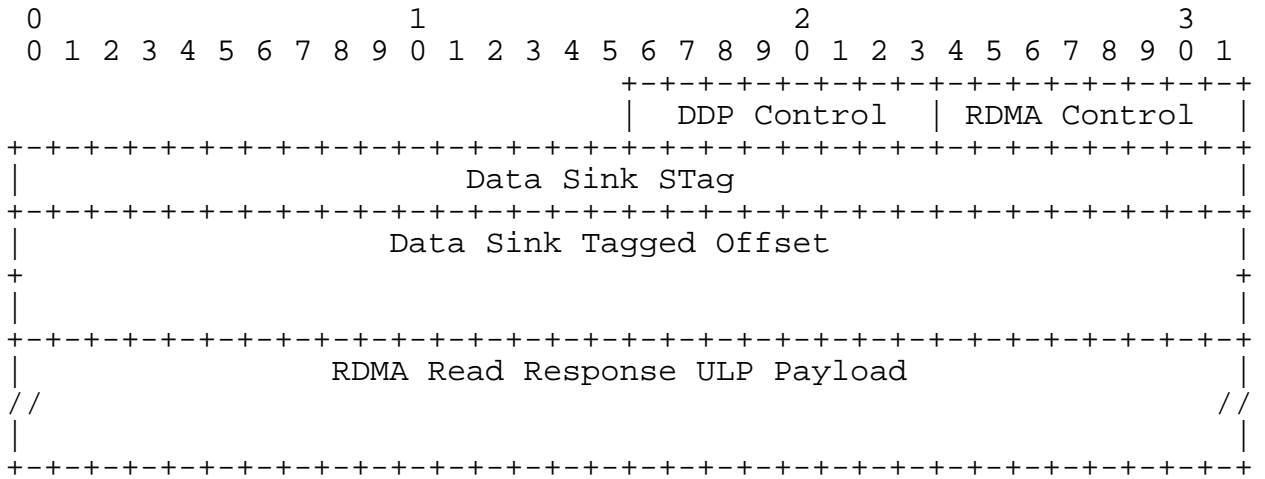


Figure 13 RDMA Read Response, DDP Segment format

12.1.4 DDP Segment for Send and Send with Solicited Event

The following figure depicts a Send and Send with Solicited Request, DDP Segment:

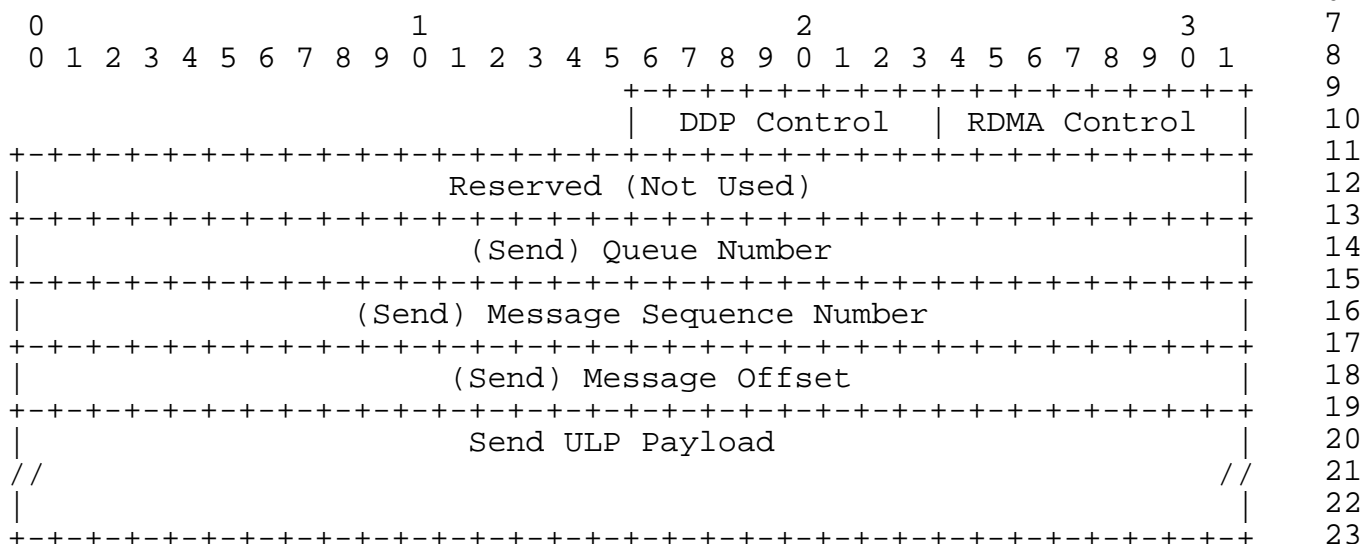


Figure 14 Send and Send with Solicited Event, DDP Segment format

12.1.5 DDP Segment for Send with Invalidate and Send with SE and Invalidate

The following figure depicts a Send with invalidate and Send with Solicited and Invalidate Request, DDP Segment:

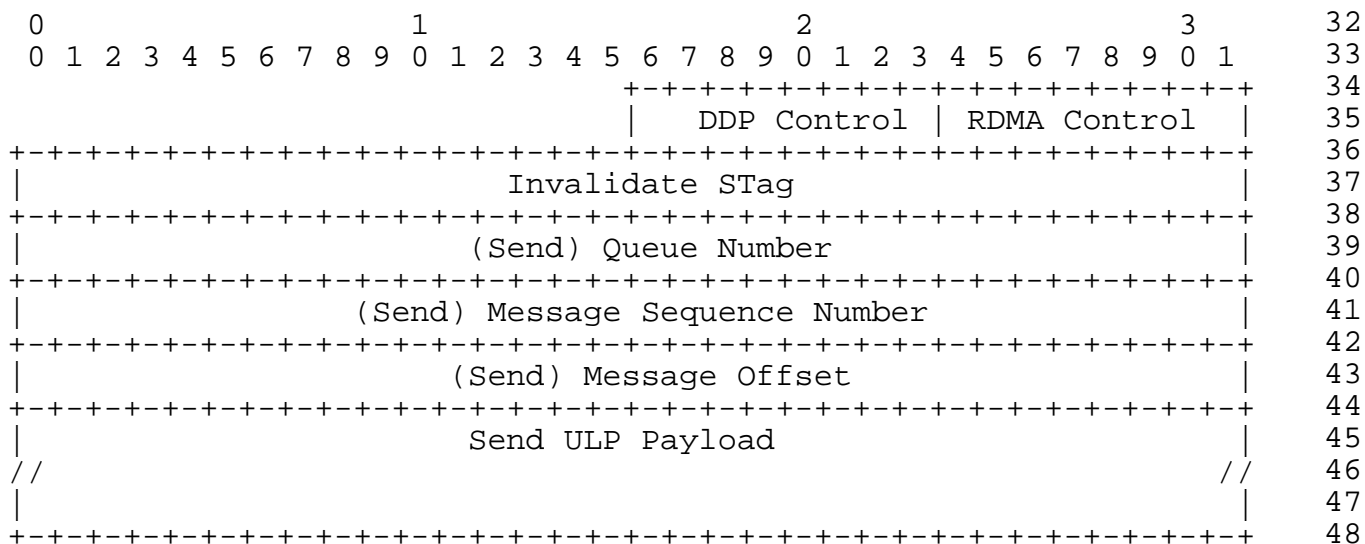


Figure 15 Send with Invalidate and Send with SE and Invalidate, DDP Segment

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

12.1.6 DDP Segment for Terminate

The following figure depicts a Terminate, DDP Segment:

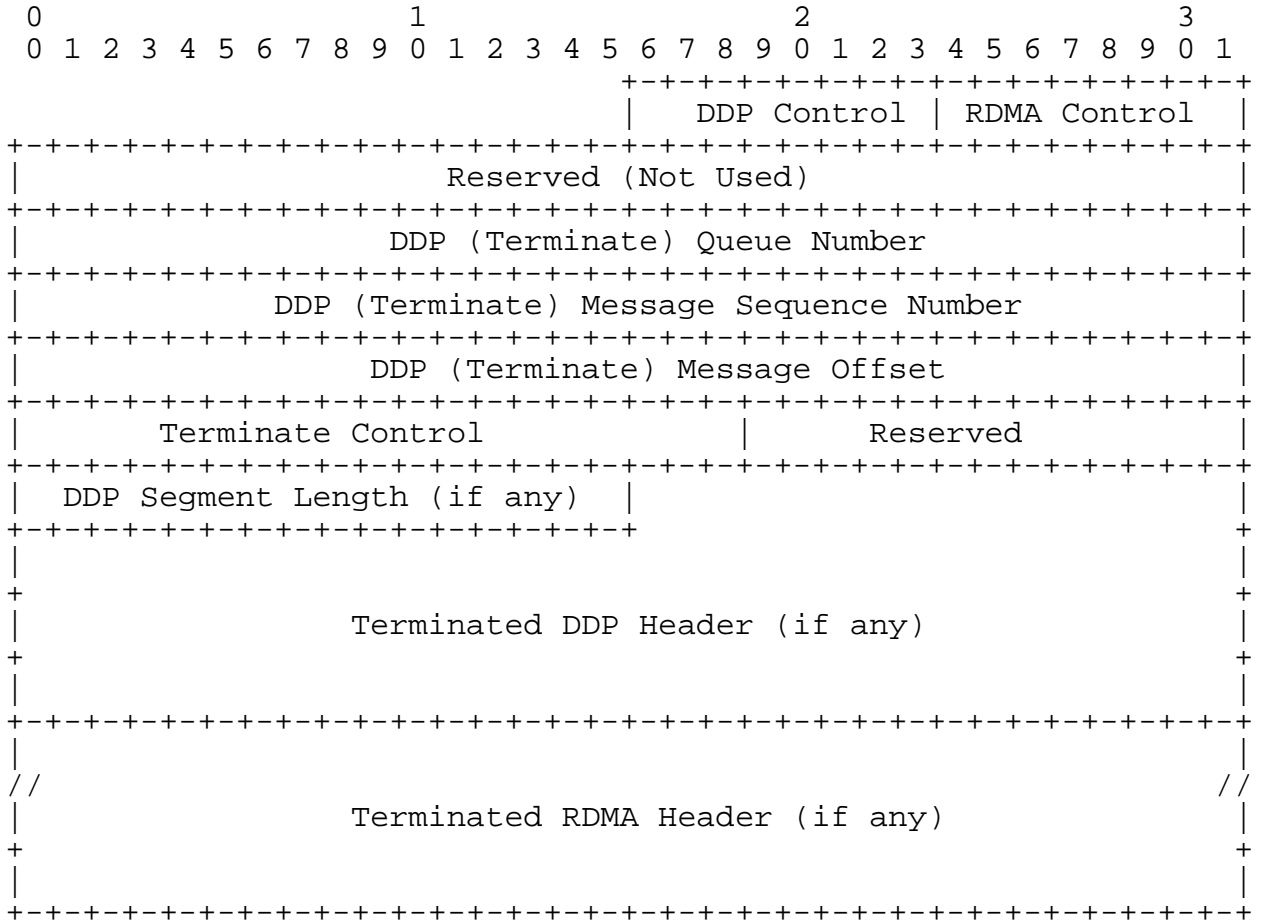


Figure 16 Terminate, DDP Segment format

12.2 Ordering and Completion Table

The following table summarizes the ordering relationships that are defined in section 7.5 Ordering and Completions from the standpoint of the local peer issuing the two Operations. Note, in the table that follows Send includes Send, Send with Invalidate, Send with Solicited Event, and Send with Solicited Event and Invalidate

First Op	Later Op	Placement guarantee at Remote Peer	Placement guarantee Local Peer	Ordering guarantee at Remote Peer
Send	Send	No placement	Not applicable	Completed in

		guarantee. If guarantee is necessary, see footnote 1.		order.	1 2 3 4 5 6
Send	RDMA Write	No placement guarantee. If guarantee is necessary, see footnote 1.	Not applicable	Not applicable	7 8 9 10 11 12
Send	RDMA Read	No placement guarantee between Send Payload and RDMA Read Request Header	RDMA Read Response Payload will not be placed at the local peer until the Send Payload is placed at the remote peer	RDMA Read Response Message will not be generated until Send has been Completed	13 14 15 16 17 18 19 20 21 22
RDMA Write	Send	No placement guarantee. If guarantee is necessary, see footnote 1.	Not applicable	Not applicable	23 24 25 26 27 28
RDMA Write	RDMA Write	No placement guarantee. If guarantee is necessary, see footnote 1.	Not applicable	Not applicable	29 30 31 32 33 34
RDMA Write	RDMA Read	No placement guarantee between RDMA Write Payload and RDMA Read Request Header	RDMA Read Response Payload will not be placed at the local peer until the RDMA Write Payload is placed at the remote peer	Not applicable	35 36 37 38 39 40 41 42 43 44 45
RDMA Read	Send	No placement guarantee between RDMA Read Request Header and Send payload	Send Payload may be placed at the remote peer before the RDMA Read Response is	Not applicable	46 47 48 49 50 51

			generated. If guarantee is necessary, see footnote 2.	
RDMA Read	RDMA Write	No placement guarantee between RDMA Read Request Header and RDMA Write payload	RDMA Write Payload may be placed at the remote peer before the RDMA Read Response is generated. If guarantee is necessary, see footnote 2.	Not applicable
RDMA Read	RDMA Read	No placement guarantee of the two RDMA Read Request Headers Additionally, there is no guarantee that the Tagged Buffers referenced in the RDMA Read will be read in order	No placement guarantee of the two RDMA Read Response Payloads.	Second RDMA Read Response will not be generated until first RDMA Read Response is generated.

Figure 17 Operation Ordering

Footnote 1: If the guarantee is necessary, a ULP may insert an RDMA Read Operation and wait for it to complete to act as a Fence.

