



OpenCAPI 3.0

Transaction Layer Specification

Version 1.0
28 January 2020

Approved
Approved for Distribution

Transaction Layer Specification

TL Specification Working Group
OpenCAPI Consortium

Version 1.0 (28 January 2020)

Copyright © OpenCAPI Consortium 2016-2020.

Printed in the United States of America April 7, 2020 (1.9.1.3p0.002).

Use of this document (including for investigative purposes, research, development, productizing, etc.) is controlled by the OpenCAPI License Agreement, which is available at <https://www.opencapi.org>

The owners of all copyrights in this document have agreed to make them available to non-OpenCAPI members under the OpenCAPI Study Only License to the same extent as if they were OpenCAPI members. The license is found here: <https://opencapi.org/license/>.

All capitalized terms in the following text have the meanings assigned to them in the OpenCAPI Intellectual Property Rights Policy (the “OpenCAPI IPR Policy”). The full Policy may be found at the OpenCAPI website.

This document and the information contained herein is provided on an “AS IS” basis AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, THE OPENCAP CONSORTIUM AS WELL AS THE AUTHORS AND DEVELOPERS OF THIS DRAFT STANDARD OR OTHER DOCUMENT HEREBY DISCLAIM ALL OTHER WARRANTIES AND CONDITIONS, EITHER EXPRESS, IMPLIED OR STATUTORY, INCLUDING BUT NOT LIMITED TO, ANY IMPLIED WARRANTIES, DUTIES OR CONDITIONS OF MERCHANTABILITY, OF FITNESS FOR A PARTICULAR PURPOSE, OF ACCURACY OR COMPLETENESS OF RESPONSES, OF RESULTS, OF WORKMANLIKE EFFORT, OF LACK OF VIRUSES, OF LACK OF NEGLIGENCE OR NON-INFRINGEMENT.

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

The limited permissions granted above are perpetual and will not be revoked by OpenCAPI or its successors or assigns while the specification in question is the current version. Upon release by the OpenCAPI Consortium of a new version, all the above rights shall automatically cease.

OpenCAPI and the OpenCAPI logo design are trademarks of the OpenCAPI Consortium.

Other company, product, and service names may be trademarks or service marks of others.

Abstract

This document details the OpenCAPI TL specification. It is the work product of the OpenCAPI Consortium TL Specification Work Group.

This document is handled in compliance with the requirements outlined in the OpenCAPI Consortium Work Group (WG) process document. Comments, questions, etc. can be submitted to membership@opencapi.org.

Approved

Participants

Brian Allison, IBM, *Chair*

Michael Siegel, IBM, *Technical Editor*

Joe Breher, Western Digital Technologies,
Inc

Harold Dozier, Micron Technology, Inc

Mark Fredrickson, IBM

Rick Hagen, NVIDIA Corporation

Paul Hartke, Xilinx, Inc

Curt Wollbrink, IBM

Contents

List of figures	7
List of tables	8
Revision log	10
About this document	11
Architecture compliance terminology	11
Conventions used in this specification	11
Bit and byte numbering	11
Representation of numbers	12
RTL notation	12
Notes	13
Engineering notes	13
Developer notes	13
Command flows and transaction diagrams	14
Command flow diagrams	14
Transaction diagrams	14
Terms	17
1. Overview	23
1.1 OpenCAPI protocol stack	24
1.2 Host operation modes	25
1.2.1 No attached device (C0, M0)	25
1.2.2 MEM-only mode (C0, M1)	25
1.2.3 Checkout mode (C1, M0)	26
1.2.4 Checkout with MEM (C1, M1)	26
1.3 Command ordering	26
1.4 Write fragmentation ordering and atomicity	26
1.4.1 Write fragmentation ordering and atomicity at the host	26
1.4.1.1 Partial write operations	26
1.4.1.2 64-, 128-, 256-byte write operations	26
1.4.2 Write fragmentation ordering and atomicity at the AFU	27
1.4.2.1 Partial write operations	27
1.4.2.2 64-, 128-, 256-byte write operations	27
1.5 OpenCAPI device PA space specification	27
1.5.1 PA-to-RA mapping rules	28
1.6 Address translation	28
1.6.1 Effective to real address translation	28
2. TL and TLX command and response specifications	30
2.1 Handling multiple responses to a single command	33
2.1.1 TLX Read request getting multiple TL responses	33
2.1.2 TLX Write request getting multiple TL responses	34

Approved

2.1.3 TL read request getting multiple TLX responses	34
2.1.4 TL write request getting multiple TLX responses	35
2.2 TL CAPP command packets	36
2.3 TLX AP command packets	46
2.4 TL CAPP response packets	66
2.5 TLX AP response packets	77
3. Virtual channel and data credit pool specification	84
3.1 Virtual channel	85
3.1.1 TLX command and response VC (TLX.vc)	85
3.1.2 TL command and response VC (TL.vc)	85
3.1.3 VC credit count specification	86
3.2 Data credit pool	86
3.2.1 TLX data DCP (TLX.dcp)	86
3.2.2 TL data DCP (TL.dcp)	86
3.2.3 DCP credit count specification	87
3.3 TL Virtual channel and service queues	88
3.3.1 Host TLX command handling	88
3.3.2 Host TLX response handling	90
3.4 Device TL virtual channel queues	90
3.5 Virtual channel dependency rules	91
4. The acTag table	93
4.1 acTag table contents	93
4.2 acTag table access	93
4.2.1 Error cases when accessing the acTag table	93
4.3 acTag entry management	93
5. DL frame format	95
5.1 DL frame control flit (64 bytes)	96
5.1.1 DL content	96
5.1.2 TL command/response content	97
5.1.3 Data transport, order, and alignment	97
6. TL and TLX template specifications	99
6.1 TLX receive and TL transmit template capability specification	100
6.2 TL receive and TLX transmit template capability specification	101
6.3 Control-flit rate capability	101
7. Error detection	103
7.1 Error events	104
8. OpenCAPI profiles	110
Appendix A. AP (TLX) command transaction diagrams	115
A.1 AFU read with no intent to cache; 128 bytes	116
A.2 AFU DMA write; 128 bytes	117

Approved

A.3 AFU DMA partial write; 8 bytes	119
Appendix B. CAPP (TL) command transaction diagrams	120
B.1 CAPP memory read; 128 bytes	120
B.2 CAPP memory write; 128 bytes	121

List of figures

Figure 1.	Big- and little-endian comparisons	12
Figure 2.	Command flow example	14
Figure 3.	Example TLX and TL transaction diagram	16
Figure 1-1.	OpenCAPI stack	24
Figure 2-1.	Address translation sequence: xlate_touch	62
Figure 3-1.	TL command flow from the VC queue to the service queue	89
Figure 3-2.	TLX command and response flow from the VC to the AFU protocol stack	91
Figure 3-3.	VC dependency graph	92
Figure A-1.	TLX and TL interaction: rd_wnitc	116
Figure A-2.	TLX and TL interaction: dma_w	117
Figure A-3.	TL and TLX interaction: dma_pr_w	119
Figure B-1.	TL and TLX transaction: rd_mem	120
Figure B-2.	TL and TLX transaction: write_mem	121

List of tables

Table 1.	Architecture terms	11
Table 2-1.	TL and TLX command operands	30
Table 2-2.	The Resp_code specification for xlite_done	37
Table 2-3.	The Resp_code specification for intrap_rdy	38
Table 2-4.	The cmd_flag specification for amo_rd	52
Table 2-5.	The cmd_flag specification for amo_rw	53
Table 2-6.	The cmd_flag specification for amo_w	55
Table 2-7.	The cmd_flag specification for xlite_touch (all forms)	60
Table 2-8.	The Resp_code specification for touch_resp	67
Table 2-9.	touch_resp Resp_code use by TLX command	68
Table 2-10.	The Resp_code specification for read_failed	69
Table 2-11.	read_failed Resp_code use by TLX command	70
Table 2-12.	The Resp_code specification of write_failed	72
Table 2-13.	write_failed Resp_code use by TLX command	74
Table 2-14.	The Resp_code specification for intrap_resp	74
Table 2-15.	intrap_resp Resp_code use by TLX command	75
Table 2-16.	The Resp_code specification for wake_host_resp	76
Table 2-17.	The Resp_code specification for mem_rd_fail	78
Table 2-18.	mem_rd_fail Resp_code use by TL command	79
Table 2-19.	The Resp_code specification for mem_wr_fail	81
Table 2-20.	mem_wr_fail Resp_code use by TL command	82
Table 3-1.	VC maximum credit count specification	86
Table 3-2.	DCP maximum credit count specification	87
Table 3-3.	Summary VC and DCP assignments	87
Table 5-1.	DL frame format showing CRC and “bad data flit” coverage	95
Table 6-1.	Template capability definitions	100
Table 6-2.	TLX receive/TL transmit template	100
Table 6-3.	TL receive/TLX transmit template	101
Table 7-1.	Error event specification	104
Table 8-1.	Feature compliance requirement notation	110
Table 8-2.	Profile specifications for TL commands	111
Table 8-3.	Profile specifications for TLX commands	111
Table 8-4.	Profile specifications for TL responses	112
Table 8-5.	Profile specifications for TLX responses	112
Table 8-6.	Profile specifications for TLX receive/TL transmit templates	113
Table 8-7.	Profile specifications for TL receive/TLX transmit templates	113
Table 8-8.	Profile specifications host operation modes	113
Table 8-9.	Profile specifications supported page size	114

Approved

Table 8-10.	Profile specifications supported dLength by TLX	114
Table 8-11.	Profile specifications supported dLength by TL	114
Table 8-12.	Profile specifications support of endianness data format by the TL	114

Approved

Revision log

Each release of this document supersedes all previously released versions. The revision log lists all significant changes made to the document since its initial release. In the rest of the document, change bars in the margin indicate that the adjacent text was modified from the previous release of this document.

Revision date	Description
28 January 2020	Release of Approved OpenCAPI TL 3.0 specification.

About this document

This document provides the architectural specification of the OpenCAPI™ transaction layer (TL and TLX).

Architecture compliance terminology

In architecture descriptions, certain terms carry meaning in addition to their normal use in English. The following terms are used within this architecture specification to describe the requirements an implementation must meet to be considered compliant.

Table 1. Architecture terms

Term	Description
invalid	Used for multi-bit fields where the contents are not reliable. The field or bus shall not be examined for any functional or error checking actions.
may	An architectural option indicating that an implementation is allowed to have this behavior or characteristic.
reserved	With respect to a field of a register or bus: <ul style="list-style-type: none"> • A reserved field shall be set to 0 by an implementation. • A reserved field shall not be examined by an implementation. With respect to a code point: <ul style="list-style-type: none"> • A reserved code point shall not be issued by a compliant implementation • A reserved code point shall cause a bounded undefined response (that is, it won't hang the system). • A reserved code point may be used in future revisions of the architecture. The architecture may specify that the use of a reserved code point is an error condition.
shall	An architectural requirement indicating a required behavior or characteristic.
uncertain	Used for single-bit fields where the contents are not reliable. The field or bus shall not be examined for any functional or error checking actions.
undefined	When the value of a field or a bus is undefined, the value may vary between implementations and may vary for a particular implementation for different actions. An implementation shall not examine a field when its value is undefined for functional purposes. However, the field may be checked for errors in those cases where an implementation includes error checking (that is, parity, ECC and so on)

Conventions used in this specification

Bit and byte numbering

Throughout this document, little-endian notation is used, which means that bits and bytes are numbered in descending order from left to right.

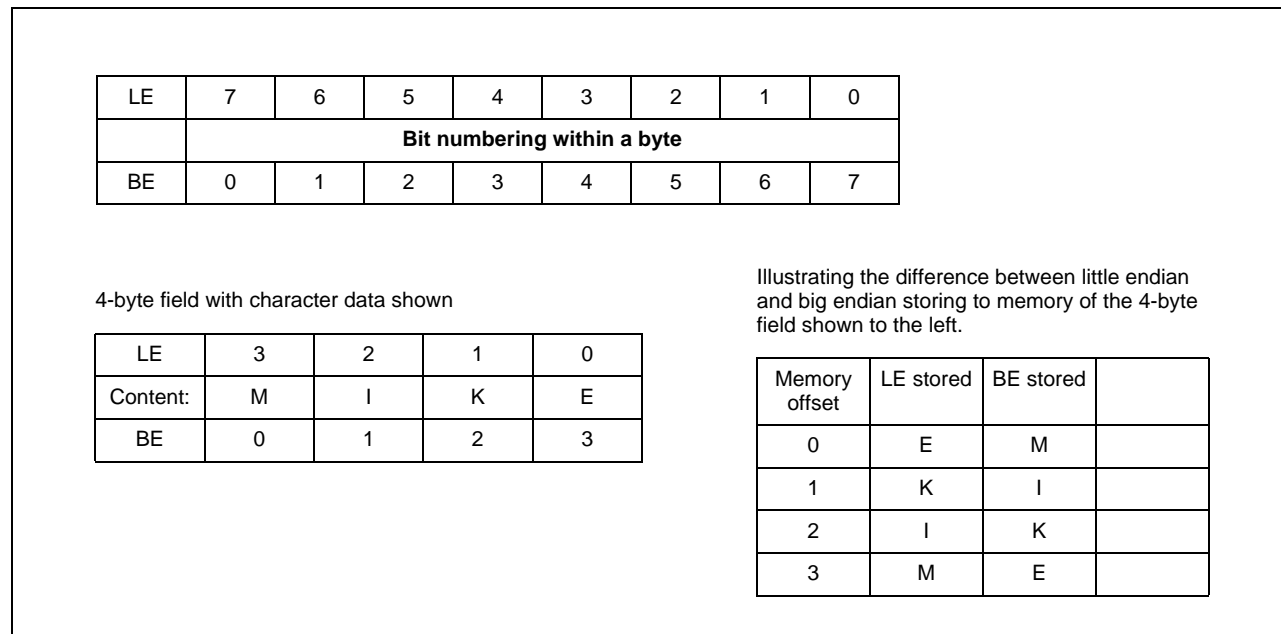
Thus, in the description of a 4-byte field, bit 31 is the most significant bit (MSb) and bit 0 is the least significant bit (LSb). The corresponding byte numbering is also shown.

MSb	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	LSb
byte 3								byte 2								byte 1								byte 0									

Approved

The big-endian and little-endian byte ordering are described in the *POWER ISA, version 3.0, Book I. Figure 1* compares big-endian and little-endian notation.

Figure 1. Big- and little-endian comparisons



Representation of numbers

The notation for bit encoding is as follows:

- Hexadecimal values are preceded by an x and enclosed in single quotation marks. For example x'0A00'. Bit numbering is little endian and, in this example, is 15 to 0.
- Binary values in sentences are shown in single quotation marks. For example '1010'. Bit numbering is little endian and, in this example, is 3 to 0.
- ⁿx means the replication of x, n times. That is, x is concatenated to itself n-1 times. ⁿ0 and ⁿ1 are special cases:
 - ⁿ0 means a field of n bits with each bit equal to 0. For example, ⁵0 is equivalent to '00000'.
 - ⁿ1 means a field of n bits with each bit equal to 1. For example, ⁵1 is equivalent to '11111'.

RTL notation

RTL notations are used to specify the architectural transformation performed by the execution of a command.

Notation	Meaning
←	Assignment.
	Concatenation.
=, ≠	Equal, not equal relations.
≥, ≤	Greater than or equal to, less than or equal to relations.

Approved

Notation	Meaning
>, <	Greater than or less than relations.
+	Two's complement addition.
-	Two's complement subtraction, unary minus
∨	Bitwise logical OR
∧	Bitwise logical AND
⊕	Bitwise logical exclusive OR
Max(x,y)	Returns x when $x \geq y$; otherwise returns y
Min(x,y)	Returns x when $x \leq y$; otherwise returns y.
{x...y}	All integer values from x through y.
A = {x...y}	Returns true when A is a member of the set of integer values in the range of x through y.

Notes

This document contains engineering and developer notes.

Engineering notes

Engineering notes provide additional implementation details and recommendations not found elsewhere. The notes might include architectural compliance requirements. That is, the text might include *Architecture compliance terminology*. These notes should be read by all implementation and verification teams to ensure architectural compliance.

Engineering note

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Proin cursus hendrerit enim, vel tempus nibh ornare ut. Quisque ac augue eu augue convallis hendrerit. Mauris iaculis viverra ipsum nec dapibus. Nunc at porta libero. Curabitur luctus ultrices augue non pulvinar. Vestibulum mattis non ipsum at venenatis. Suspendisse euismod, neque et suscipit luctus, odio metus semper lectus, quis volutpat est libero quis nunc. Vivamus rutrum mauris sed tristique malesuada. Vivamus at augue vitae nisl cursus feugiat. Pellentesque efficitur sed nisi in dapibus. Curabitur vestibulum cursus arcu, ut mattis nisl.

Developer notes

Developer notes are used to document the reasoning and discussions that led to the current version of the architecture. These notes might also include recommended changes for future versions of the architecture, or warnings of approaches that have failed in the past. These notes should be read by verification teams and contributors to the architecture.

Developer note

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Proin cursus hendrerit enim, vel tempus nibh ornare ut. Quisque ac augue eu augue convallis hendrerit. Mauris iaculis viverra ipsum nec dapibus. Nunc at porta libero. Curabitur luctus ultrices augue non pulvinar. Vestibulum mattis non ipsum at venenatis. Suspendisse euismod, neque et suscipit luctus, odio metus semper lectus, quis volutpat est libero quis nunc. Vivamus rutrum mauris sed tristique malesuada. Vivamus at augue vitae nisl cursus feugiat. Pellentesque efficitur sed nisi in dapibus. Curabitur vestibulum cursus arcu, ut mattis nisl.

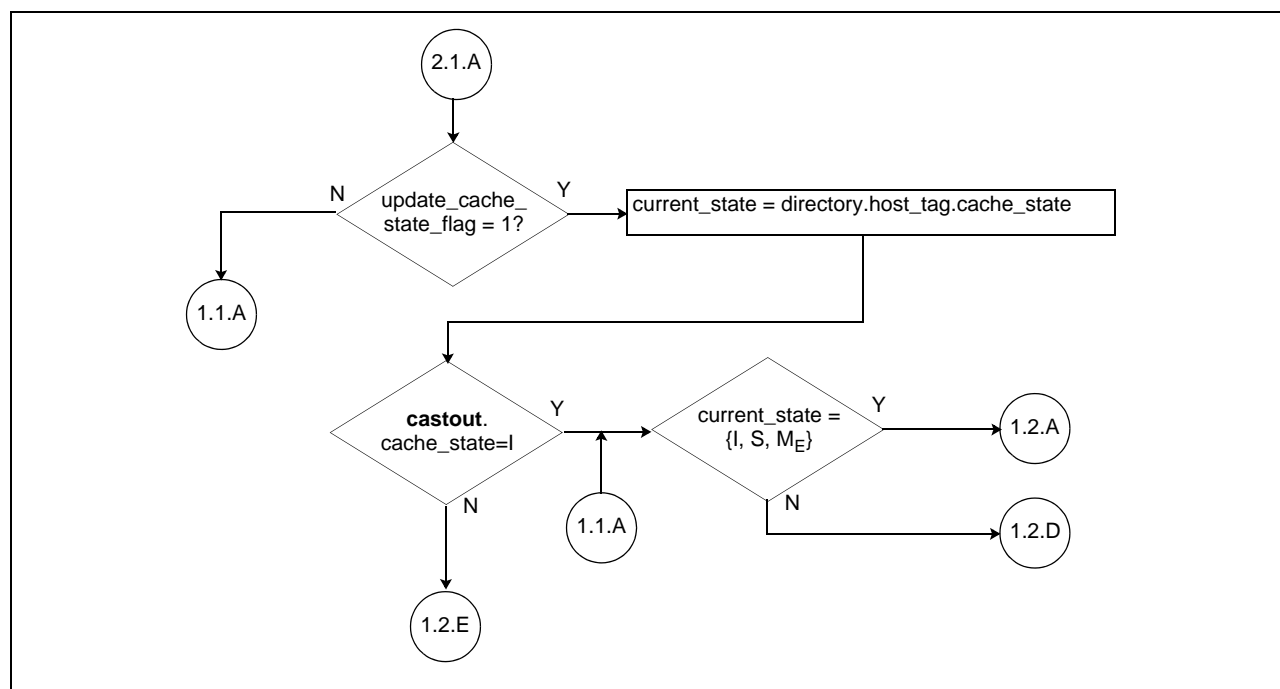
Command flows and transaction diagrams

Command flow diagrams

Command-flow diagrams show interactions within and across the different levels of the OpenCAPI protocol stack. Command flows use diamonds for decision blocks and rectangles for actions taken. Circles are used for on-page and off-page connectors and indicate a from-to direction based on the text content of the circle.

In *Figure 2*, a simple decision block with a state change and an off-page connector is shown. The text within the off-page connector has the format of “source page”.destination page”.”instance”. The off-page connectors shown in the figure is on page 1 of the figure¹ and is connecting to page 2 of the figure. On page 2, identical off-page connectors can be found. The instance indication allows for multiple connections to be shown between two pages. Connector 2.1.A illustrates a connection from page 2 to page 1 of *Figure 2*. An off-page connector can also be used to “connect” two spots on the same page as illustrated by connector 1.1.A. The direction of the arrow, into or out of a connector, decision block, or assignment-action block, indicates the direction of the sequence within the flow diagram.

Figure 2. Command flow example



Transaction diagrams

Transaction diagrams show the interaction between the TL and TLX layers and provide some illustrative notes for actions taken at the host protocol layer and the attached functional unit (AFU) protocol layer. In *Figure 3* on page 16, the diagram is broken into three vertical sections. From left to right, these are the AFU protocol layer notes, transactions between the TL and TLX layers, which are typically command and response packets, and the host protocol layer notes. Arrows indicate the direction in which the packet or action flows; for example, towards or away from the host (TL) layer.

1. All multi-page figures contain a “page n of y” notation in the figure description.

Approved

Circles are used for on-page and off-page connectors and indicate a from-to direction based on the text content of the circle. The text content is specified in the same manner for transaction diagrams as previously described for command flows. In addition to the specification of how connectors are used in command flows, in transaction diagrams, when a connector is used without an arrow, the transaction shown is one of multiple possible transaction outcomes. The use of this technique reduces the size of the transaction figure because the preceding set of transactions do not have to be repeated.

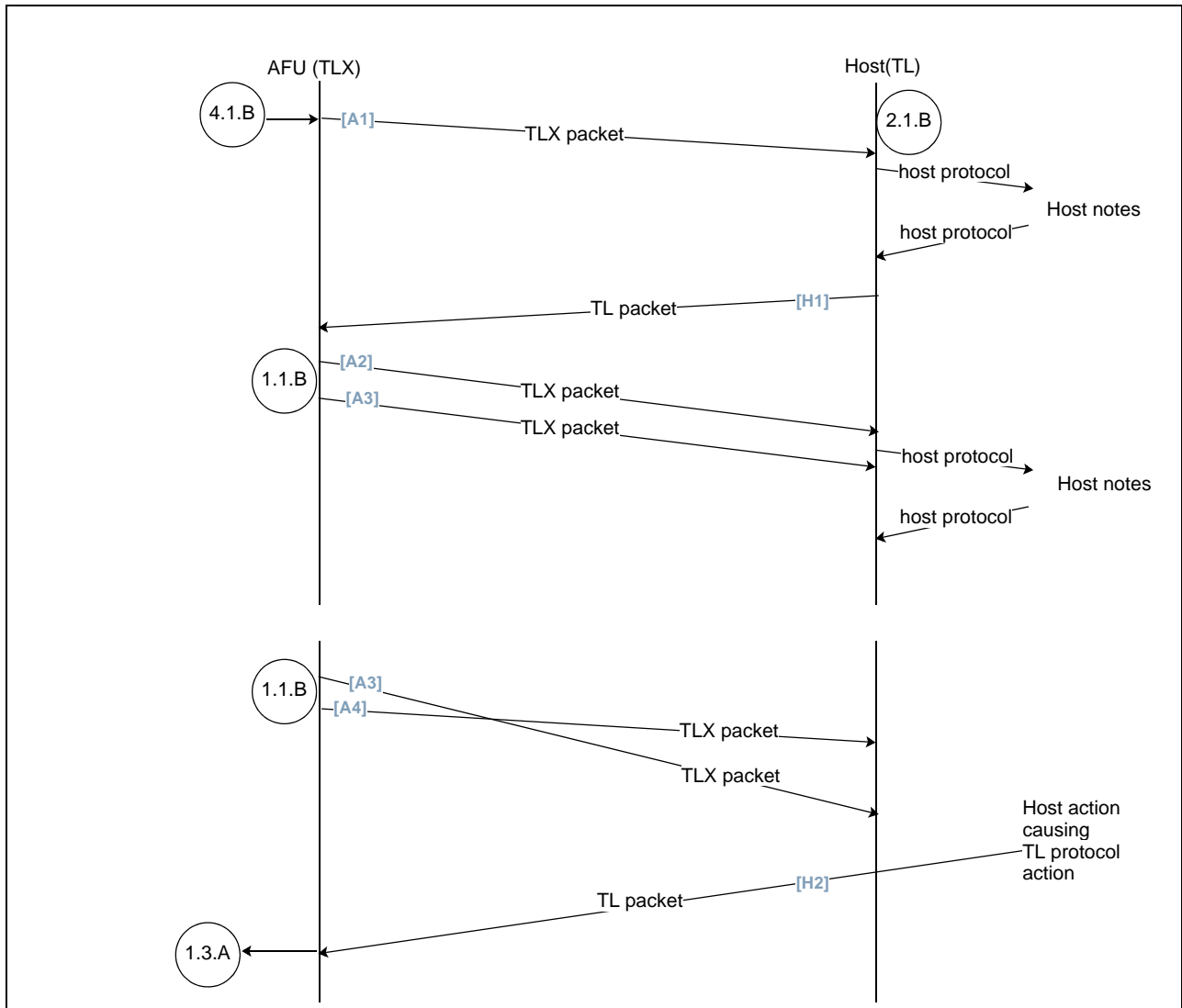
In *Figure 3*, connector 1.1.B illustrates an on-page connection without an arrow to indicate a different transaction outcome. The prior events are assumed to have occurred when looking at the second instance of the 1.1.B connector. In the second case, one TLX packet has passed a previously issued TLX packet; this is something that can occur when two packets use different virtual channels. Connector 1.3.A shows an off-page connection to page 3, and connector 4.1.B shows an off-page connection from page 4.

Arrow numbering is included in transaction diagrams to simplify references to transactions. The form of the arrow references indicates the source of the transaction (AFU or Host) and the instance of the arrow. As seen in *Figure 3*, [A1] is the first arrow from the TLX packet transaction and [H1] is the first TL transaction.

A break in the vertical lines indicates where a new transaction illustration starts or ends.

Approved

Figure 3. Example TLX and TL transaction diagram



Terms

The following terms are used in this document.

{EA, address context}	<p>A short hand notation to indicate an EA and <i>address context</i> pair.</p> <ul style="list-style-type: none"> Used when specifying an entry in an AFU L1 directory Use in discussions about address translation from EA to either an RA or PA.
acLookup(acTag)	<p>This is a function call used in command flows and transaction diagrams. It converts an acTag found in a TLX command packet into the address context (ac) used by the host's platform architecture to authenticate and provide the function requested by the TLX command.</p> <p>The result of an acLookup provides the error state of the address context provided. The state is shown as addressContext.state in the flows. The states are:</p> <ol style="list-style-type: none"> Good. The address context provided is valid. Invalid acTag. The acTag entry in the acTag table is not valid, or the acTag is specified outside the acTag table range. See <i>Table 7-1</i> on page 104. Invalid address context. The BDF and PASID associated with the acTag are invalid. The address context returned by the look up is not valid. See <i>Table 7-1</i> on page 104. <p>The function description is host specific.</p>
ACK	Acknowledgment.
address context	<p>(ac or addressContext). Address context is the information associated with a particular BDF and PASID pair. The association is formed by actions specified by the host's platform architecture.</p> <p>For TLX commands, the acTag and the acTag table provide the BDF and PASID. See <i>Section 4 The acTag table</i> on page 93 for additional details.</p>
address context space	A PASID paired with a BDF uniquely identifies the address space associated with a request. In OpenCAPI, a request is a TLX command.
address tenure	In a split transaction bus protocol, the commands and addresses are sent on the bus by the master before any data that might be associated with the transaction is moved. After the address tenure is completed, the status of the completion is examined. The data, if any is specified, is sent conditionally based on the status.
AFU	Attached functional unit. Architecturally, AFU refers to an end point unit or resource. Communication from the processor to the AFU goes through a protocol stack, transaction layer (TL), data link layer (DL), and physical medium layer (PHY). Command and data packets at the AFU interface are specified by the AFU command/data interface, which is the interface between the AFU protocol stack and the AFU.
AFU protocol	<p>AFU protocol layer. This layer currently consists of:</p> <ul style="list-style-type: none"> AFU_C protocol layer AFU_M protocol
AFU_C	<p>A processing element that is able to generate and receive commands to obtain data in a checked-out (non-cached) state.</p> <p>It uses the AFU command/data interface to communicate with the AFU_C protocol stack. All addressing to the AFU_C protocol uses an EA only. It uses the AFU_C protocol stack to send and receive commands through the TLX.</p> <p>See <i>AFU type on page 25</i> for the different sub-types of an AFU_C.</p>
AFU_C protocol	AFU _C protocol layer. This protocol specifies the sequences on the AFU command/data interface and the OpenCAPI packet interface (TLX boundary) for an AFU _C -defined AFU.
AFU_M	<p>A processing element that receives commands to either provide or receive data. This element is a memory storage device and may be mapped to the system's memory address range.</p> <p>The attributes of the memory held by an AFU_M are managed by the operating system.</p> <p>It uses the AFU interface to communicate with the AFU_M protocol stack. All addressing to the AFU_M uses a PA only.</p> <p>See <i>AFU type on page 25</i> for the different subtypes of an AFU_M.</p>
AFU_M protocol	AFU _M protocol layer. This protocol specifies the sequences on the AFU command/data interface and the OpenCAPI packet interface (TLX boundary) for an AFU _M -defined AFU.

alias	<p>When address translation from one address type to another results in a many to one mapping, the set of addresses that map into the single address are referred to as alias of each other.</p> <p>During address translation, an alias is formed when two different addresses translate into the same address. For example:</p> <ul style="list-style-type: none"> • Two or more physical addresses (PA) of an OpenCAPI device map to the same host real address (RA). • Two or more host RA map to a single attached OpenCAPI device's physical address.
AMO	Atomic memory operation. This operation performs an atomic update to a naturally aligned memory location. In some cases, this type of operation returns the original value of the memory location.
AP	Attached processor. Synonymous with AFU.
ATC	Address translation cache. The architecture describes a model for both a host ATC and an AFU ATC. See <i>Section 1.6 Address translation</i> on page 28.
BAR	Base Address Register.
back-off event	An event that causes a retry of an operation at some future time. The architecture specifies one type of back off event: long. The back off duration is controlled by a configuration space register specified in the OpenCAPI platform architecture.
BE	Byte enable.
CAPI	Coherent Accelerator Processor Interface.
CAPP	Coherent accelerator processor proxy.
command packet	TL/TLX construct. Contains command information for TL-to-TLX and TLX-to-TL communication.
convert2PA(RA)	This is a function call used in command flows and transaction diagrams. This converts an RA seen on the host processor bus into a PA used by the attached OpenCAPI device. The mapping of an RA to a device PA is device and host platform dependent.
CRC	Cyclic redundancy check.
data carrier	Data is transported between the TL and TLX in data carriers, which are defined as 64-byte data flits.
DCP	<p>Data credit pool. Each command or response specified with immediate data consumes one or more data credits.</p> <p>To add a command or response to a DL frame's control flit, both the VC credit and the DCP credit must be atomically obtained. That is, you must have both to proceed to insert the command or response into the DL frame.</p> <p>Adding a command or response specified with immediate data to a DL control flit defines the order the data is sent towards its destination.</p> <p>See <i>Section 5.1.3 Data transport, order, and alignment</i> on page 97 for full details.</p>
dError	Data error.
Device	The device refers to hardware and software attached via an OpenCAPI interface comprised of the PHYX, DLX, TLX Framer/Parser, TLX, AFU protocol stack, AFU protocol layer AFU interface and the AFU itself. See <i>Figure 1-1 OpenCAPI stack</i> on page 24.
DL	OpenCAPI data link layer found on the host processor.
dLength	Data length (dL).
DLX	OpenCAPI data link layer found on the external OpenCAPI device.
DMA	Direct memory access. A technique for using a special-purpose controller to generate the source and destination addresses for a memory or I/O transfer.
dP, dPart	Data part (dP).
EA	Effective address. This is the address as seen by a program. Some host architectures refer to this as a virtual address (VA). Mapping from an EA to an RA requires address translation services.
ECC	Error correction code. A code appended to a data block that can detect and correct bit errors within the block.

Approved

Flit	An acronym for FLOW control digITs. Typically used in networking to specify the smaller pieces that a larger network layer packet is broken into. See FLITs . In this architecture specification, a flit is associated with the specification of a DL frame and is defined as a 64-byte unit of data. Control and data flits are specified.
flit-cycle	The amount of time it takes 64-bytes to be either sent or received at the DL/TL or DLX/TLX interface.
host	The host refers to the host processor attached via an OpenCAPI link. It is comprised of the OpenCAPI PHY, DL, TL Framer/Parser, TL, the Host bus protocol stack interface and the hosts processors and other components that are implementation dependent on the host connected. See <i>Figure 1-1 OpenCAPI stack</i> on page 24.
host bus protocol layer	Specifies the sequences on both the host bus and at the host bus protocol layer and the OpenCAPI packet interface to: <ul style="list-style-type: none"> • Respond to snooped host bus commands from the OpenCAPI packet to the OpenCAPI transaction layer to initiate action at the target AFU. • Master commands on the host bus, per the specification found in the OpenCAPI packet, from the OpenCAPI transaction layer. Respond back to the source AFU at the conclusion of the host bus operation via an OpenCAPI packet to the TL layer.
immediate data	Data associated with a command or response. Immediate data is the data specified for a write operation (the command and the data travel in the same direction). A read response has immediate data (the response and the data travel in the same direction). A read command does not have immediate data; the data arrives with the response.
inbound	The direction from the attached OpenCAPI device towards the attached processor chip.
LRU	Least recently used. A policy for a caching algorithm that removes from the cache the item that has the longest elapsed time since its last access. An algorithm used to identify and make available the cache space that contains the data that was least recently used.
MEM	The memory-mapped owner of the line. The owner could be the memory controller or an the owner of a memory-mapped I/O space. Some coherency protocols refer to this as a point of coherency (POC).
minimum signed integer value	4-byte value: x'8000_0000' 8-byte value: x'8000_0000_0000_0000'
MMIO	Memory-mapped input/output. Refers to the mapping of the address space required by an I/O device for Load or Store operations into the system's address space.
mnemonic specification	<p>The general format of a mnemonic for either commands or responses is based on a base command/response type and "dotted" subtypes.</p> <p>The following subtypes are currently specified:</p> <ul style="list-style-type: none"> .be Byte enable field specified (dot-be). .d Data transfer (dot-d). Used only for intrap_req commands. .n Used for commands that require address translation (dot-n). If the address translation results in a miss in the ATC, the results of the address translation are used for the current operation, but are not loaded into the ATC. <p>An implementation may:</p> <ul style="list-style-type: none"> • Ignore this directive. • Store the results in a TLB. <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p style="text-align: center;">Engineering Note</p> <p>The dot-n form is expected to be used with a host implementation that has a multi-level ATC. This form of the command allows <i>warming up</i> the higher levels of the ATC hierarchy without installing into the more resource-precious level 1 ATC.</p> </div>
MRU	Most recently used. One of the results of an <i>LRU</i> algorithm. The cache entry that has the shortest amount of elapsed time since its last access.
NACK	Negative acknowledgment.

Approved

naturally aligned data block	<p>A data block containing L bytes is naturally aligned when the address specifying the location of the data block is an integer multiple of the length of the block.</p> <p>Where: $i = \{0, 1, 2, 3, \dots\}$ $L = \{64, 128, 256\}$</p> <p>A naturally aligned data block is located from byte address $i * L$ through $(i + 1) * L - 1$.</p> <p>A command's address specification may not be aligned as specified above. An unaligned address points to a naturally aligned data block of length L, with a starting address of $\text{adr}(63:(\log_2 L)) \parallel (\log_2 L)0$.</p>
nMMU	<p>An abstraction of a host implementation-dependent construct that performs page-table walks based on the underlying page-table architecture as specified by the host architecture and host's platform architecture.</p> <p>In the command flows and transaction diagrams, it returns:</p> <ul style="list-style-type: none"> • nMMU_response.status = 0 when an address translation is successful. It returns page_size and access permissions. • nMMU_response.status <> 0 when software must be invoked to complete the requested address translation.
null control flit	<p>A null control flit is defined as using template x'00'. The 6-slot packet contains a 1-slot null command, and the remaining five slots are undefined. A return credit response found in slots 0 and 1 may be used to return credits. See <i>Section 6 TL and TLX template specifications</i> on page 99 for the specification of template x'00'.</p>
outbound	The direction from the processor chip to the attached OpenCAPI device.
PA	<p>Physical address. This refers to the address space owned by an AFU_M device. The host converts the RA to the AFU_M device's physical address space using configuration settings in the host that are determined during initialization of the attached OpenCAPI device.</p> <p><i>A PA is not the result of address translation of an EA as might be the case in some host architectures. The host maps the device's PA into its own (RA) address space.</i></p>
packet	<p>TL/TLX unit of information. A command packet contains commands. A response packet contains response information. See the specification of command and response packets in <i>Section 2 TL and TLX command and response specifications</i> on page 30.</p> <p>Data is transferred in address-aligned:</p> <ul style="list-style-type: none"> • 64-byte data flits
PHY	<p>The PHY layer interfaces to the DL and the network.</p> <p>This is the bit stream level specifying the electrical and optical transmission medium as well as the network interconnect topology.</p> <p>The current specification for the network is a point-to-point connection.</p>
PHYX	<p>On the OpenCAPI device, the PHYX layer interfaces to the DLX and the network. This is the bit stream level that specifies the electrical and optical transmission medium as well as the network interconnect topology. The current specification for the network is a point-to-point connection.</p>
pL, pLength	Partial length.
POC	Point of coherency. See definition of <i>MEM</i> .
RA	Real address. A real address is the result of address translation of an EA. Some host architectures refer to this as a physical address; this specification reserves the term physical address for other purposes. See the definition of <i>PA</i> .
Reserved/R	Indicates that a field or bit specification is reserved. A reserved field is set to zero and shall not be examined by an implementation. See <i>Architecture compliance terminology</i> on page 11.
response packet	TL construct that contains response information to commands. Used for TL-to-TLX and TLX-to-TL communication.
responder	TL or TLX that accepts a command, services the command, and sends back a response TL/TLX packet that provides data, when required, and status of the service to the command.
requester	TL or TLX that issues a command. The requester collects all responses returned by the responder, if any, to determine the status of the service provided by the command. When the command is posted, responses are not returned.
RTL	Register transfer language.

Approved

segment	When used in reference to data, a segment refers to a naturally aligned 64-byte portion of a data transfer. For example, a 256-byte data transfer contains four segments.
service queue	<p>The members of a service queue are an ordered set of commands. The commands are selected by applying a hash against the VC, BDF, PASID and stream_id associated with the command. The hash results in the selection of a specific service queue. The hash is both implementation and command dependent. The hash is command dependent because not all commands are specified with a BDF, PASID and stream_id. Commands that are not specified with a VC do not enter a service queue.</p> <p>Per VC, the following operands may be included in the hash:</p> <p>TLX.vc.0 This VC is used for most responses, and the hash is the VC.</p> <p>TLX.vc.3 Contains various read and write TLX commands as well as assign_actag.</p> <p>The assign_actag command is serviced before entering into a service queue. All other commands are sorted using a hash based on the VC, BDF, PASID and stream_id.</p> <p>When the hash specified is used for a VC, the hash is perfect and the resulting service queue is identical to the definition of a <i>virtual queue</i>.</p> <p>When an implementation removes hash terms from the VC-specific specification, the hash is not perfect.</p> <p>There is at least one service queue per VC supported by the implementation.</p>
slot	A slot is a 28-bit granule used to specify a TL or TLX command or response packet.
SUE	Special uncorrectable error. Refers to error detection and attempted correction to a block of data. A SUE indicates that an error was detected upstream from the present error detection logic. The use of SUE indications aids in determining error origination as part of a first error incident reporting scheme.
TL	<p>OpenCAPI transaction layer found on the host processor.</p> <ul style="list-style-type: none"> • Interfaces to the DL and the protocol layer. Responsible for command-packet formation and response-packet handling and formation. Ensures that the order of data sent to the DL matches the command- and response-packet order sent to the DL. • Manages data flits. Associates the data with the command or response packet that was received prior to the arrival of the data. The command- and response-packets contain data descriptors that enable this association. • Performs flow control. • Performs error handling and control. • Manages all <i>service queues</i> associated with each virtual channel. Order is retained within virtual channels.
TLB	Translation lookaside buffer. An on-chip cache that holds the translation of an effective address (EA) to a real address (RA). A TLB caches page-table entries for the most recently accessed pages, thereby eliminating the necessity to access the page table from memory during load-store operations.
TLX	<p>OpenCAPI transaction layer found on the external OpenCAPI device.</p> <ul style="list-style-type: none"> • Interfaces to the DLX and the protocol layer. Responsible for command packet formation and response packet handling and formation. Ensures that the order of data sent to the DLX matches the command and response packet order sent to the DLX. • Manages data flits and associates the data with the command or response packet that was received prior to the arrival of the data. The command and response packets contain data descriptors that enable this association. • Flow control. • Error handling and control.
UE	Uncorrectable error. Refers to error detection and attempted correction to a block of data. An uncorrectable error indicates that an error was detected and the attempted correction failed.
VC	Virtual channel. See <i>Section 3 Virtual channel and data credit pool specification</i> on page 84, and the specification of all commands and responses in this section.
virtual queue	<p>The specification of a <i>service queue</i> describes the VC-specific hash required to form a service queue from a virtual queue.</p> <p>The members of a virtual queue are an ordered set of commands received from a VC. That is, the ordering of the commands found in the VC shall be retained when adding commands from the VC to a virtual queue.</p>
warming up	The process of loading or populating a cache with a set of valid data.
write class command	A command that is used to write data to a destination. The source of a write class command is also the source of the data.