Footnote 2: If the guarantee is necessary, a ULP may wait for the RDMA Read Operation to complete before performing the Send.

13 Authors Addresses

Paul R. Culley
Hewlett-Packard Company
20555 SH 249
Houston, Tx. USA 77070-2698
Phone: 281-514-5543
Email: paul.culley@hp.com

Dave Garcia
Hewlett-Packard Company
19333 Vallco Parkway
Cupertino, Ca. USA 95014
Phone: 408.285.6116
Email: dave.garcia@hp.com

Jeff Hilland
Hewlett-Packard Company
20555 SH 249
Houston, Tx. USA 77070-2698
Phone: 281-514-9489
Email: jeff.hilland@hp.com

Renato J. Recio
IBM Corp.
11501 Burnett Road
Austin, Tx. USA 78758
Phone: 512-838-3685
Email: recio@us.ibm.com

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

14 Acknowledgments

Dwight Barron

Hewlett-Packard Company
20555 SH 249
Houston, Tx. USA 77070-2698
Phone: 281-514-2769
Email: dwight.barron@compaq.com

John Carrier

Adaptec, Inc.
691 S. Milpitas Blvd.
Milpitas, CA 95035 USA
Phone: +1 (360) 378-8526
Email: john_carrier@adaptec.com

Ted Compton

EMC Corporation
Research Triangle Park, NC 27709, USA
Phone: 919-248-6075
Email: compton_ted@emc.com

Uri Elzur

Broadcom Corporation
16215 Alton Parkway
Irvine, California 92619-7013 USA
Phone: +1 (949) 585-6432
Email: Uri@Broadcom.com

Hari Ghadia

Adaptec, Inc.
691 S. Milpitas Blvd.,
Milpitas, CA 95035 USA
Phone: +1 (408) 957-5608
Email: hari_ghadia@adaptec.com

Howard C. Herbert

Intel Corporation
MS CH7-404
5000 West Chandler Blvd.
Chandler, Arizona 85226
Phone: 480-554-3116
Email: howard.c.herbert@intel.com

Mike Ko

IBM
650 Harry Rd.
San Jose, CA 95120
Phone: (408) 927-2085
Email: mako@us.ibm.com

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

Mike Krause
Hewlett-Packard Company
43LN
19410 Homestead Road
Cupertino, CA 95014 USA
Phone: 408-447-3191
Email: krause@cup.hp.com

Dave Minturn
Intel Corporation
MS JF1-210
5200 North East Elam Young Parkway
Hillsboro, Oregon 97124
Phone: 503-712-4106
Email: dave.b.minturn@intel.com

Mike Penna
Broadcom Corporation
16215 Alton Parkway
Irvine, California 92619-7013 USA
Phone: +1 (949) 926-7149
Email: MPenna@Broadcom.com

Jim Pinkerton
Microsoft, Inc.
One Microsoft Way
Redmond, WA, USA 98052
Email: jpink@microsoft.com

Hemal Shah
Intel Corporation
MS PTL1
1501 South Mopac Expressway, #400
Austin, Texas 78746
Phone: 512-732-3963
Email: hemal.shah@intel.com

Allyn Romanow
Cisco Systems
170 W Tasman Drive
San Jose, CA 95134 USA
Phone: +1 408 525 8836
Email: allyn@cisco.com

Tom Talpey
Network Appliance
375 Totten Pond Road
Waltham, MA 02451 USA
Phone: +1 (781) 768-5329
EMail: thomas.talpey@netapp.com

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

Patricia Thaler
Agilent Technologies, Inc.
1101 Creekside Ridge Drive, #100
M/S-RG10
Roseville, CA 95678
Phone: +1-916-788-5662
email: pat_thaler@agilent.com

Jim Wendt
Hewlett-Packard Company
8000 Foothills Boulevard MS 5668
Roseville, CA 95747-5668 USA
Phone: +1 916 785 5198
Email: jim_wendt@hp.com

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

15 Full Copyright Statement

This document and the information contained herein is provided on an "AS IS" basis and ADAPTEC INC., AGILENT TECHNOLOGIES INC., BROADCOM CORPORATION, CISCO SYSTEMS INC., EMC CORPORATION, HEWLETT-PACKARD COMPANY, INTERNATIONAL BUSINESS MACHINES CORPORATION, INTEL CORPORATION, MICROSOFT CORPORATION, AND NETWORK APPLIANCE INC. DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Copyright (c) 2002 ADAPTEC INC., BROADCOM CORPORATION, CISCO SYSTEMS INC., EMC CORPORATION, HEWLETT-PACKARD COMPANY, INTERNATIONAL BUSINESS MACHINES CORPORATION, INTEL CORPORATION, MICROSOFT CORPORATION, NETWORK APPLIANCE INC., All Rights Reserved.