Approved

<p>xlate_result = adr_xlate(EA, addressContext)</p>	<p>This is a function call used in command flows and transaction diagrams. This returns the host's results from an address translation. The function returns an RA (xlate_result.RA) and a status (xlate_result.status). The status returned is:</p> <ol style="list-style-type: none">1. Complete. Address translation completed successfully with an RA provided. The ATC may have been updated with the result.2. rty_req. Indicates that the address translation could not be completed at this time. The operation may be attempted at a later time.3. xlate_pending. Indicates that the address translation could not be completed. The ATC did not contain the translation and software was invoked. An asynchronous xlate_done TL command is sent when the software actions have completed. <p>Engineering note</p> <p>The Resp_code=xlate_pending is sent in a read_failed, touch_resp, or write_failed response packets. These TL responses shall precede the xlate_done command in the TL.vc.0 virtual channel.</p>
--	--

1. Overview

The OpenCAPI transaction layer specifies the control and response packets passed between a host and an OpenCAPI device. The transaction layer implemented on the host is referred to as the TL. The transaction layer implemented on the OpenCAPI device is referred to as the TLX.

On the host, the transaction layer converts:

- Host-specific protocol requests into transaction-layer-defined commands.
- TLX commands into host-specific protocol requests. When the host protocol completes, it provides responses to the TLX commands when required.
- TLX responses into responses for host-initiated requests.

On the OpenCAPI device, the transaction layer converts:

- AFU-specific protocol requests into transaction-layer-defined commands.
- TL commands into AFU-specific protocol requests. When the AFU protocol completes, it provides responses to the TL commands when required.
- TL responses into responses for AFU-initiated requests.

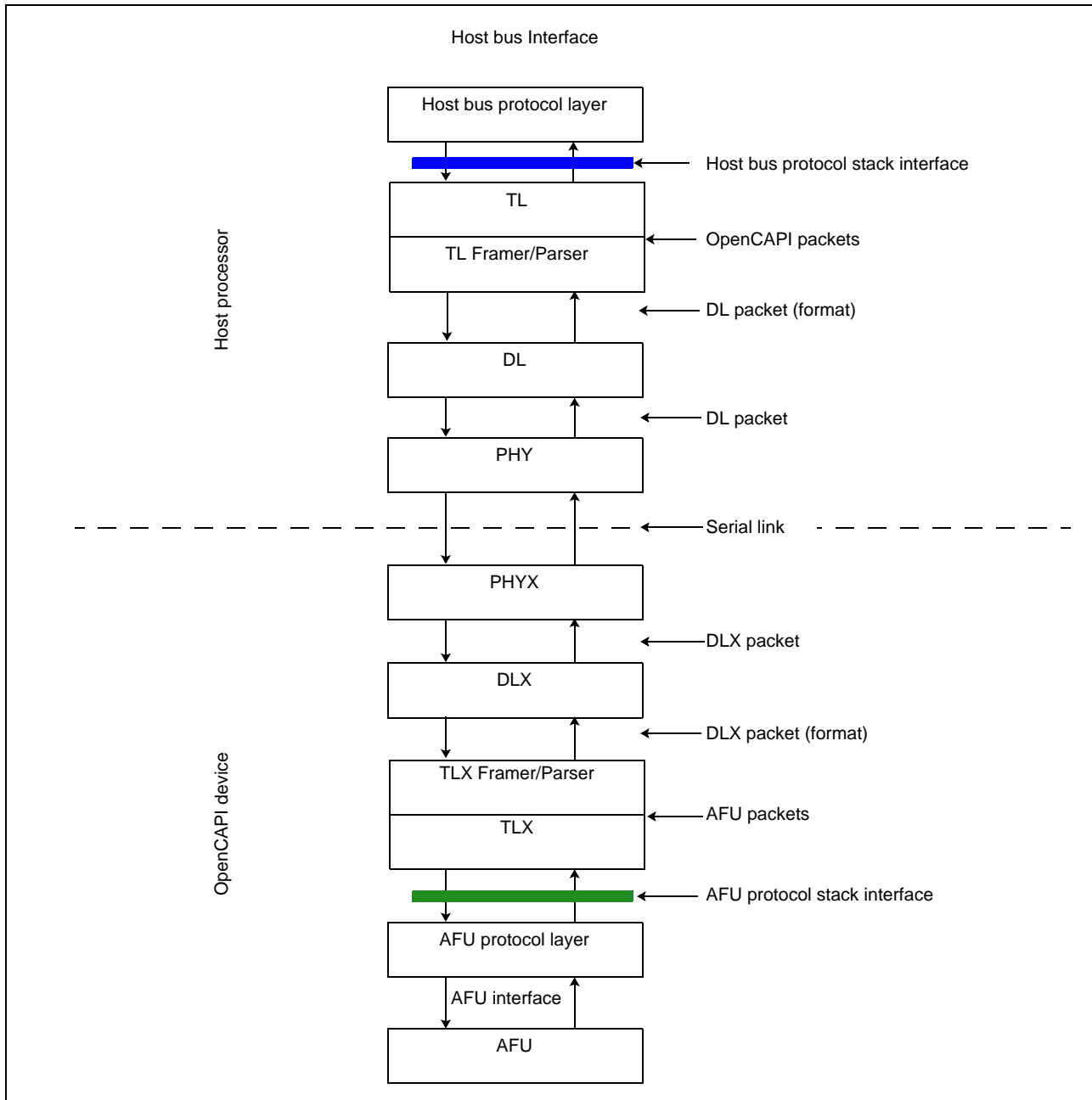
Working together, the TL and TLX provide a standard method to bridge between a host protocol architecture and an AFU protocol architecture. This is accomplished by the exchange of command and response packets specified by the OpenCAPI transaction layer specification.

Approved

1.1 OpenCAPI protocol stack

Figure 1-1 on page 24 shows the OpenCAPI protocol layers.

Figure 1-1. OpenCAPI stack



Approved

1.2 Host operation modes

The interface between the host and the AFU can be implemented with varying levels of complexity. Interoperability with an AFU implementation is dependent on the operation mode supported by the host and the requirements of the AFU.

The various combinations of AFU capabilities are broken into two subclasses, **AFU_C** and **AFU_M**, and each subclass is broken into two types.

AFU type	Description
AFU _{C0}	(C0 or none). There is no visible-to-the-host processing element. The host never sees any commands sourced by the TLX. While a processor element might not be visible to the host, it may still be present. If it is present, it shall not cache any lines in any coherent data valid state and shall not rely on the host's coherency protocol for correct operation.
AFU _{C1}	(C1 or type 1 processing element). A processing element with no cache. An AFU _{C1} may issue TLX commands to the host. It uses an EA to access host system memory. The host provides address translation and access to system memory.
AFU _{M0}	(M0 or none). There is no host system address space mapped to this device. That is, host system address space shall not be mapped to this device. Configuration space may be specified for this device.
AFU _{M1}	(M1 or type 1 MEM). A range of host system address space shall be assigned to this device. This address range shall be accessible only through the host (TL-to-TLX interactions). The host shall use the PA to access data. <div style="border: 1px solid black; padding: 5px; margin: 10px auto; width: fit-content;"> <p style="text-align: center;">Engineering Note</p> <p>The address range assigned to this type of device may be limited to MMIO space only or may include memory that is managed by the operating system; for example, memory that is backed by DRAM and that can be migrated to disk as needed.</p> </div>

The following sections describe the combinations of AFU_C and AFU_M devices on a single OpenCAPI device.

1.2.1 No attached device (C0, M0)

No device is attached to the OpenCAPI interface. No transactions occur.

1.2.2 MEM-only mode (C0, M1)

In this mode of operation, the AFU appears to be a memory controller with an address space mapped into the host's system address space. Access by the host uses the PA. See the specification of an AFU_M.

Developer note

While a processor element might not be visible (C0), it might still be present. Examples of this type of configuration that meet the above requirements are:

- An encrypted memory device. In this device, the data is encrypted/decrypted when the data is written to the M1 or read from the M1. The processing element that performs the encryption/decryption is not visible to the host. Topologically, the processing element is between the memory and the host.
- A memory cache device. The cache is managed by a processing element. The cache exists to reduce latency, and the cache states are not related to the host's coherency protocol. Data might be fetched or stored into the cache. The host cannot tell this is happening except for an improvement in performance.

1.2.3 Checkout mode (C1, M0)

In this mode of operation, the AFU appears as a processing element without a cache. It may have a non-coherent scratch pad memory, which is used for local processing only. Access to system memory is permitted as coherent, no-intent-to-cache actions. The AFU shall use an EA for these requests. The host shall perform address translation to enable access to system memory.

1.2.4 Checkout with MEM (C1, M1)

In this mode of operation, the AFU appears as a processing element without a cache. It may have a non-coherent scratch pad memory, which is used for local processing only. Access to system memory is permitted as coherent, no-intent-to-cache actions. The AFU shall use an EA for these requests. The host shall perform address translation to enable access to system memory.

In addition, the AFU provides a memory controller function with an address space mapped into the host's system address space. Access by the host uses a PA. See the specification of an AFU_M .

1.3 Command ordering

Ordering within a VC is maintained through the TL/TLX, but it is not assured after the command has moved to the upper protocol layers (host and AFU) as described in *Section 3 Virtual channel and data credit pool specification* on page 84.

1.4 Write fragmentation ordering and atomicity

1.4.1 Write fragmentation ordering and atomicity at the host

Write commands issued by the AFU may be fragmented by the host. The following sections specify the atomicity of the fragments and the order in which the updates become globally visible.

1.4.1.1 Partial write operations

These are TLX commands found in the following command classifications: `pr_dma_write`, `atomics.r`, `atomics.rw`, and `atomics.w`.

Minimum guaranteed write atomicity is specified as 16 bytes when aligned on a 16-byte address boundary. When the partial write operation is not specified with a naturally aligned address, atomicity may be reduced to a single byte. Data shall be globally visible in increasing address order.

1.4.1.2 64-, 128-, 256-byte write operations

These are TLX commands restricted to 64-, 128-, or 256-byte naturally aligned write operations. These are TLX commands found in the following command classification: `dma_write`.

Minimum guaranteed write atomicity is specified as 64 bytes. When the write operation specifies 128 or 256 bytes and the host fragments the write operation, there is no ordering guarantee for the data segments written.

Developer note

Applications are expected to use 64-byte or smaller writes for synchronizing required events; for example, when writing a semaphore. Larger data transfer sizes are expected to be used for data transfer efficiency.

1.4.2 Write fragmentation ordering and atomicity at the AFU

Write commands issued by the host may be fragmented by the AFU. The following sections specify the atomicity of the fragments and the order in which the updates become globally visible.

1.4.2.1 Partial write operations

These are TL commands found in the following command classifications: `pr_mem_write` and `configuration`.

Minimum guaranteed write atomicity is specified as 16 bytes when aligned on a 16-byte address boundary. When the partial write operation is not specified with a naturally aligned address, atomicity may be reduced to a single byte. Data shall be globally visible in increasing address order.

Developer note

The architecture does not currently provide TL commands that are able to specify partial write operations where the data and address are not naturally aligned.

1.4.2.2 64-, 128-, 256-byte write operations

These are TL commands restricted to 64-, 128-, or 256-byte naturally aligned write operations. These are TL commands found in the following command classification: `mem_write`.

Minimum guaranteed write atomicity is specified as 64 bytes. When the write operation specifies 128 or 256 bytes and the AFU fragments the write operation, there is no ordering guarantee for the data segments written.

1.5 OpenCAPI device PA space specification

An OpenCAPI device may have the following three PA spaces specified:

1. Configuration space shall be specified for the device.
2. System memory space may be specified for the device.
3. MMIO space may be specified for the device.

The configuration space is accessed by using the **`config_read`** or **`config_write`** commands. The PA specified for this space is separate from the system memory space and the MMIO space. The host may:

- Provide a configuration address BAR to access this space using a direct access load/store model.
- Provide an MMIO register set to access this space using an indirect access method.

Approved

The system memory space is memory space owned by the OpenCAPI device that is mapped to the host's system memory. The PA for system memory space is defined to start at offset 0. The host differentiates between the different system memory spaces of different OpenCAPI devices by providing a configuration address BAR for each attached device.

The MMIO space shares the PA space used by the system memory space. It is specified by a fixed offset from PA 0 which is specified in the OpenCAPI device's configuration space. The host differentiates the MMIO spaces of different OpenCAPI devices by providing a configuration address BAR for each attached device. Access to MMIO space is sensitive to the operand length and the command specified. The device literature should provide information on how to correctly access MMIO space.

- A device may not support all operand lengths provided by the architecture when accessing a specific address found in MMIO space. If an MMIO access does not use a correct operand size for the address specified, an unsupported-operand-length Resp_code shall result.
- A device may not support all commands provided by the architecture when accessing a specific address found in MMIO space. If an MMIO access does not use a correct command for the address specified, a Failed Resp_code shall result.
- All accesses to MMIO space shall result in a single response from the device. That is, when a dLength of 128 or 256 bytes is permitted, the device shall respond with the same dLength used in the command.

System and MMIO spaces are expected to be contiguous based on the configured starting PA and size. Access to unimplemented addresses results in the following:

- Read access to an unimplemented PA shall return all 1s data.
- Write access to an unimplemented PA shall result in discarded data.

1.5.1 PA-to-RA mapping rules

Real addresses (RA) are mapped into the physical address (PA) space specified for a device. This eliminates any requirement placed on the OpenCAPI device to have knowledge of the host's real address space or how the OpenCAPI device's PA space is mapped into it. The following rules place restrictions on the OpenCAPI device's specification of its PA space.

1. No address aliasing for PA-to-RA translation. That is, a PA for any specific device attached to an OpenCAPI link (PA + unique interface) translates into a unique RA. The address translation is specified by address ranges configured by software.
2. No address aliasing for RA-to-PA translation. The host protocol is provided with a single POC for each RA.

1.6 Address translation

1.6.1 Effective to real address translation

Effective to real address translation is specified by the host's architecture. The TL architecture model assumes a host address translation cache (ATC) that holds valid effective (EA) to real address (RA) translations. The ATC contains, at a minimum, a valid indication, the page size, the page size aligned starting effective address (EA), the address context of the translation in the form of the BDF and PASID, page write permission (W) and the host's corresponding real address.

Approved

The architecture supports multiple page sizes. See the specification of `log2_page_size` on page 32 and page size capability recommendations found in *Table 8-9 Profile specifications supported page size* on page 114.

An implementation may choose to provide additional fields, or may replace some of the fields listed above with other host specific content. Since the architecture assumes that the contents of an ATC entry contain the fields specified by the TL architectural model, any implementation specific alterations shall be done in such a manner that the differences are not externally observable.

The architecture model does not require, but allows for, a multi-level ATC. Higher level ATC might have a smaller capacity and have faster access than a lower level ATC that have more capacity and longer access latency. The structure of a host's ATC is outside the scope of this architecture.

The TL architectural model assumes that TLX commands with an EA specified go through effective to real address translation before execution on the host protocol bus. See *Section 3.3 TL Virtual channel and service queues* on page 88 for additional details.

An AFU can warm up the host's ATC by using the TLX `xlate_touch` command. See the command description for additional details.

Approved

2. TL and TLX command and response specifications

This section specifies all command and response types originated in the TL and the TLX. Commands originating in the TL are referred to as CAPP command packets (CAPP_cmd). Commands originating in the TLX are referred to as AP command packets (AP_cmd). Responses originating in the TL are referred to as CAPP response packets (CAPP_response). Responses originating in the TLX are referred to as AP response packets (AP_response)

In the subsections of this chapter descriptions use the following format:

Command descriptive name	mnemonic	Assigned opcode
command classification	VC used, DCP used (immediate data)	28-bit slot count

Table 2-1 lists the command operands used in the TL and TLX command and response specifications. See *Terms* on page 17 for definitions of terms used in these specifications.

Table 2-1. TL and TLX command operands (Page 1 of 4)

Operand mnemonic	Field width	Description
acTag	12	Address context tag. The address context tag is managed by the AFU. The acTag is used as an index into a host table that contains the BDF and PASID associated with the acTag. The OpenCAPI device learns its Bus number during a config_write , T=0 operations. The function and device numbers are assigned by the attached OpenCAPI device's implementation and cannot be modified by any configuration actions. The OpenCAPI device shall be assigned at least one PASID, and may be assigned more than one PASID, by host software during the initialization and operation of the device. The BDF and PASID are used for address translation authorization and operation validation.
AFUtag	16	<p>Unique handle specifying the AFU and command instance. Provided by the AFU that is requesting command services of the TLX. A TL response to a single TLX command may be broken into multiple TL response packets. When this occurs, all responses associated with the TLX command shall return the same AFU Tag value.</p> <div style="border: 1px solid black; padding: 10px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>The TL shall not use the AFUtag for any purpose other than as data to complete the contents of a response packet, or when forming an xlate_done or intrap_rdy TL command packet. Any retirement rules specified by a device implementation for the AFUtag shall not be checked by the TL.</p> <p>AFU tag retirement recommendations:</p> <ul style="list-style-type: none"> For non-posted commands, the AFU should not reuse an AFUtag until all responses for the command have been received. </div>
BDF	16	Bus device function. This is the identifier of a TLX requester. See acTag on page 30 for additional details.
Byte enable	64	(BE) This field is found in commands with dot-be mnemonic specifications. Valid only for write class commands.

Approved

Table 2-1. TL and TLX command operands (Page 2 of 4)

Operand mnemonic	Field width	Description
CAPPTag	16	<p>Unique handle specifying the host CAPP and command instance. Provided by the CAPP that is requesting command services of the TL.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p>Engineering Note</p> <p>The TLX shall not use the CAPPTag for any purpose other than as data to complete the contents of a response packet. Any retirement rules specified by the host implementation for the CAPPTag shall not be checked by the TLX.</p> <p>CAPPTag retirement recommendations:</p> <ul style="list-style-type: none"> • For non-posted commands, the CAPP should not reuse a CAPPTag until all responses for the command have been received. </div>
cmd_flag	4	Specifies execution behavior for commands and responses specified with this field. The command or response specification includes the behavior specification for the cmd_flag when the field is specified.
cmd_opcode	8	Specifies the operation to be performed.
credit_return	48	Specifies the number of credits returned to the VC and DCP credit pools. The credits are returned in fixed subfield locations in a 2-slot (56-bit) TL or TLX response packet. See the specification for return_tl_x_credits and return_tl_credits for the format of the field. Each VC credit allows for a single command or response to be sent in the virtual channel. Each DCP credit allows the sending of one <i>data carrier</i> .
dLength	2	<p>Data length (dL). Indicates the number of data bytes associated with a command or response packet. This 2-bit field indicates a length of:</p> <p>00 Reserved.</p> <p>01 64 bytes. This field value shall be used in a response packet when the command is a partial read or write operation.</p> <p>10 128 bytes. Reserved when the command is a partial read or write operation.</p> <p>11 256 bytes. Reserved when the command is a partial read or write operation.</p> <p>When the dLength field in the response packet does not match the full amount of data requested by the command, the dPart field is used to indicate the offset within the <i>naturally aligned data block</i> specified by the command's address. For example, in the multiple responses to a single TLX read command, the AFUtag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command.</p> <p>For multiple responses to a single command, there is no order requirement placed by the architecture. That is, continuing with the above example, the TLX may see the values of dPart returned in any order.</p> <p>Support for 256 bytes is optional for the TLX and AFU. See <i>Table 8-10 Profile specifications supported dLength by TLX</i> on page 114.</p>

Approved

Table 2-1. TL and TLX command operands (Page 3 of 4)

Operand mnemonic	Field width	Description
dPart	2	<p>Data part (dP(1:0) or dPart(1:0)). Indicates the data content of the current response packet. Read requests can be 64, 128, or 256 bytes in length. This field indicates the starting offset from the naturally aligned data block specified by the address provided in the read command. The amount of data transferred due to this response packet is found in the dLength field.</p> <p>00 Offset at 0 bytes. This field value shall be used for response packets when the command is a partial read or write operation.</p> <p>01 Offset at 64 bytes. This field value <i>shall not</i> be used when the dLength specifies 128 or 256 bytes. Reserved when the command is a partial read or write.</p> <p>10 Offset at 128 bytes. This field value <i>shall not</i> be used when the dLength specifies 256 bytes. Reserved when the command is a partial read or write.</p> <p>11 Offset at 192 bytes. This field value <i>shall not</i> be used when the dLength specifies 128 or 256 bytes. Reserved when the command is a partial read or write.</p> <p>The presence of this field in a command allows for multiple responses to be returned for a command. For example, a 256-byte read command such as rd_wn1tc may result in four responses with the dPart field taking on all four states.</p> <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p style="text-align: center;">Developer note</p> <p>The constraints placed on dPart are to ensure that responses specify naturally aligned data blocks. This is intended to simplify the design and the verification state space.</p> </div>
E	1	<p>Operand endianness. Used for mem_atomics.* and atomics.* class commands to specify the endianness of the operands. For bitwise logical operations, the endianness of the operands does not change the result. The field is specified as follows:</p> <p>0 Operands are little endian.</p> <p>1 Operands are big endian.</p>
EA	52, 59, 64	<p>Effective address (also referred to as the VA or virtual address by some host architectures). Length specification is dependent on the command issued and is noted in the command specification.</p>
log ₂ _page_size	6	<p>Log₂ value of the page size determined by the host when executing xlate_touch. The value of the page size touched in bytes is specified as $2^{\text{bin2dec}(\text{log}_2_page_size)}$.</p> <p>See Table 8-9 Profile specifications supported page size on page 114. Values corresponding to a page size of 1-byte to 2K-bytes are reserved.</p>
Object_handle	64/68	<p>Used by message class commands.</p> <p>For TL commands, the object handle is specified by the OpenCAPI device manufacturer and is loaded into a table maintained by the device software. It is accessed by the host based on a method specified by the host's OpenCAPI platform architecture.</p> <p>For TLX commands, the object handle is specified by the host architecture and is loaded into a table held in the device's MMIO space. The OpenCAPI device manufacturer specifies the location of the MMIO space, and it is provided to the host through the device software.</p>
PA	59, 64	<p>Physical address. Translation from the host RA to the AFU_M PA is performed by the host and configured during device initialization.</p> <p>The AFU's configuration space provides the information about the topology of the physical address space held by the OpenCAPI device. Types of address spaces are:</p> <ul style="list-style-type: none"> • Address space shared with the system. This excludes MMIO space. • MMIO space. • Configuration address space. This address space may be directly memory mapped and use a simple load/store model, or it may be accessed using indirect address methods. The choice is host dependent and is transparent to the OpenCAPI device and TL protocol.
PASID	20	<p>This term identifies the user process associated with a request. In OpenCAPI, a request is a TLX command. See acTag on page 30 for additional details.</p>

Table 2-1. TL and TLX command operands (Page 4 of 4)

Operand mnemonic	Field width	Description
pLength	3	(pL) Partial length. Specifies the number of data bytes specified for a partial write command. The address specified shall be naturally aligned based on the pLength specified. 000 1 byte. Reserved when the command is an amo* . 001 2 bytes. Reserved when the command is an amo* . 010 4 bytes. Reserved when the command is amo_rw and the operation is specified as a Fetch and swap. That is the command flag is {x'8'..x'A'}. 011 8 bytes. Reserved when the command is amo_rw and the operation is specified as a Fetch and swap. That is the command flag is {x'8'..x'A'}. 100 16 bytes. Reserved when the command is an amo* . 101 32 bytes. Reserved when the command is an amo* . 110 Specifies 4-byte operands when the command is amo_rw and the operation is specified as a Fetch and swap. That is, the command flag is {x'8'..x'A'}. Otherwise, this field is reserved. 111 Specifies 8-byte operands when the command is amo_rw and the operation is specified as a Fetch and swap. That is, the command flag is {x'8'..x'A'}. Otherwise, this field is reserved.
Resp_code	4	Response code. On a failed transaction, this field is found in a response packet reporting the reason the transaction failed. See the response packet for encoding and specifications. "Done" is not typically an included encoding because the response packet used is different for a failed transaction. For example, in response to a rd_writc AP command, the read_failed (TL response) is sent when the read is not able to complete successfully. The read_response (TL response) is used to indicate a successful completion and that data is associated with the response. A response code of "done" is implied with the read_response .
stream_id	4	Stream identifier used by the AP. This is used as part of the virtual channel, virtual queue, service queue specification.
T	1	Configuration read or write command type. 0 Indicates a type 0 configuration read or write command. A config_write , T=0 shall be used by the AFU to learn its bus number. For config_read with TL=0, the bus number is unchecked. 1 The operation shall result in a mem_wr_fail or mem_rd_fail TLX response with a Resp_code = Failed.

2.1 Handling multiple responses to a single command

As noted in the description of some responses, a single command may receive multiple responses. This might be due to a mismatch between the host's and OpenCAPI device's maximum data length specification.

For example, the host's or the device's internal bus protocol might be limited to atomically accessing 64 bytes of data. Read and write cases are examined in the following sections.

2.1.1 TLX Read request getting multiple TL responses

A 128-byte read request by the OpenCAPI device may be broken into two 64-byte read requests on the host protocol bus. This results in two TL responses returning data to the OpenCAPI device. The responses are not returned in any specified order. After all the responses are returned to the requester (the OpenCAPI device in this example), the requester examines the responses.

- When all responses indicate success, the command has completed successfully. In this example, each response provides 64 bytes of data, fulfilling the OpenCAPI device's request for 128 bytes.
- When all responses indicate failure, the command has failed. In this example, no data has been returned.

Approved

- When one response indicates success and the other indicates a failure, the command has failed. In this example, only 64 bytes of the requested 128 bytes have been returned. Because this is a read request, the data may be discarded. Depending on the Resp_code and the TL response, the entire operation, or just the failing portion may be retried. Refer to the specification of the TL response for when the operation may be retried.

The following TLX read commands may receive multiple responses.

- **rd_wnitc, rd_wnitc.n**

These read commands, when receiving multiple TL responses, shall see only the following responses²:

- **read_response, read_failed**

2.1.2 TLX Write request getting multiple TL responses

A 128-byte write request by an OpenCAPI device may be broken into two 64-byte operations on the host protocol bus. This results in two TL responses indicating the status of the write operation in the host. The responses are not returned in any specified order. After all the responses are returned to the device, the device examines the responses.

- When all responses indicate success, the command has completed successfully. The write operation has completed and the changes to the specified memory locations are globally visible.
- When all responses indicate failure, the command has failed. The locations in memory specified by the command may have been modified by the failed operation. That is, the data at the locations may be unmodified, may contain undefined data, or may contain SUE data. The Resp_code field in the fail response indicates what might have occurred at the memory location specified by the write command.
- When one response indicates success and the other indicates failure, the command has failed. Only the data corresponding to the 64-byte block specified by the successful response has completed its operation in the host and the changes to the specified memory location are globally visible. The data corresponding to the 64-byte block specified by the failed response shall contain SUE data. Depending on the Resp_code and the TL response, the failing portion may be retried and the successful portion shall not be retried. Refer to the specification of the TL response for when the operation may be retried.

The following TLX write commands may receive multiple responses:

- **dma_w, dma_w.n**

These write commands, when receiving multiple TL responses, shall see only the following responses.

- **write_response, write_failed**

2.1.3 TL read request getting multiple TLX responses.

A 128-byte read request by the host to an OpenCAPI device may be broken into two 64-byte read requests at the OpenCAPI device. This results in two TLX responses returning data to the host. Further, the responses are not returned in any specified order. After all responses are returned to the host, the host examines the responses.

- When all responses indicate success, the command has completed successfully. In this example, each response provides 64-bytes of data, fulfilling the host's request for 128 bytes.

2. Not all responses apply to all commands. See the command descriptions for applicable responses.

Approved

- When all responses indicate failure, the command has failed. No data has been returned.
- When one response indicates success and the other indicates failure, the command has failed. In this example, only 64-bytes of the requested 128 bytes have been returned. The data obtained may be discarded. Depending on the Resp_code and the TLX response, the entire operation or just the failing portion may be retried.

The following TL read commands may receive multiple responses:

- **rd_mem**

These read commands, when receiving multiple TLX responses, shall see only the following responses:

- **mem_rd_response, mem_rd_fail**

2.1.4 TL write request getting multiple TLX responses

A 128-byte write request by the host to an OpenCAPI device may be broken into two 64-byte write requests at the OpenCAPI device. This results in two TLX responses indicating the completion status of the write operation in the OpenCAPI device. Further, the responses are not returned in any specified order. After all the responses are returned to the host, the host examines the responses.

- When all responses indicate success, the command has completed successfully. The write operation has completed and the changes specified by the memory locations are globally visible,
- When all responses indicate failure, the command has failed. The locations in memory specified by the command may have been modified by the failed operation. That is, the data at the locations may be unmodified, may contain undefined data, or may contain SUE data. The Resp_code field in the fail response indicates what might have occurred at the memory location specified by the write command.
- When one response indicates success and the other indicates failure, the command has failed. Only the data corresponding to the 64-byte block specified by the successful response has completed its operation in the OpenCAPI device and the changes to the specified memory location are globally visible. The data corresponding to the 64-byte block specified by the failed response may be unmodified, may contain undefined data, or may contain SUE data. The contents of the data block is dependent on the address of the command. Depending on the Resp_code and the TLX response, the entire operation or just the failing portion may be retried. Refer to the specification of the TL response for when the operation may be retried and the state of the data block when the response indicates a failure.

The following TL write commands may receive multiple responses:

- **write_mem**

These commands, when receiving multiple TLX responses, shall see only the following responses:

- **mem_wr_response, mem_wr_fail.**

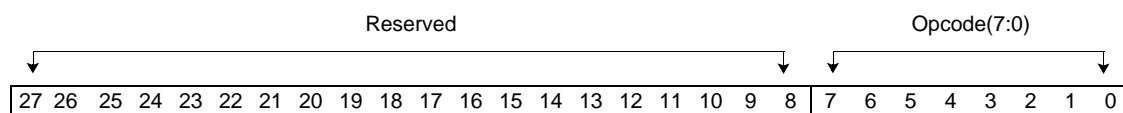
Approved

2.2 TL CAPP command packets

TL commands are sent from the host to the AFU. An alphabetical list of the TL commands follows; each command is hyperlinked to its specification. In this section, the TL command specifications are in opcode order.

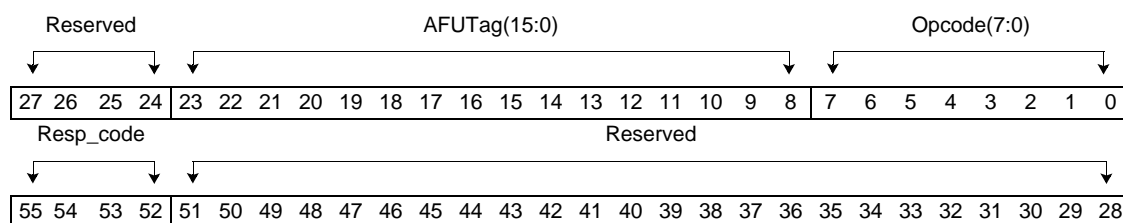
config_read **config_write** **intrp_rdy**
nop **pr_rd_mem** **pr_wr_mem**
rd_mem **write_mem** **write_mem.be**
xlate_done

No operation	nop	'0000 0000'
NA	NA	1



This command has no operands and performs no action. It is discarded at the TLX.

Address translation completed	xlate_done	'0001 1000'
async notification	TL.vc.0	2



The host is sending an asynchronous notification that an address translation requested by a prior TLX command has completed with the indicated response code. The remaining fields of the command identify the prior TLX command.

The TL is required to maintain the order of the matching **read_failed**, **write_failed**, or **touch_resp** response packets carrying the AFUtag and the Resp_code = intrp_pending when loading the VC. That is, the TL shall ensure that the response packets precede the **xlate_done** command in the VC.

The following illustrates how **xlate_done** is used:

1. The device issues a command that requires an address translation that the host is unable to complete.
2. The host responds with
 - a. a **read_failed** response with a Resp_code = xlate_pending. See *Table 2-11 read_failed resp_code use by tlx command* on page 70 for a list of the commands the device might have issued in step 1.

Approved

- b. a **write_failed** response with a `Resp_code = xlate_pending`. See *Table 2-13 write_failed resp_code use by tlx command* on page 74 for a list of the commands the device might have issued in step 1.
 - c. a **touch_resp** response with a `Resp_code = xlate_pending`. See *Table 2-9 touch_resp resp_code use by tlx command* on page 68 for a list of the commands the device might have issued in step 1.
3. Once the host has completed the address translation, the host issues **xlate_done** indicating if the device should retry or abort the operation.

The `Resp_code` field is specified in *Table 2-2*.

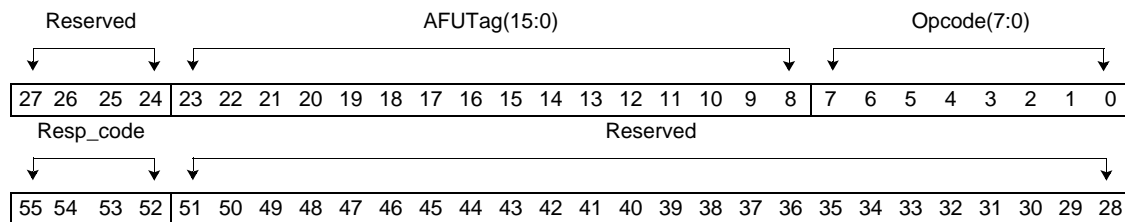
Table 2-2. The Resp_code specification for xlate_done

Resp_code encode	Description
'0000'	Completed. Address translation completed successfully
'0001'	Reserved.
'0010'	Retry request (<code>rtly_req</code>). Indicates that the address translation could not be completed at this time. The AFU may make an address translation attempt at a later time. This is a long back-off event.
'0011' - '1110'	Reserved.
'1111'	Translation address error (<code>adr_error</code>). Indicates that the address translation requested resulted in an address translation error.

Note: The errors specified by `Resp_code` do not include the fatal error conditions described in *Table 7-1* on page 104.

This command is posted.

Interrupt ready	intrp_rdy	'0001 1010'
async notification	TL.vc.0	2



The host is sending an asynchronous status notification for a previously attempted interrupt. The AFU determines its actions based on the `Resp_code` received. The `AFUTag` field of the command identifies the prior TLX command.

The TL is required to maintain the order of the matching **intrp_resp** carrying the `AFUTag` and the `Resp_code = intrp_pending` when loading the VC. That is, the TL shall ensure that the response packets precede the **intrp_rdy** command in the VC.

This command is used by

1. **intrp_req:**

- a. The device issues an **intrp_req** using the `cmd_flag` and `Object_handle` specified in the device's MMIO space.

Approved

- b. The host protocol, using the *Object_handle* for some form of address translation, is unable to complete. The host returns an **intrp_resp** with a *Resp_code* of *intrp_pending* (4).
- c. Once the host has completed the address translation, the host issues **intrp_rdy** indicating if the device should retry or abort the **intrp_req** operation.

The *Resp_code* field is specified in *Table 2-3*.

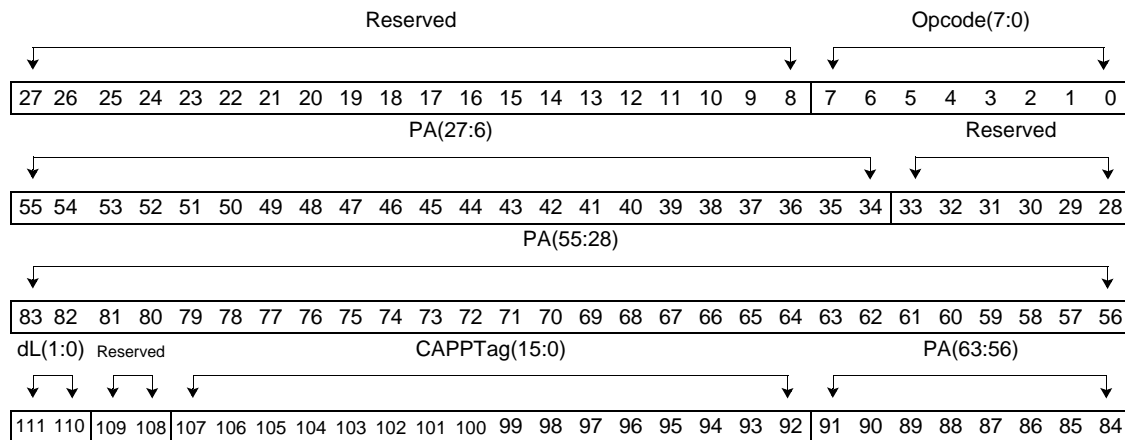
Table 2-3. The Resp_code specification for intrp_rdy

Resp_code encode	Description
'0000'	Ready to service the interrupt. The AFU may retry the prior intrp_req , or intrp_req.d command.
'0001'	Reserved.
'0010'	Retry request (<i>rtv_req</i>). Indicates that the host is unable to service the interrupt at this time. The AFU may retry the prior interrupt request at a later time as specified by its long back off event timer. This is a long back-off event.
'0011' - '1101'	Reserved.
'1110'	<p>Failed. The host is unable to service the interrupt specified by the prior command. Any future attempt specifying the same interrupt parameters shall fail.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p>Engineering Note</p> <p>Note that intrp_req is specified in a way that the <i>cmd_flag</i> and <i>Object_handle</i> are host specific and the values used are found in the device's configuration space. A device correctly using the <i>cmd_flag</i> and <i>Object_handle</i> should not normally see this <i>Resp_code</i>. A malicious device using values other than those provided in its MMIO space may see a failed <i>Resp_code</i>.</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for this <i>Resp_code</i>. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved.
Note: The errors specified by <i>Resp_code</i> do not include the fatal error conditions described in <i>Table 7-1</i> on page 104.	

This command is posted.

Approved

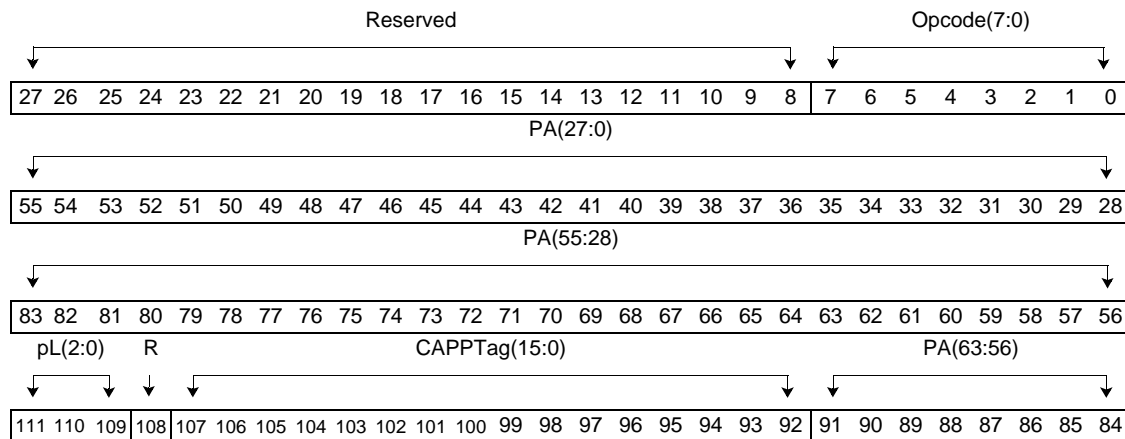
Read memory	rd_mem	'0010 0000'
mem_read	TL.vc.1	4



The host is requesting data to be read from the AFU memory. The data is a *naturally aligned data block* with a length specified by the dLength field (dL).

The response to this command is **mem_rd_response**, or **mem_rd_fail**. The **mem_rd_fail** response indicates the operation failed. See the response packet for encoding and specifications.

Partial memory read	pr_rd_mem	'0010 1000'
pr_mem_read	TL.vc.1	4

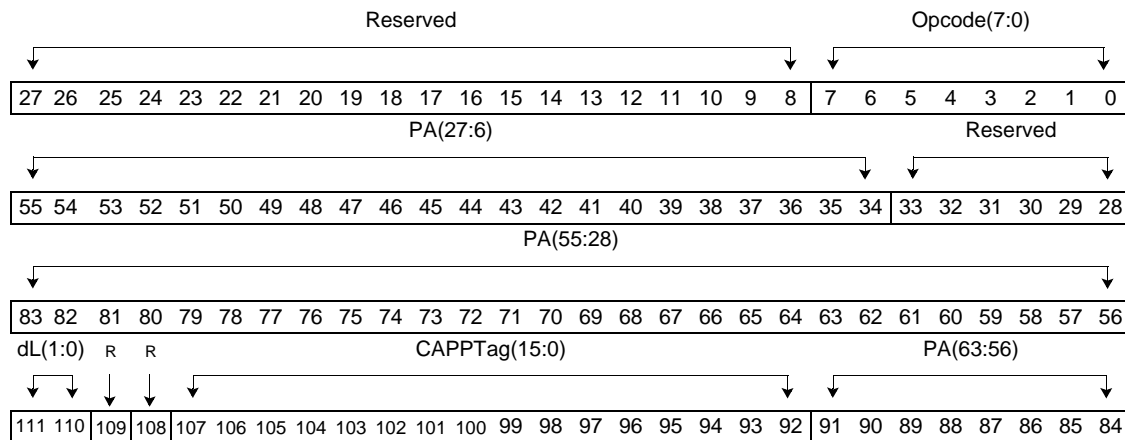


The host is requesting data to be read from the AFU memory. The number of bytes transferred is specified by pLength field(pL), and the starting address shall be naturally aligned based on the number of bytes requested. The pLength field limits the transfer size to 2^n bytes where $n = \{0..5\}$.

The response to this command is **mem_rd_response**, or **mem_rd_fail**. When a **mem_rd_response** is received, the data is found in the 64-byte data flit (only one is returned for this command). The data is address aligned as specified in *Section 5.1.3 Data transport, order, and alignment* on page 97. The **mem_rd_fail** response indicates the operation failed. See the response packet for encoding and specifications.

Approved

Write memory	write_mem	'1000 0001'
mem_write	TL.vc.1, TL.dcp.1	4



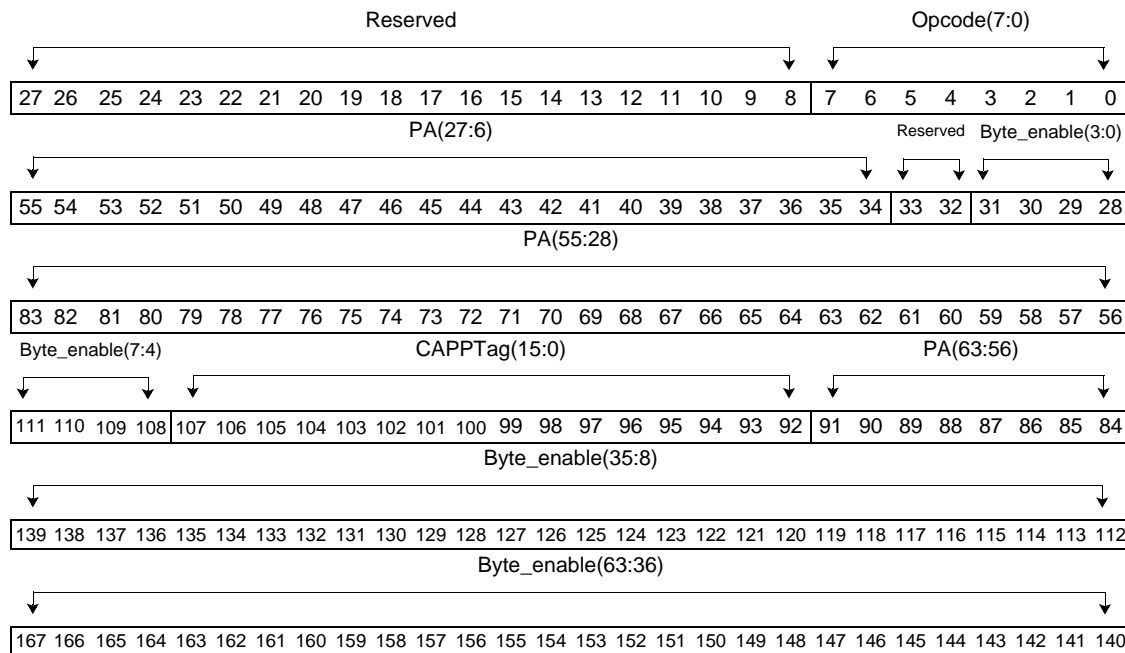
The host is writing a 64-, 128-, or 256-byte block of data to the AFU memory. The starting address is specified by the PA field and shall be naturally aligned based on the length of the data as specified by the dLength (dL) field.

This command is specified with immediate data. Credits for both the VC and DCP shall be obtained before this command is serviced by the TL.

The **mem_wr_response** and **mem_wr_fail** responses to this command indicate the status of the write to memory operation. The **mem_wr_response** indicates a successful completion of the operation. The **mem_wr_fail** response indicates that the operation failed. See the response packet description for encoding and specifications.

Approved

Byte enable memory write	write_mem.be	'1000 0010'
pr_mem_write	TL.vc.1, TL.dcp.1	6



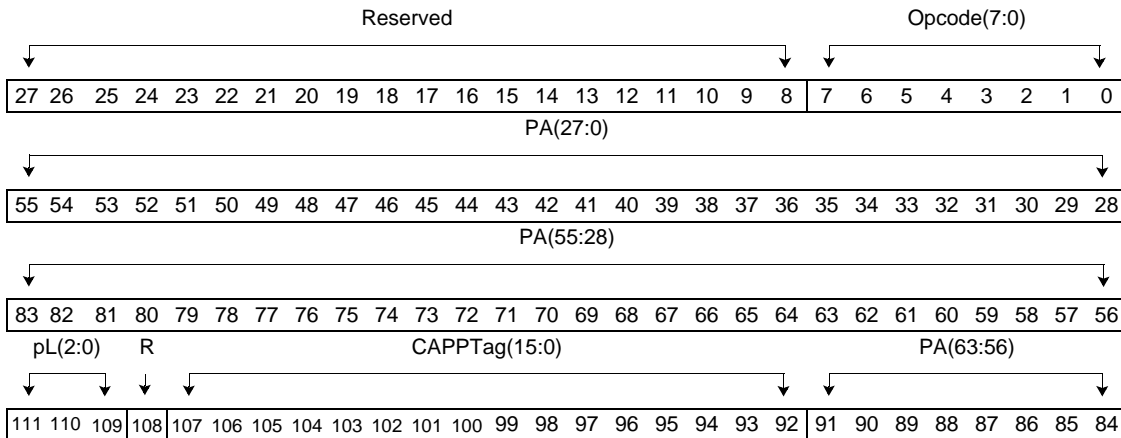
The host is writing a 64-byte data block to the AFU memory using a 64-bit byte-enable field. Each bit corresponds to one byte of the 64-byte aligned data block specified by the PA, where bit 0 of the BE determines if byte 0 of the data is written. When BE(n) is set to 1, byte n is written where n={0..63}.

This command is specified with immediate data. The data shall be sent using a single 64-byte data flit. Credits for both the VC and DCP shall be obtained before this command is serviced by the TL.

The **mem_wr_response** and **mem_wr_fail** responses to this command indicate the status of the write to memory operation. The **mem_wr_response** indicates a successful completion of the operation. The **mem_wr_fail** response indicates that the operation failed. See the response packet for encoding and specifications.

Approved

Partial cache line memory write	pr_wr_mem	'1000 0110'
pr_mem_write	TL.vc.1, TL.dcp.1	4



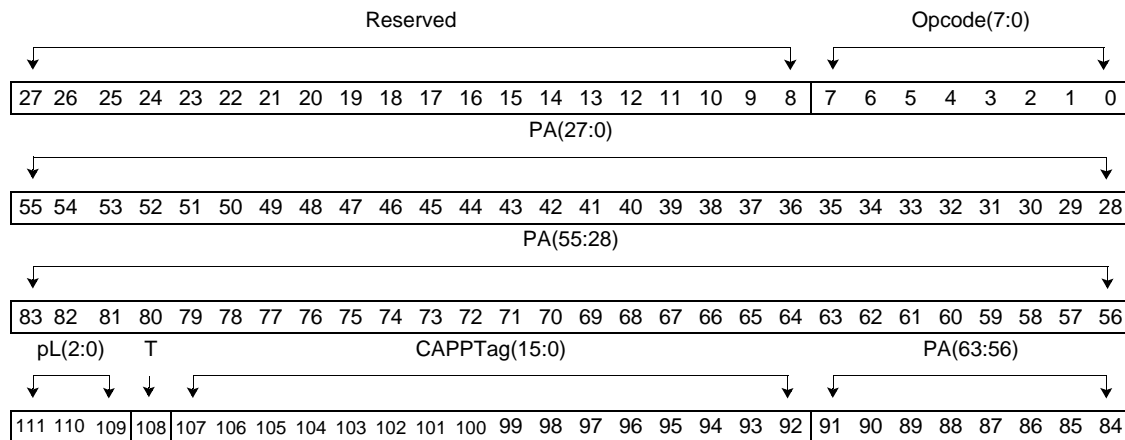
The host is writing the data to the AFU memory. The starting address is specified by the PA and shall be naturally aligned based on the length of the data as specified by the pLength (pL) field. The combination of the address and the pLength shall not cross a 64-byte address boundary.

This command is specified with immediate data. The data is address aligned as specified in *Section 5.1.3 Data transport, order, and alignment* on page 97. Credits for both the VC and DCP shall be obtained before this command is serviced by the TL.

The **mem_wr_response** and **mem_wr_fail** responses to this command indicate the status of the write to memory operation. The **mem_wr_response** indicates a successful completion of the operation. The **mem_wr_fail** response indicates that the operation failed. See the response packet for encoding and specifications.

Approved

Configuration read	config_read	'1110 0000'
configuration	TL.vc.1	4



The host is issuing a read to the AFU's configuration address space. The number of bytes transferred is specified by pLength (pL) field. The starting address shall be naturally aligned based on the number of bytes requested. For this command, the pLength value is limited to a transfer size of 1, 2, or 4 bytes. That is, the specification of pLength shall limit the transfer size to 2^n bytes where $n = \{0..2\}$.

The PA field is defined as follows:

PA bits	Description
63:32	Reserved. Shall be set to ³² 0.
31:24	Bus number (7:0).
23:19	Device number (4:0).
18:16	Function number (2:0).
15:12	Reserved. Shall be set to ⁴ 0.
11:2	Register number.
1:0	Byte offset within the register.

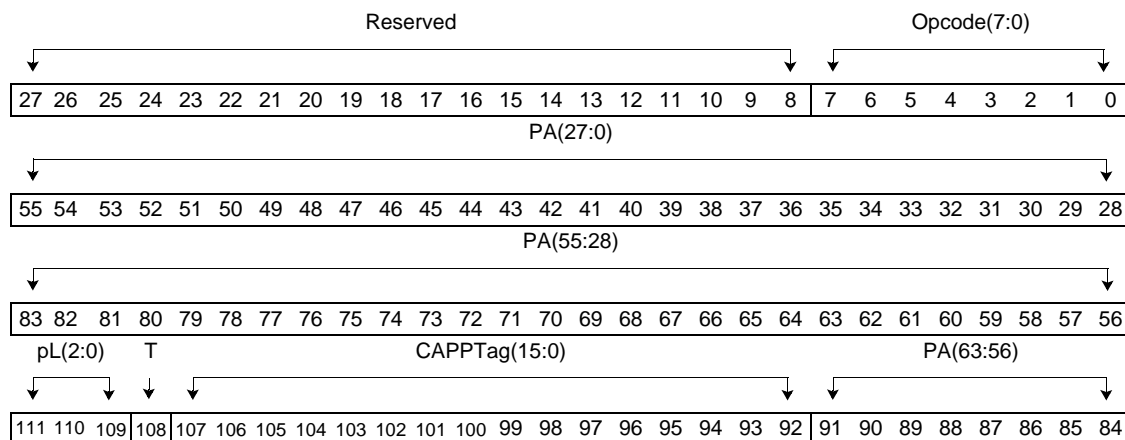
The T field, defined in *Table 2-1* on page 30, specifies the configuration type of the command.

The response to this command is **mem_rd_response**, or **mem_rd_fail**. When a **mem_rd_response** is received, the data is found in the 64-byte data flit (only one is returned for this command). The data is address aligned as specified in *Section 5.1.3 Data transport, order, and alignment* on page 97.

The **mem_rd_fail** response indicates the operation failed. If the device or function number are not recognized by the AFU, the operation shall fail with a Resp_code = Failed. See the response packet for encoding and specifications.

Approved

Configuration write	config_write	'1110 0001'
configuration	TL.vc.1, TL.dcp.1	4



The host is issuing a write to the AFU's configuration address space. The number of bytes transferred is specified by pLength (pL) field. The starting address shall be naturally aligned based on the number of bytes requested. For this command, pLength is limited to a transfer size of 1, 2, or 4 bytes. That is, the specification of pLength shall limit the transfer size to 2^n bytes where $n = \{0..2\}$.

The PA field is defined as follows:

PA bits	Description
63:32	Reserved. Shall be set to ³² 0.
31:24	Bus number (7:0).
23:19	Device number (4:0).
18:16	Function number (2:0).
15:12	Reserved. Shall be set to ⁴ 0.
11:2	Register number.
1:0	Byte offset within the register.

The T field, defined in *Table 2-1* on page 30, specifies the configuration type of the command. When T is set to 0, the AFU learns its bus number located in the PA field. The device and function number are assigned by the attached OpenCAPI device and are not modified by any configuration actions. If the device or function numbers are not recognized, the operation shall fail and the data is discarded. The failure shall be reported using a TLX **mem_wr_fail** response with a Resp_code= Failed.

This command is specified with immediate data. The data is address aligned as specified in *Section 5.1.3 Data transport, order, and alignment* on page 97. Credits for both the VC and DCP shall be obtained before this command is serviced by the TL.

Approved

The **mem_wr_response** and **mem_wr_fail** responses to this command indicate the status of the write to memory operation. The **mem_wr_response** indicates a successful completion of the operation. The **mem_wr_fail** response indicates that the operation failed. See the response packet for encoding and specifications.

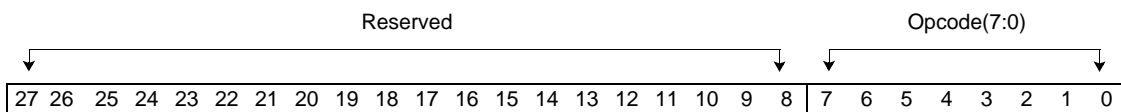
Approved

2.3 TLX AP command packets

TLX commands are sent from the AFU to the host. An alphabetical list of the TLX commands follows; each command is hyperlinked to its specification. In this section, the TLX command specifications are in opcode order.

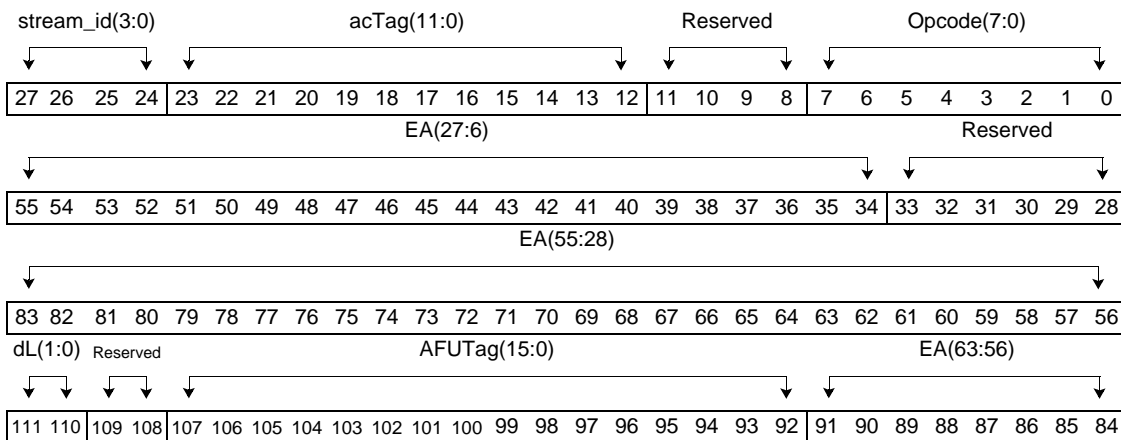
amo_rd	amo_rd.n	amo_rw	amo_rw.n
amo_w	amo_w.n	assign_actag	dma_pr_w
dma_pr_w.n	dma_w	dma_w.be	dma_w.n
dma_w.be.n	intrp_req	intrp_req.d	nop
pr_rd_wnitc	pr_rd_wnitc.n	rd_wnitc	rd_wnitc.n
wake_host_thread	xlate_touch	xlate_touch.n	

No operation	nop	'0000 0000'
NA	NA	1



This command has no operands and performs no action. It is discarded at the TL.

Read with no intent to cache	rd_wnitc rd_wnitc.n	'0001 0000' '0001 0100'
dma_read	TLX.vc.3	4



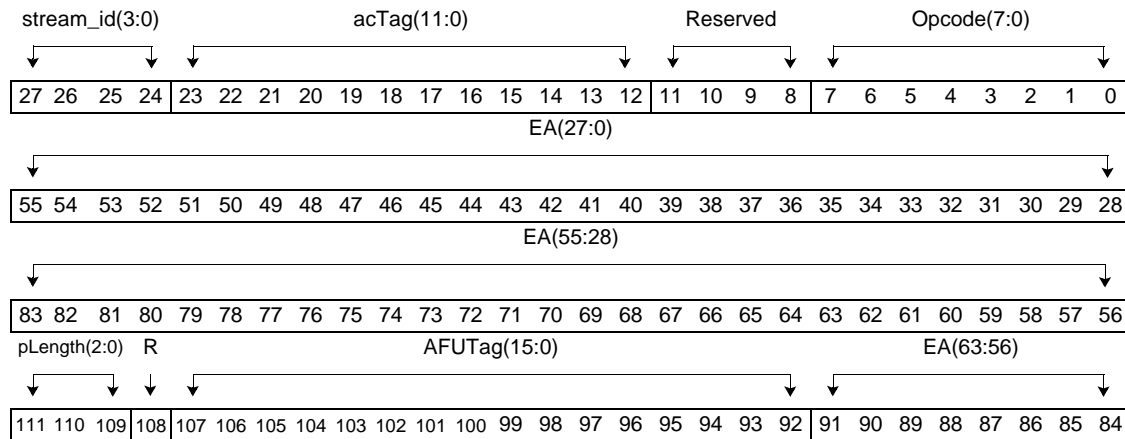
The AFU is requesting to read data with no intent to cache at the address specified by the EA. The data is a *naturally aligned data block* with a length specified by the dLength field (dL).

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

Approved

The response to this command is **read_response**, or **read_failed**. Multiple responses to a single **rd_wnitc** may occur. The **read_failed** response indicates the operation failed. See the response packet for additional details.

Partial read with no intent to cache	pr_rd_wnitc pr_rd_wnitc.n	'0001 0010' '0001 0110'
pr_dma_read	TLX.vc.3	4



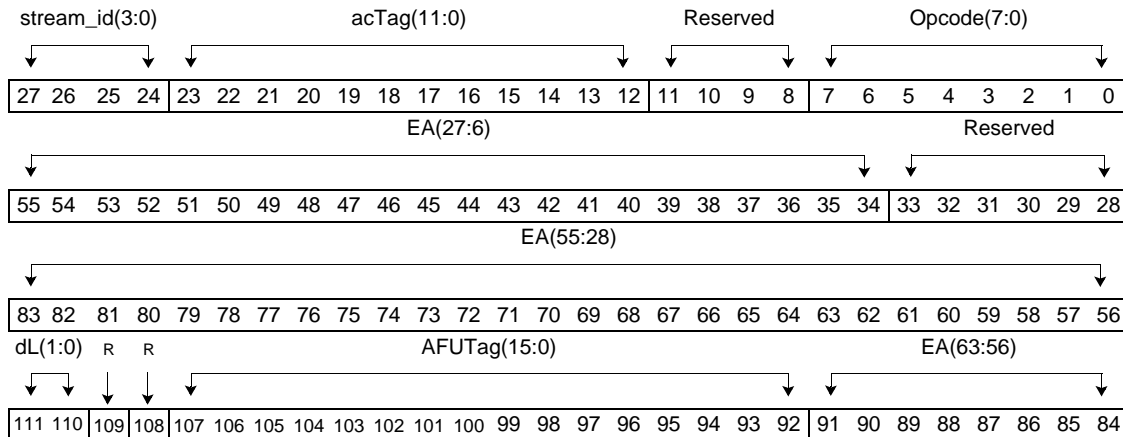
The AFU is requesting to read a partial cache line of data with no intent to cache at the address specified by the EA. The starting address shall be naturally aligned based on the length of the data specified by the pLength field. The pLength restricts this command to lengths of powers of 2 ranging from 1 to 32 bytes.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

The response to this command is **read_response**, or **read_failed**. When a **read_response** is received, the data is address aligned (address bits 5:0) in the 64-byte data flit (only one is returned for this command). The **read_failed** response indicates the operation failed. See the response packet for additional details.

Approved

DMA Write	dma_w dma_w.n	'0010 0000' '0010 0100'
dma_write	TLX.vc.3, TLX.dcp.3	4



The AFU is requesting to write data at the address specified by the EA. The starting address is specified by the EA and shall be naturally aligned based on the length of the data as specified by the dLength field.

This command is specified with immediate data. Credits for both the VC and DCP shall be obtained before this command is serviced by the TLX.

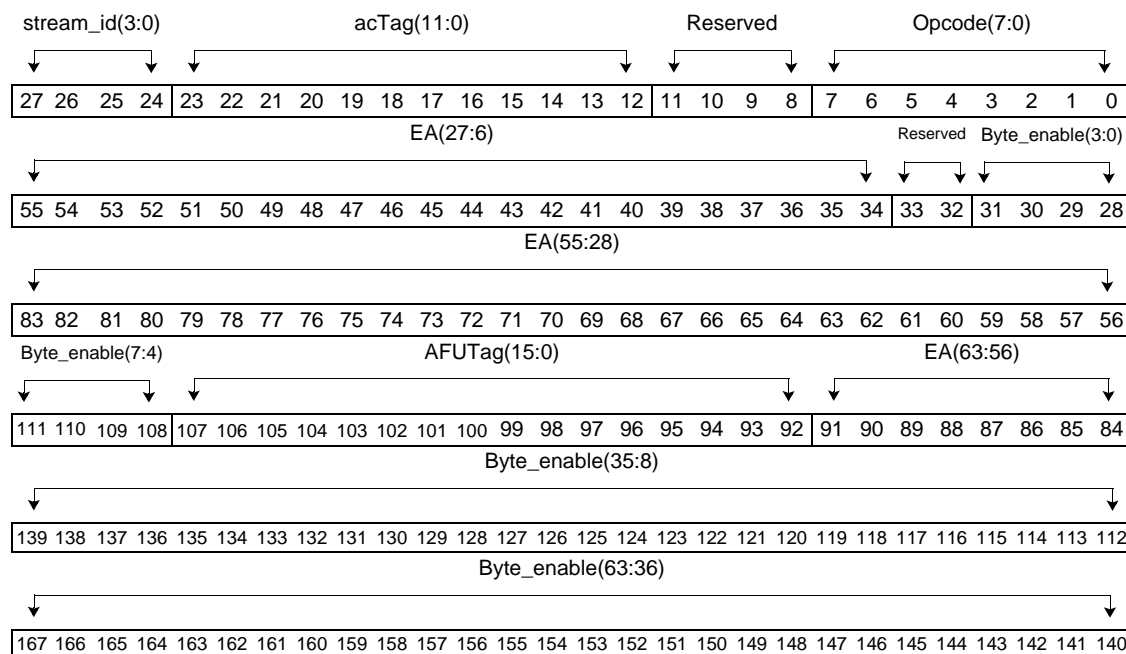
The AFU TLX shall not service this command unless all data specified by dLength is available to be sent.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

The host shall respond with either a **write_response** or a **write_failed** response packet. The **write_failed** response indicates that the operation failed. See the response packet description for additional details.

Approved

Byte enable DMA Write	dma_w.be dma_w.be.n	'0010 1000' '0010 1100'
pr_dma_write	TLX.vc.3, TLX.dcp.3	6



The AFU is writing data at the address specified by the EA using a 64-bit byte-enable field. Each bit corresponds to one byte of the 64-byte aligned data block specified by the EA, where bit 0 of the BE determines if byte 0 of the data is written. When BE(n) is set to 1, byte n is written where n={0..63}.

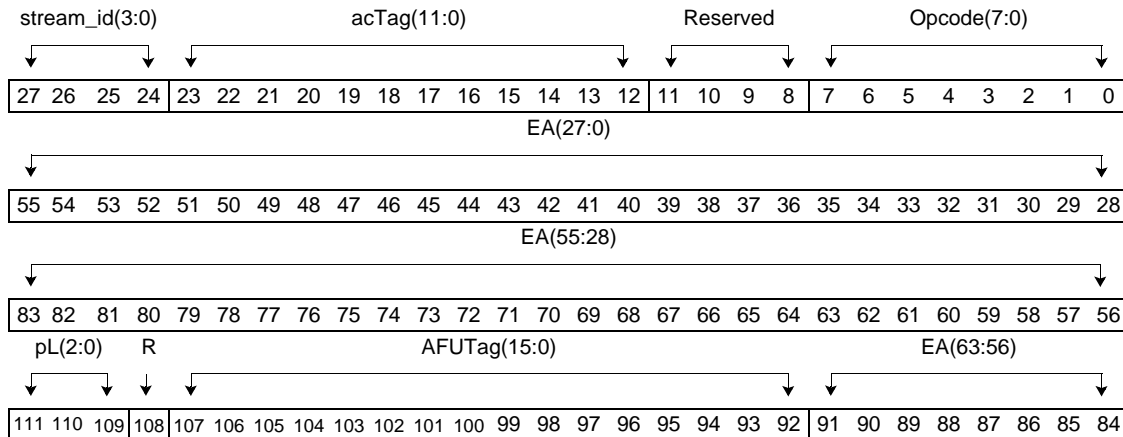
This command is specified with immediate data. The data shall be sent in a 64-byte data flit. Credits for both the VC and DCP shall be obtained before this command is serviced by the TLX.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

The host shall respond with either a **write_response** or a **write_failed** response packet. The **write_failed** response indicates that the operation failed. See the response packet description for additional details.

Approved

DMA parital write	dma_pr_w dma_pr_w.n	'0011 0000' '0011 0100'
pr_dma_write	TLX.vc.3, TLX.dcp.3	4



The AFU is requesting to write data starting at the address specified by the EA. The starting address shall be naturally aligned based on the length of the data as specified by the pLength (pL) field. The pLength restricts this command to lengths of powers of 2 ranging from 1 to 32 bytes. The combination of the EA and the pLength shall not cross a 64-byte address boundary.

Only a single data carrier is associated with this command.

This command is specified with immediate data. The data is address aligned as specified in *Section 5.1.3 Data transport, order, and alignment* on page 97. Credits for both the VC and DCP shall be obtained before this command is serviced by the TLX.

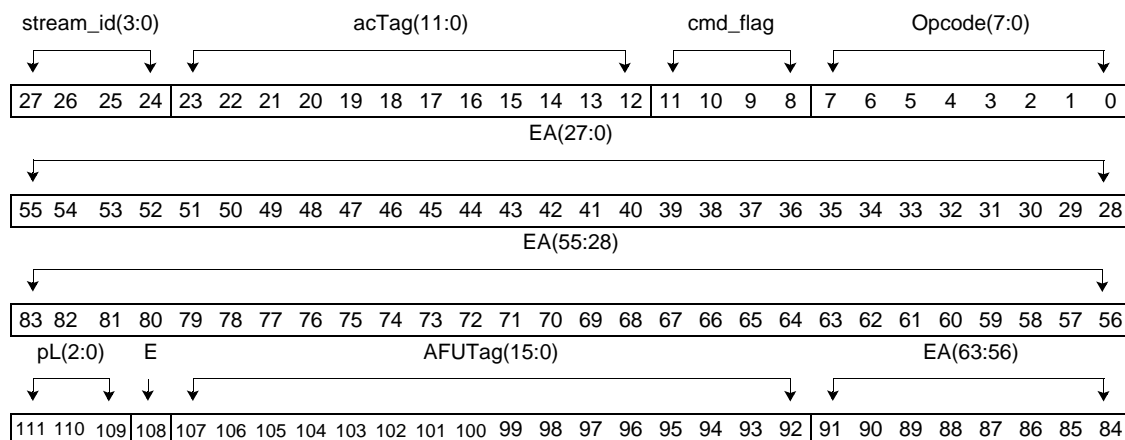
The AFU TLX shall not service this command unless all data specified by pLength is available to be sent.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

The host shall respond with either a **write_response** or a **write_failed** response packet. The **write_failed** response indicates that the operation failed. See the response packet description for additional details.

Approved

AMO read	amo_rd amo_rd.n	'0011 1000' '0011 1100'
atomics.r	TLX.vc.3	4



The AFU is requesting an atomic memory operation specified by the `cmd_flag`. All operands for this request are found in memory as specified by the EA.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

There shall be a single response to this command. The response to this command shall be one of the following TL responses: **read_response**, or **read_failed**. The **read_failed** response indicates that the operation failed. See the response packet specification for additional details.

Operation:

The operand length, as specified by the pLength (pL) field is restricted to 4- and 8-byte operands. That is, the pLength shall be specified as {'010', '011'}; all other values are reserved. Two signed integer operands are specified. The first operand "A" is found at the address specified by the command. The second operand "A2" is found at the address specified with an offset specified by the width of the operands and the operation; that is, as specified by pLength and by the command flag.

- For Fetch and increment bounded and Fetch and increment equal (that is, `cmd_flag = {'1100', '1101'}`), A2 is found at the address specified *plus* the width of the operand.
- For Fetch and decrement bounded (that is, `cmd_flag = {'1110'}`), A2 is found at the address specified *minus* the width of the operand.

The specification of the address is constrained to be naturally aligned. In addition:

- It cannot target locations at $32n - 2^{\text{bin}2^{\text{dec}(pL)}}$, where $n = 1, 2, 3, \dots$ (Fetch and increment bounded and Fetch and increment equal; that is, `cmd_flag = {'1100', '1101'}`).
- It cannot target locations at $32n$, where $n = 0, 1, 2, 3, \dots$ (Fetch and decrement bounded; that is, `cmd_flag = {'1110'}`).

The original value from the memory location specified by the command, or the 4- or 8-byte minimum signed integer value, is returned with **read_response**.

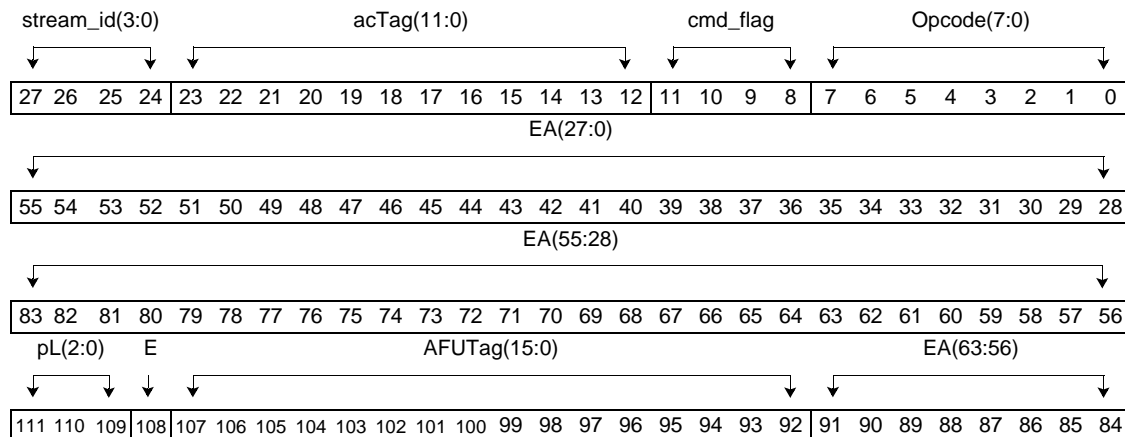
Approved

The operation performed is specified by the cmd_flag. The endianness of the operands is specified by the E bit.

Table 2-4. The cmd_flag specification for **amo_rd**

cmd_flag	Operation name and description
'0000' through '1010'	Reserved
'1100'	Fetch and increment bounded t ← A; If A != A2 then {A ← A+1; return t} else {return minimum signed integer value}
'1101'	Fetch and increment equal t ← A; If A = A2 then {A ← A+1; return t} else {return minimum signed integer value}
'1110'	Fetch and decrement bounded t ← A; If A != A2 then {A ← A-1; return t} else {return minimum signed integer value}
'1111'	Reserved

AMO read write	amo_rw amo_rw.n	'0100 0000' '0100 0100'
atomics.rw	TLX.vc.3, TLX.dcp.3	4



The AFU is requesting an atomic memory operation specified by the cmd_flag. For this request, operands are provided with the command and are found in memory as specified by the EA.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

This command is specified with immediate data. Credits for both the VC and DCP shall be obtained before this command is serviced by the TLX.

The AFU TLX shall not service this command unless all data specified by pLength is available to be sent.

Approved

There shall be a single response to this command. The response to this command shall be one of the following TL responses: **read_response** or **read_failed**. The **read_failed** response indicates that the operation failed. See the response packet specification for additional details.

Operation:

The command address specified shall be naturally aligned based on the operand length. The operand length is restricted to 4- and 8-byte operands. The pLength (pL) shall be specified as {'010', '011'} for all cmd_flag operations with the exception of fetch and swap operations where the cmd_flag is specified as {x'8'...x'A'} and pLength shall be specified as {'110', '111'}. Refer to the specification of *pLength on page 33*.

Operations specified by the cmd_flag use either two or three operands; additional classification of the operands can be found in the description of the operation in *Table 2-5*. The command's address specifies the location of a first operand, "A".

Operand A is operated on by the second operand, "V", which is provided as the command's write data. Operand V is aligned within:

- a 64-byte data flit based on address bits 5:0 specified by the command.

A third operand "W" is provided for fetch and swap operations. Operand W is placed in the same data carrier as operand V. Operand W shall be aligned within:

- the 64-byte data flit based on the following equation:

$$\text{alignment (5:0)} \leftarrow \text{EA(5:4)} \parallel (\text{EA(3:0)} + '1000')$$

Any carryout from bit 3 is ignored.

The original value from the memory location specified by the command shall be returned with a **read_response**.

The endianness of the operands is specified by the E bit. The value of E might not affect the result of the operation specified by the cmd_flag. This is noted in the operation description found in *Table 2-5*.

Table 2-5. The cmd_flag specification for amo_rw (Page 1 of 2)

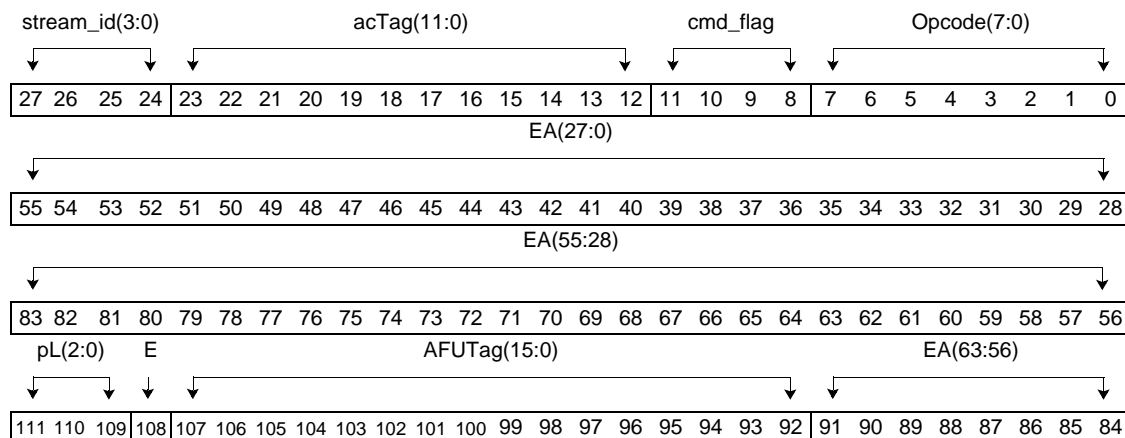
cmd_flag	Operation name and description
'0000'	Fetch and Add Operands are unsigned integers. $t \leftarrow A; A \leftarrow A + V; \text{return } t$ Overflow conditions are not reported.
'0001'	Fetch and XOR Operands are bit vectors. E has no effect on the operation. $t \leftarrow A; A \leftarrow V \oplus A; \text{return } t$
'0010'	Fetch and OR Operands are bit vectors. E has no effect on the operation. $t \leftarrow A; A \leftarrow V \vee A; \text{return } t$
'0011'	Fetch and AND Operands are bit vectors. E has no effect on the operation. $t \leftarrow A; A \leftarrow V \wedge A; \text{return } t$
'0100'	Fetch and maximum unsigned Operands are unsigned integers. $t \leftarrow A; A \leftarrow \text{Max}(A, V); \text{return } t$ A is unchanged when A is greater than or equal to V.

Approved

Table 2-5. The cmd_flag specification for *amo_rw* (Page 2 of 2)

cmd_flag	Operation name and description
'0101'	Fetch and maximum signed Operands are signed two's complement integers. $t \leftarrow A; A \leftarrow \text{Max}(A, V);$ return t A is unchanged when A is greater than or equal to V .
'0110'	Fetch and minimum unsigned Operands are unsigned integers. $t \leftarrow A; A \leftarrow \text{Min}(A, V);$ return t A is unchanged when A is less than or equal to V .
'0111'	Fetch and minimum signed Operands are signed two's complement integers. $t \leftarrow A; A \leftarrow \text{Min}(A, V);$ return t A is unchanged when A is less than or equal to V .
'1000'	Fetch and swap Operands are bit vectors. E has no effect on the operation. V is not used. $t \leftarrow A; A \leftarrow W;$ return t
'1001'	Fetch and swap equal Operands are bit vectors. E has no effect on the operation. $t \leftarrow A;$ When $V = A$, then $A \leftarrow W;$ return t
'1010'	Fetch and swap not equal Operands are bit vectors. E has no effect on the operation. $t \leftarrow A;$ when $V \neq A$, then $A \leftarrow W;$ return t
'1011' through '1111'	Reserved

AMO write	amo_w amo_w.n	'0100 1000' '0100 1100'
atomics.w	TLX.vc.3, TLX.dcp.3	4



The AFU is requesting an atomic memory operation specified by the cmd_flag. For this request, operands are provided with the command and are found in memory as specified by the EA.

This command is specified with immediate data. Credits for both the VC and DCP shall be obtained before this command is serviced by the TLX.

Approved

The AFU TLX shall not service this command unless all data specified by pLength (pL) is available to be sent.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

There shall be a single response to this command. The host shall respond with either a **write_response** or a **write_failed** response packet. The **write_failed** response indicates that the operation failed. See the response packet description for additional details.

Operation:

The command's address shall be naturally aligned based on the operand length. The operand length, as specified by the pLength (pL) field, is restricted to 4- and 8-byte operands. That is, the pLength shall be {'010', '011'}. All other values of pLength are reserved.

The number of operands specified by this command is determined by an examination of the cmd_flag. Two or three operands may be specified as shown in *Table 2-6* on page 55. The operands are designated as "A," "A2," and "V". The command's address specifies the location of each operand as follows:

- Operand A is located in memory at the address specified by the EA and shall be naturally aligned. When the cmd_flag indicates the use of operand A2, the address of operand A is further constrained and shall not target locations at $32n \cdot 2^{\text{bin2dec}(\text{pL})}$, where $n = 1, 2, 3 \dots$
- Operand A2 is located in memory at the address specified by the EA plus an offset specified by the width of the operands.
- Operand V is provided as the command's write data. Operand V shall be aligned within:
 - a 64-byte data flit based on address bits 5:0 specified by the command

The endianness of the operands is specified by the E bit. The value of E might not affect the result of the operation specified by the cmd_flag. This is noted in the operation description found in *Table 2-6*.

*Table 2-6. The cmd_flag specification for **amo_w** (Page 1 of 2)*

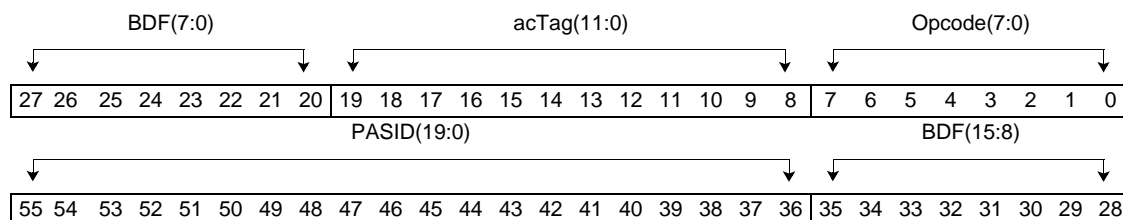
cmd_flag	Operation name and description
'0000'	Store and Add Operands are unsigned integers. $A \leftarrow A + V$ Overflow conditions are not reported.
'0001'	Store and XOR Operands are bit vectors. E has no effect on the operation. $A \leftarrow V \oplus A$
'0010'	Store and OR Operands are bit vectors. E has no effect on the operation. $A \leftarrow V \vee A$
'0011'	Store and AND Operands are bit vectors. E has no effect on the operation. $A \leftarrow V \wedge A$
'0100'	Store and maximum unsigned. Operands are unsigned integers. $A \leftarrow \text{Max}(A, V)$ A is unmodified when A is greater than or equal to V.

Approved

Table 2-6. The *cmd_flag* specification for *amo_w* (Page 2 of 2)

cmd_flag	Operation name and description
'0101'	Store and maximum signed Operands are signed two's complement integers. $A \leftarrow \text{Max}(A, V)$ A is unmodified when A is greater than or equal to V.
'0110'	Store and minimum unsigned Operands are unsigned integers. $A \leftarrow \text{Min}(A, V)$ A is unmodified when A is less than or equal to V.
'0111'	Store and minimum signed Operands are signed two's complement integers. $A \leftarrow \text{Min}(A, V)$ A is unmodified when A is less than or equal to V.
'1000 through '1011'	Reserved.
'1100'	Store and compare twin Operands are bit vectors. E has no effect on the operation. When $A = A2$, then $(A \leftarrow V, A2 \leftarrow V)$
'1101' through '1111'	Reserved.

acTag Assignment	assign_actag	'0101 0000'
acTag mgmt	TLX.vc.3	2



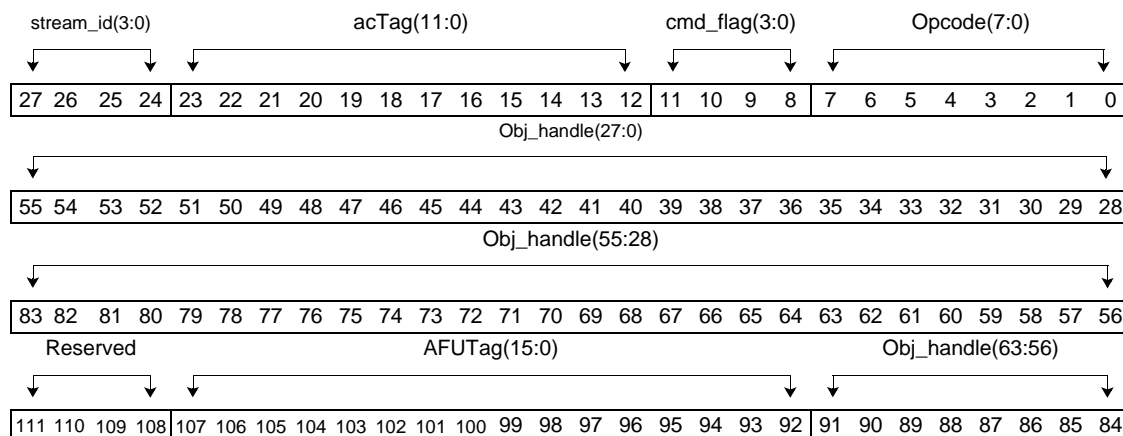
This command is used by the attached OpenCAPI device to assign an acTag value to a BDF and PASID combination. The OpenCAPI device uses this command to manage the contents of the acTag table. See *Section 4 The acTag table* on page 93 for the use of this command, the acTag table, and the management requirements placed on the OpenCAPI device.

This command is serviced when it reaches the head of the VC in the TL. It is not added to a *service queue*.

This command is posted. No response is sent for this command.

Approved

Interrupt Request	intrp_req	'0101 1000'
message	TLX.vc.3	4



This command is used to request interrupt service on the host. No data is transferred with this request. The specification of the object handle and the **cmd_flag** is found in the host's platform architecture.

The response to this command is **intrp_resp**.

Engineering Note

The AFUtag is passed back to the TLX in a response packet and has no control function in the TL as described in *Table 2-1 TL and TLX command operands* on page 30. The **intrp_resp**'s response code specification of **intrp_pending** indicates to the TLX that a subsequent **intrp_rdy** command is sent when the host is ready to service the interrupt. The **intrp_rdy** command contains the AFUtag that is specified in the original **intrp_req** command sent. While there are no requirements placed on the AFU to reserve the AFUtag used by the **intrp_req** command until the **intrp_rdy** command is received by the TLX, it is strongly recommended that an AFU implementation do so. It is an AFU implementation choice to use the AFUtag to precisely determine which interrupt to retry.

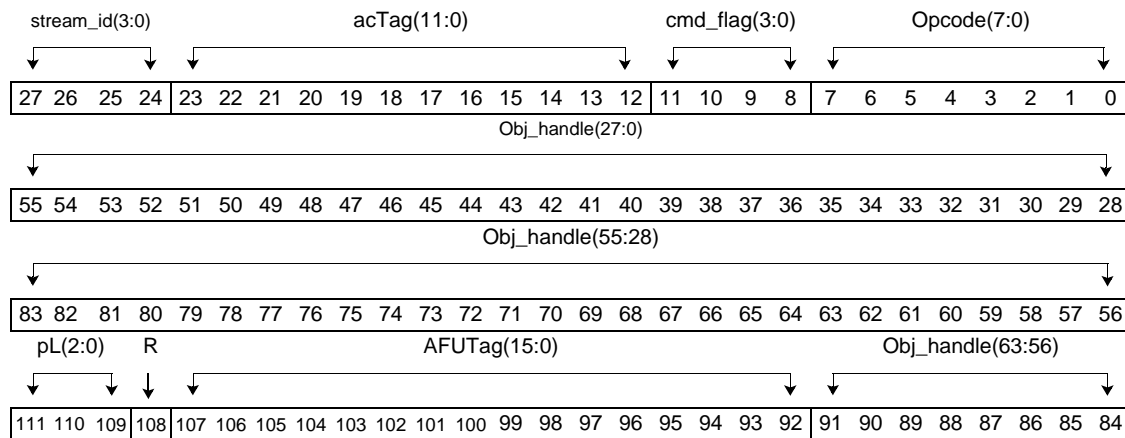
Developer Note

Both the command flag and the object handle fields specified for this command are specified by the host's platform architecture.

The attached OpenCAPI device provides MMIO space where the combination of the object handle and the command flag associated with this command are located. The manufacturer of the OpenCAPI device determines the number of command-flag and object-handle combination entries supported based on what is supported by the host and the function provided by the device.

Approved

Interrupt Request	intrp_req.d	'0101 1010'
message	TLX.vc.3, TLX.dcp.3	4



This command is used to request interrupt service on the host. Data, with the length specified by the pLength (pL) field, is transferred with this request. The data shall be sent in a 64-byte data flit. The alignment of the data within the data flit is specified by the host's platform architecture. The specification of the object handle and the command flag is found in the host's platform architecture.

The response to this command is **intrp_resp**.

Engineering Note

The AFUtag is passed back to the TLX in a response packet and has no control function in the TL as described in *Table 2-1* on page 30. The **intrp_resp**'s response code specification of **intrp_pending** indicates to the TLX that a subsequent **intrp_rdy** command is sent when the host is ready to service the interrupt. The **intrp_rdy** command contains the AFUtag that is specified in the original **intrp_req.d** command sent. While there are no requirements placed on the AFU to reserve the AFUtag used by the **intrp_req.d** command until the **intrp_rdy** command is received by the TLX, it is strongly recommended that an AFU implementation do so. It is an AFU implementation choice to use the AFUtag to precisely determine which interrupt to retry.

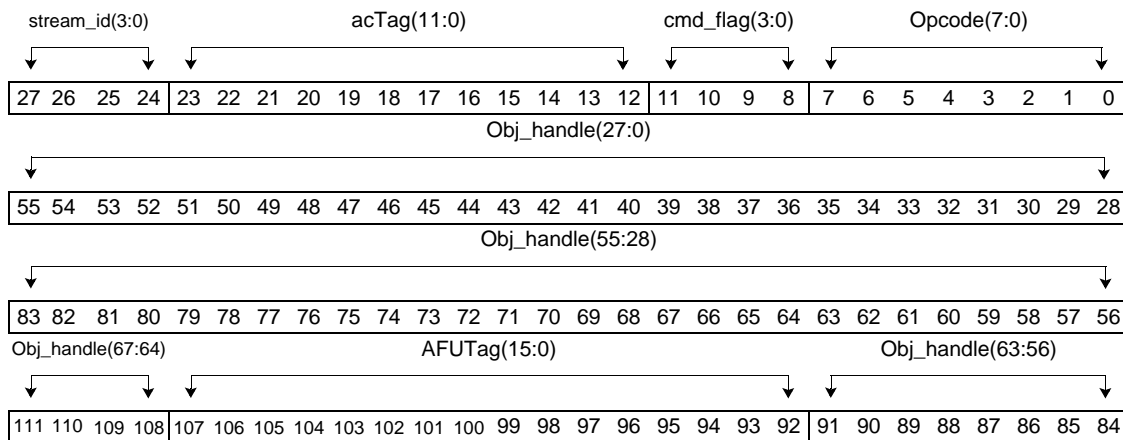
Developer Note

The command flag, data, and the object handle fields specified for this command are specified by the host's platform architecture.

The attached OpenCAPI device provides MMIO space where the combination of the object handle, data, and the command flag associated with this command are located. The manufacturer of the OpenCAPI device determines the number of command-flag and object-handle combination entries supported based on what is supported by the host and the function provided by the device.

Approved

Wake host thread	wake_host_thread	'0101 1100'
message	TLX.vc.3	4



This command is used to wake a thread on the host. The specification of the object handle and the command flag is found in the host's platform architecture.

Results are returned to the AFU using **wake_host_resp**.

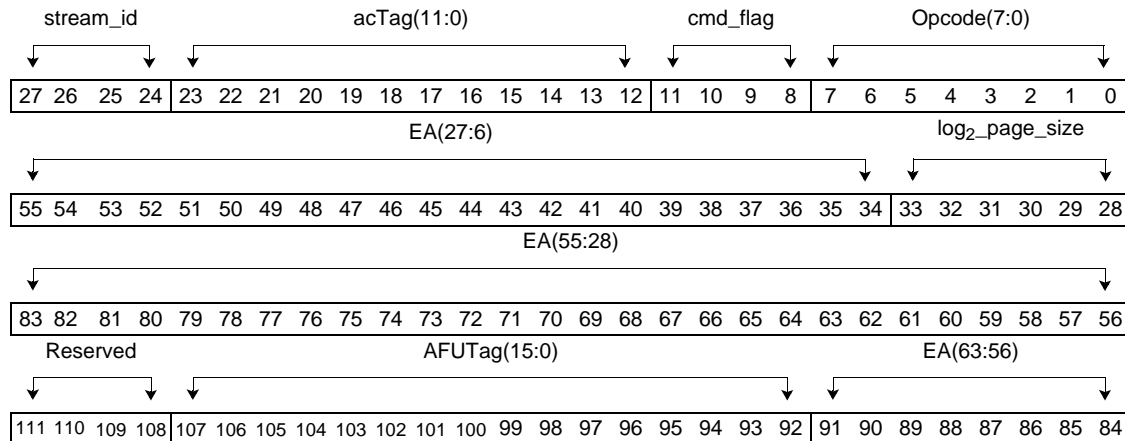
Developer Note

See the Developer note found in the description of **intrp_req** for details on the specification of the object handle and command flag and the requirements this specification places on the OpenCAPI device.

wake_host_resp indicates if the operation was successful, or if an interrupt is required.

Approved

Address translation prefetch	xlate_touch xlate_touch.n	'0111 1000' 0111 1100'
address translation management	TLX.vc.3	4



This command is used to request address translation prefetch for the address (EA) specified. The address can specify any 64-byte aligned address (EA). The `log2_page_size` field specifies the page size for an age-out ATC entry request and is reserved for an address translation request. See the specification of the command flag, bit 0, in *Table 2-7*.

- The dot-n form indicates that the results of the address translation may not be installed into the host's ATC as part of ATC miss handling.

Table 2-7 provides the specification of the `cmd_flag` field. *Figure 2-1* on page 62 provides the architectural description of the command's operation.

*Table 2-7. The cmd_flag specification for **xlate_touch** (all forms) (Page 1 of 2)*

cmd_flag bit	Description
3	Reserved.
2	0 Light-weight touch (lwt). Address translation stops and returns status if software intervention is required to complete the address translation request. Software intervention shall not be initiated.
	1 Heavy-weight touch (hwt). Address translation invokes software intervention if required to complete the address translation request. Status is returned immediately. The result of the software intervention is reported to the AFU using xlate_done .
1	0 Read-only access requested (ro). Read permission is requested by this address translation request. Write permission may be obtained.
	1 Write access requested (w). Write permission is requested by this address translation request.

Approved

Table 2-7. The *cmd_flag* specification for ***xlate_touch*** (all forms) (Page 2 of 2)

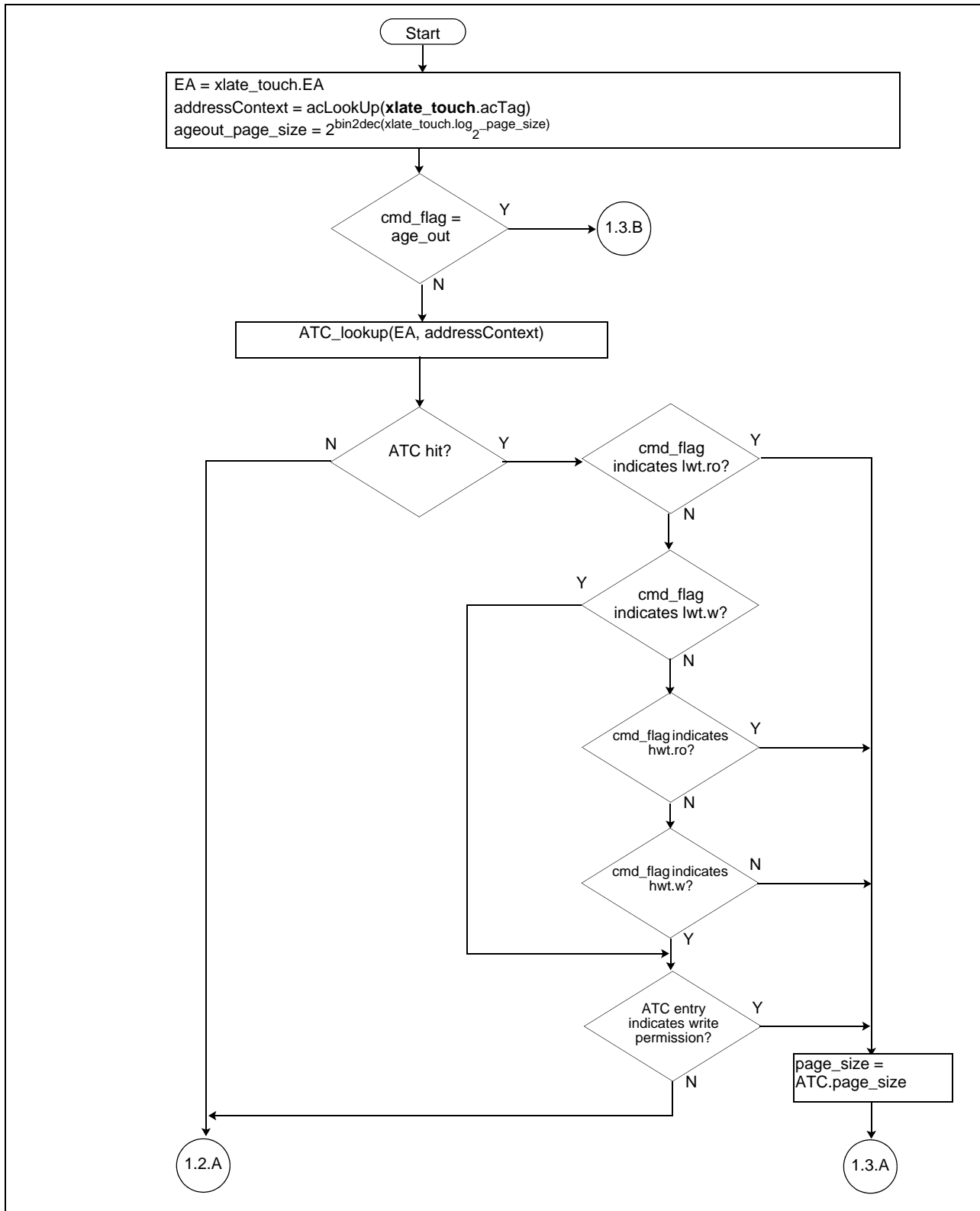
cmd_flag bit	Description
0	<p>0 Address translation request (<i>xlate</i>). <i>log₂_page_size</i> is reserved.</p> <p>1 Age-out ATC entry request (<i>age_out</i>). <i>log₂_page_size</i> specifies the page size of the entry to be aged out.</p> <div style="border: 1px solid black; padding: 5px;"> <p style="text-align: center;">Engineering note</p> <p><i>xlate_touch</i> can be used to update the LRU mechanism of the host's ATC. <i>cmd_flag</i>(0) can be used to provide hints to the host.</p> <p>0 Address translation is invoked. If an entry is found in the ATC, or one is added, that entry is marked MRU</p> <p>1 Address translation is invoked. If an entry is found in the ATC, that entry is marked as LRU.</p> <p>It is determined by the host implementation if early aging causes immediate invalidation of the matching ATC entries, the entries are marked as LRU, or no action is taken. A host implementation may chose to ignore the LRU hints described above.</p> <p>The page size associated with an EA is provided in the <i>touch_resp</i> of a previous <i>xlate_touch</i>. The OpenCAPI device shall retain this information when making an age-out request.</p> </div>
<p>cmd_flag encode specification:</p> <p>0000 <i>xlate</i>, lwt.ro 0001 <i>age out</i> 0010 <i>xlate</i>, lwt_w 0011 Reserved 0100 <i>xlate</i>, hwt.ro 0101 Reserved 0110 <i>xlate</i>, hwt.w 0111 Reserved</p> <p>The use of reserved code points results in fatal errors. See <i>Table 7-1 Error event specification</i> on page 104. The use of a dot-n form and age out is an error. See <i>Age out specified for <i>xlate_touch.n</i></i> on page 104 for details.</p>	

Developer Note

Figure 2-1 on page 62 does not show validation of the *addressContext* or the impacts of other hardware-driven events that might terminate this operation. See the specification of ***touch_resp*** for details of the result specification.

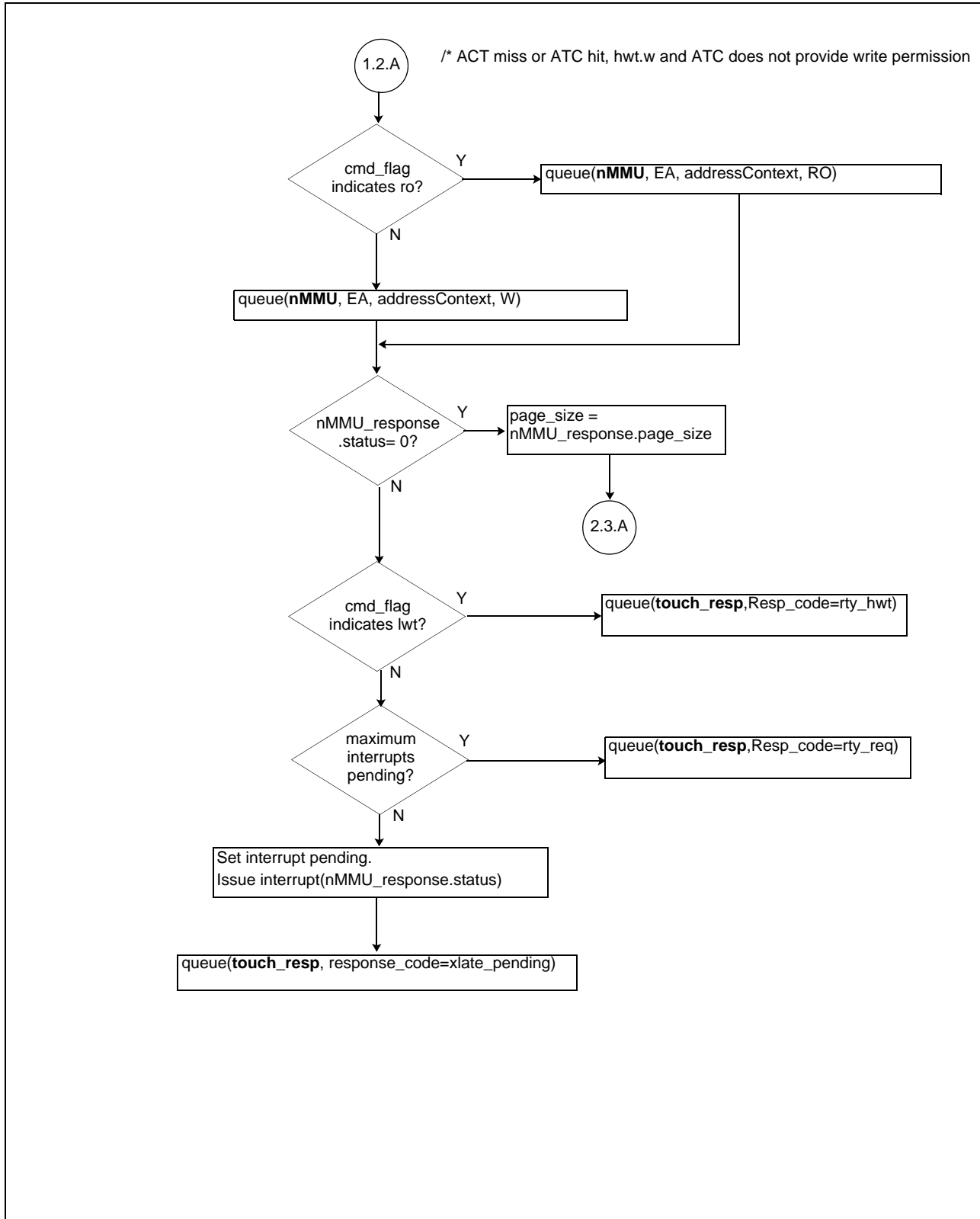
Approved

Figure 2-1. Address translation sequence: *xlate_touch* (Page 1 of 3)



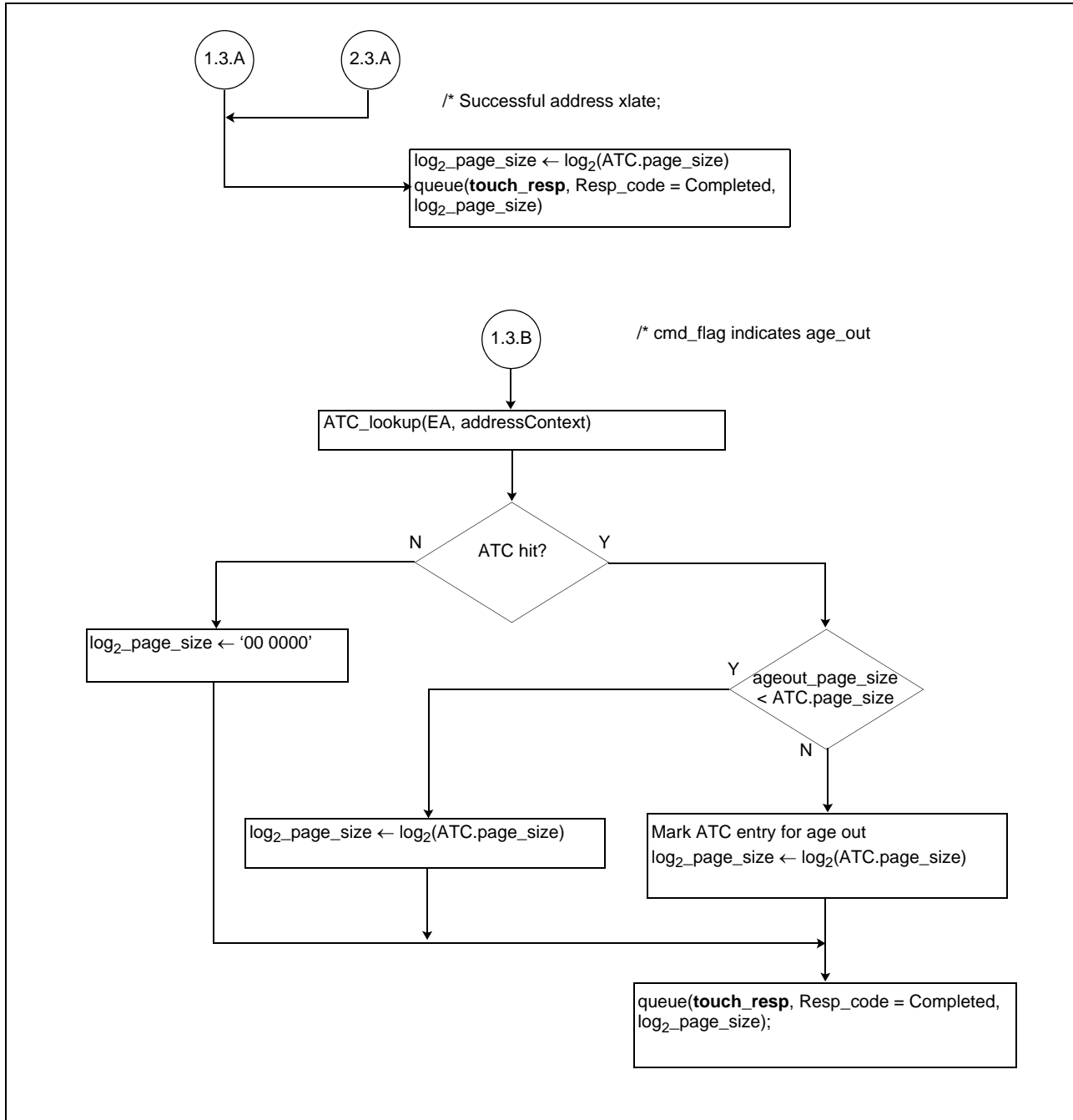
Approved

Figure 2-1. Address translation sequence: *xlate_touch* (Page 2 of 3)



Approved

Figure 2-1. Address translation sequence: *xlate_touch* (Page 3 of 3)



Approved

Results are returned to the AFU using **touch_resp**.

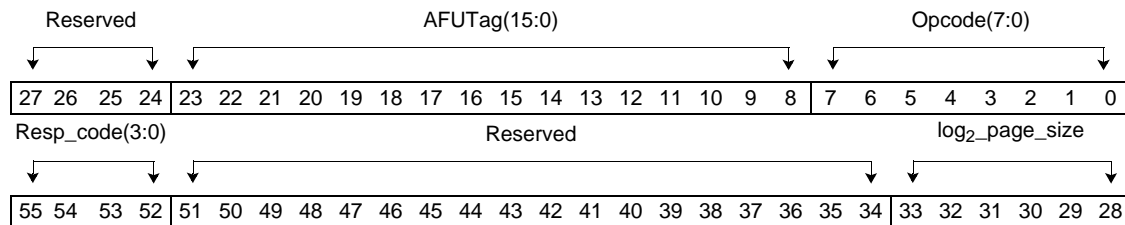
Engineering Note

The AFUtag is passed back to the TLX in a response packet and has no control function in the TL as described in *Table 2-1* on page 30. The **touch_resp**'s response code specification of `xlate_pending` indicates to the TLX that a subsequent **xlate_done** command is sent when the host is ready to service the translation request.

- The **xlate_done** command contains the AFUtag that is specified in the original **xlate_touch** command sent. While there are no requirements placed on the AFU to reserve the AFUtag used by the **xlate_touch** command until the **xlate_done** command is received by the TLX, it is strongly recommended that an AFU implementation do so. It is an AFU implementation choice to use the AFUtag to precisely determine which address translation to retry. The alternative is likely to be less efficient.
- **xlate_done** uses TL.vc.0. The implementation shall ensure that the **touch_resp** carrying the response code of `xlate_pending` is added to the VC prior to the **xlate_done**.

Approved

touch response	touch_resp	'0000 0010'
address translation management	TL.vc.0	2



This is a response to an **xlate_touch** command. `log2_page_size` is valid only when the `Resp_code` = Completed, and when `xlate_touch` specifies address translation (`xlate`), otherwise the field is reserved. The `Resp_code` field is specified in *Table 2-8*.

Table 2-8. The Resp_code specification for touch_resp

Resp_code encode	Description
'0000'	Completed. Address translation completed successfully.
'0001'	Retry using the heavy-weight touch specification (<code>rtw_hwt</code>). The translation could not be completed using the light-weight touch (<code>lwt</code>) specified by the <code>xlate_touch</code> command.
'0010'	Retry request (<code>rtw_req</code>). Indicates that the address translation could not be completed at this time. An address translation attempt may be made a later time. This is a long <i>back-off event</i> .
'0011'	Reserved.
'0100'	Translation pending (<code>xlate_pending</code>). Indicates that the address translation could not be completed. The ATC did not contain the translation, and software was invoked. An asynchronous <code>xlate_done</code> TL command shall be sent when software actions have completed. It is strongly recommended that the device wait for <code>xlate_done</code> to be received before retrying the operation. However, using a retry back off mechanism is permitted to determine when to retry the command. Such an implementation shall examine <code>xlate_done</code> for the results of the address translation and take action based on those results. <ul style="list-style-type: none"> <code>xlate_done</code> uses TL.vc.0. The implementation shall ensure that the <code>touch_resp</code> carrying the response code of <code>xlate_pending</code> is added to the VC prior to the <code>xlate_done</code>.
'0101' - '1011'	Reserved.
'1100'	Reserved.
'1101'	Reserved.
'1110'	Failed. The operation has failed and cannot be recovered. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the <code>Resp_code</code> = Failed. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved.

Note: The errors specified by `Resp_code` do not include the fatal error conditions described in *Table 7-1* on page 104.

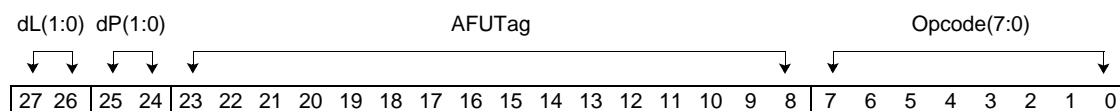
touch_resp responds to the TLX commands found in *Table 2-9*. For each command only the `Resp_codes` indicated with a Y may be used. `Resp_code` with N shall not be used.

Approved

Table 2-9. *touch_resp* Resp_code use by TLX command

TLX command	Completed (0)	rty_hwt (1)	rty_req (2)	xlate_pending (4)	Failed (14)
xlate_touch	Y	Y	Y	Y	Y
xlate_touch.n	Y	Y	Y	Y	Y

Read response	read_response	'0000 0100'
Read data return	TL.vc.0, TL.dcp.0	1



In response to a non-cacheable read command initiated by the AFU, the host is returning data. The AFU can determine which command to associate the data with by using the AFUtag provided with the command and returned with the response.

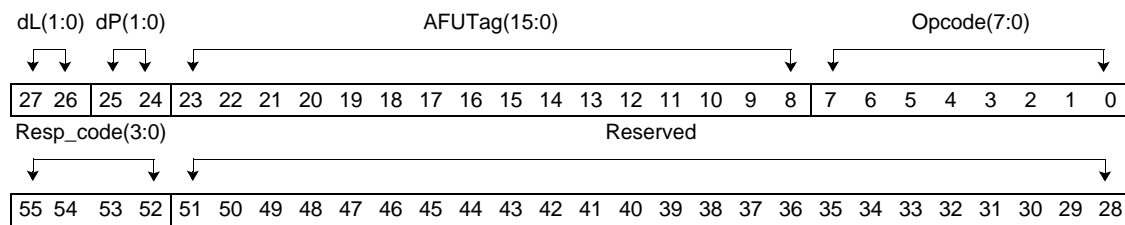
The dLength (dL) field indicates the amount of data contained in this response. Multiple read response packets may be received for a single read command. When the dLength field in the response does not match the full amount of data requested by the command, the dPart (dP) field is used to indicate the offset within the *naturally aligned data block* specified by the address in the read command. For multiple responses to a single command, the AFUtag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TLX may see the values of dPart returned in any order. When multiple responses are received for a read command, a combination of **read_response** and **read_failed** responses may be received. Taken together, all responses shall contain a combination of dPart and dLength to cover the command's dLength specification.

The dLength and dPart fields shall be specified as 64 bytes and offset at 0 for **pr_rd_wnitc** and any of the commands classified as mem_atomics that return data. A single response covers the entire operation. Data is aligned within the data flit based on the command's address bits 5:0.

This response is specified with immediate data. Credits for both the VC and DCP shall be obtained before this response is serviced by the TL.

Approved

Read failed response	read_failed	'0000 0101'
Read data return	TL.vc.0	2



In response to a read command initiated by the AFU, the host is indicating that the read failed. The AFU determines which command to associate the failure with by using the AFUTag provided with the command and returned with the response.

The dLength and dPart fields specify how much of the read operation is being reported. Multiple read response packets may be received for a single read command. When the dLength field in the response does not match the full amount of data requested by the command, the dPart field is used to indicate the offset within the *naturally aligned data block* specified by the address in the read command. For multiple responses to a single command, the AFUTag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TLX may see the values of dPart returned in any order.

- When multiple responses are received for **rd_wnitc**, a combination of **read_response** and **read_failed** responses may be received. Taken together, all responses shall contain a combination of dPart and dLength to cover the command's dLength specification.

The dLength and dPart fields shall be specified as 64 bytes and offset at 0 for **pr_rd_wnitc**, **amo_rd**, **amo_rw**, and all dot variants of these commands. A single response shall be returned for these commands.

The Resp_code field indicates the type of failure being reported. The Resp_code field is specified in *Table 2-10*.

*Table 2-10. The Resp_code specification for **read_failed** (Page 1 of 2)*

Resp_code encode	Description
'0000' - '0001'	Reserved.
'0010'	Retry request (rty_req). Use of this code point might be due to an event in the host that may require software intervention, or may indicate that address translation could not be completed for the command at this time, or may be due to a hardware recovery mechanism. The operation may be retried by the device. This is a long back-off event.
'0011'	Reserved.
'0100'	Translation pending (xlite_pending). Indicates that the address translation could not be completed. The ATC did not contain the translation, and software was invoked. An asynchronous xlite_done TL command shall be sent when software actions have completed. It is strongly recommended that an AFU wait for the xlite_done before retrying the command. However, using a retry back off timer is permitted to determine when to retry the command. Such an implementation shall examine xlite_done for the results of the address translation and take action based on those results. <ul style="list-style-type: none"> • xlite_done uses TL.vc.0. The implementation shall ensure that the read_failed carrying the response code of xlite_pending is added to the VC prior to the xlite_done.

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

Approved

Table 2-10. The *Resp_code* specification for **read_failed** (Page 2 of 2)

Resp_code encode	Description
'0101'	Reserved
'0110'	Reserved.
'0111'	Reserved.
'1000'	<p>Data error (dError). The host's protocol stack operation has completed. The data obtained by the host has been corrupted and is not correctable. This may be a recoverable error by retrying the operation. See the device documentation and <i>Section 2.1.1</i> for additional information.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering note</p> <p>A dError condition may also be reported using read_response as appropriate for the TLX command, and shall indicate that the data is bad using the bad data indication in the control flit as specified for the data carrier used.</p> </div>
'1001'	Unsupported operand length. The operation specifies an operand length that is not supported by the device. A retry of the operation shall not be successful.
'1010'	Reserved.
'1011'	Bad address specification. The address specified, by the command is not naturally aligned on a boundary specified by the operand length. Additional restrictions for address specification are specified in the operation descriptions of the TLX command amo_rd on page 51. A retry of the operation shall not be successful.
'1100'	Reserved.
'1101'	Reserved.
'1110'	<p>Failed. The operation has failed and cannot be recovered. This code point indicates that the state of the host due to the error occurrence does not allow a successful retry of the operation.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the <i>Resp_code</i> = Failed. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved.

Note: The errors specified by *Resp_code* do not include the fatal error conditions described in *Table 7-1* on page 104.

read_failed responds to the TLX commands found in *Table 2-11*. For each command only the *Resp_codes* indicated with a Y may be used. *Resp_code* indicated with an N shall not be used.

Table 2-11. **read_failed** *Resp_code* use by TLX command (Page 1 of 2)

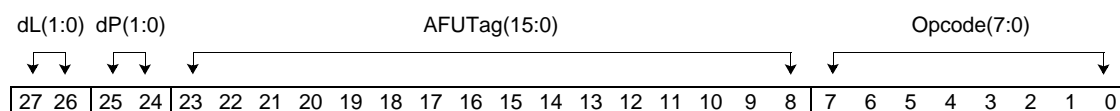
TLX command	rtly_req (2)	xlate_pending (4)	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
rd_wnitc	Y	Y	Y	N	Y	Y
pr_rd_wnitc	Y	Y	Y	N	Y	Y
rd_wnitc.n	Y	Y	Y	N	Y	Y
pr_rd_wnitc	Y	Y	Y	N	Y	Y

Approved

Table 2-11. *read_failed* Resp_code use by TLX command (Page 2 of 2)

TLX command	rty_req (2)	xlate_pending (4)	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
amo_rd	Y	Y	Y	N	Y	Y
amo_rd.n	Y	Y	Y	N	Y	Y
amo_rw	Y	Y	Y	N	Y	Y
amo_rw.n	Y	Y	Y	N	Y	Y

Write response	write_response	'0000 1000'
write response	TL.vc.0	1



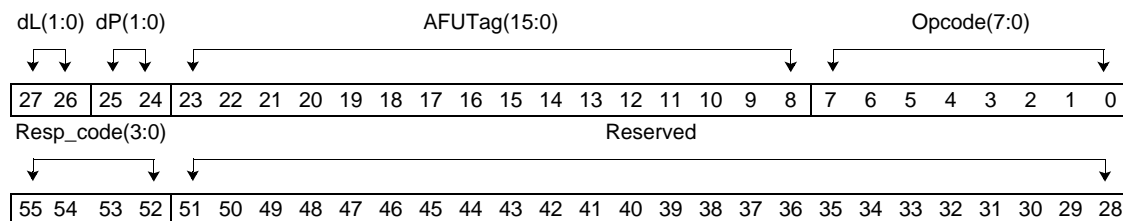
This packet is used in response to a write command (that is, **dma_w**, **dma_w.be**, **dma_pr_w**, **amo_w**) operation that has succeeded. The AFU determines which command to associate with this response by using the AFUTag provided with the command and returned with the response. Data specified by this response is global visible. That is, a subsequent read shall see the new data.

For **dma_w**, the dLength (dL) and dPart (dP) fields specify how much of the write operation is being reported. A single response may cover the entire operation. For a single response, the dLength must match the dLength specified by the command, and dPart must indicate a starting offset of 0. Multiple write response packets may be received for a single write command. When the dLength field in the response does not match the full amount of data requested by the command, the dPart field is used to indicate the offset from the starting address specified in the write command. For multiple responses to a single command, the AFUTag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TLX may see the values of dPart returned in any order. When multiple responses are received for a write command, a combination of **write_response** and **write_failed** responses may be received. Taken together, all responses shall contain a combination of dLength and dPart to cover the command's dLength specification.

For **dma_pr_w** and **amo_w** commands, only one response is expected; dLength and dPart shall be specified as 64 bytes and offset at 0.

Approved

Write failed response	write_failed	'0000 1001'
write response	TL.vc.0	2



In response to a write command initiated by the AFU, the host is indicating that the write failed. The AFU determines which command to associate the failure with by using the AFUtag provided with the command and returned with the response.

The dLength and dPart fields specify how much of the write operation is being reported. A single response may cover the entire operation. For a single response, the dLength must match the dLength specified by the command, and dPart must indicate a starting offset of 0. Multiple write response packets may be received for a single write command. When the dLength field in the response does not match the full amount of data requested by the command, the dPart field is used to indicate the offset from the starting address specified in the write command. For multiple responses to a single command, the AFUtag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TLX may see the values of dPart returned in any order. When multiple responses are received for a write command, a combination of **write_response** and **write_failed** responses may be received. Taken together, all responses shall contain a combination of dLength and dPart to cover the command's dLength specification.

This response shall be returned when the operation fails.

For **dma_w.be**, **dma_pr_w**, **amo_w**, and all dot variants of these commands, only one response shall be returned; dLength and dPart shall be specified as 64 bytes and offset at 0.

The Resp_code field indicates the type of failure being reported. The Resp_code field is specified in *Table 2-12*.

Table 2-12. The Resp_code specification of write_failed (Page 1 of 2)

Resp_code encode	Description
'0000' - '0001'	Reserved.
'0010'	Retry request (rty_req). Use of this code point might be due to an event in the host that may require software intervention, or may indicate that address translation could not be completed at this time, or may be due to a hardware recovery mechanism. The operation may be retried by the device. This is a long back-off event.
'0011'	Reserved.

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

Approved

Table 2-12. The Resp_code specification of **write_failed** (Page 2 of 2)

Resp_code encode	Description
'0100'	<p>Translation pending (xlate_pending). Indicates that the address translation could not be completed. The ATC did not contain the translation, and software was invoked. An asynchronous xlate_done TL command shall be sent when software actions have completed. It is strongly recommended that an AFU wait for the xlate_done before retrying the command. However, using a retry back off timer is permitted to determine when to retry the command. Such an implementation shall examine xlate_done for the results of the address translation and take action based on those results.</p> <ul style="list-style-type: none"> xlate_done uses TL.vc.0. The implementation shall ensure that the write_failed carrying the response code of xlate_pending is added to the VC prior to the xlate_done.
'0101' - '0110'	Reserved.
'0111'	Reserved.
'1000'	<p>Data error (dError). The host's protocol stack operation completed. The received data was UE data, or might have been marked bad in the control flit associated with the data transfer, or might have been damaged in the host. Changes, if any, to the memory location specified by the response are globally visible. The memory location shall contain SUE data.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p>Engineering note</p> <p>If an implementation is unable to modify the memory location specified by the command to contain SUE data, the implementation shall not report a dError. The implementation shall report a Failed.</p> </div> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p>Engineering note</p> <p>A dError condition may also be reported by the consumer of the data. That is, the reporting of the dError condition may be delayed until the data is consumed by a read operation. This requires that the actions taken when the error condition is detected either shall cause the memory location to contain SUE data or shall use an alternate method to report the data is invalid prior to or when it is consumed. When either of these methods are used, the host may response with write_response instead of a write_failed.</p> </div>
'1001'	Unsupported operand length. The operation specifies an operand length that is not supported by the device. A retry of the operation shall not be successful.
'1010'	Reserved.
'1011'	Bad address specification. The address specified is not naturally aligned on a boundary specified by the operand length. A retry of the operation shall not be successful.
'1100'	Reserved
'1101'	Reserved.
'1110'	<p>Failed. The operation has failed and cannot be recovered. This code point indicates that the state of the host due to the error occurrence does not allow a successful retry of the operation. This includes the following:</p> <ul style="list-style-type: none"> A dError event was detected and the implementation is unable to modify the memory location specified by the command to contain SUE data. Changes, if any to the memory location specified by the response are globally visible. The memory location may be unmodified, or may contain undefined data. Any other failure detected by the host that is not included in any of the specified response codes. The failure may cause the modification of the memory location specified by the command. Changes, if any to the memory location specified by the response are globally visible. The memory location may be unmodified, may contain undefined data, or may contain SUE data. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p>Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the Resp_code = Failed. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

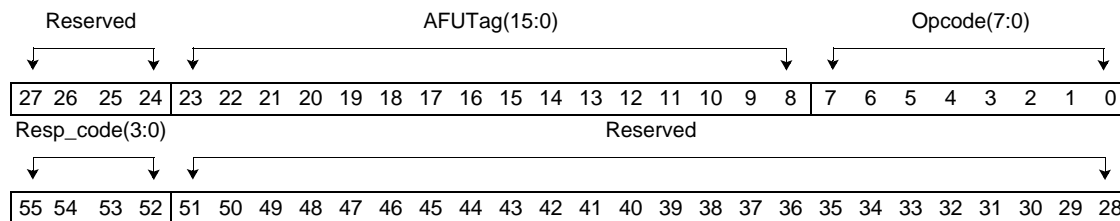
Approved

write_failed responds to the TLX commands found in *Table 2-13*. For each command only the *Resp_codes* indicated with a Y may be used. *Resp_code* indicated with an N shall not be used.

Table 2-13. write_failed Resp_code use by TLX command

TLX command	rty_req (2)	xlate_pending (4)	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
dma_w	Y	Y	Y	N	Y	Y
dma_w.n	Y	Y	Y	N	Y	Y
dma_w.be	Y	Y	Y	N	N	Y
dma_w.be.n	Y	Y	Y	N	N	Y
dma_pr_w	Y	Y	Y	N	Y	Y
dma_pr_w.n	Y	Y	Y	N	Y	Y
amo_w	Y	Y	Y	Y	Y	Y
amo_w.n	Y	Y	Y	Y	Y	Y

Interrupt response	intrp_resp	'0000 1100'
message response	TL.vc.0	2



This packet is used in response to **intrp_req**, and **intrp_req.d** commands.

The response code indicates that the interrupt was successfully initiated or provides error status. The *Resp_code* field is specified in *Table 2-14*.

Table 2-14. The Resp_code specification for intrp_resp

Resp_code encode	Description (Page 1 of 2)
'0000'	Interrupt request accepted.
'0001'	Reserved.
'0010'	Retry request (rty_req). Use of this code point might be due to an event in the host that may require software intervention, or may indicate that address translation could not be completed at this time, or may be due to a hardware recovery mechanism. The operation may be retried by the device. This is a long back-off event.
'0011'	Reserved.

Note: The errors specified by *Resp_code* do not include the fatal error conditions described in *Table 7-1* on page 104.

Approved

Table 2-14. The *Resp_code* specification for *intrp_resp*

Resp_code encode	Description (Page 2 of 2)
'0100'	Interrupt resources pending (<i>intrp_pending</i>). Indicates that the operation could not be completed at this time requiring additional software intervention. Software intervention has been successfully invoked. An asynchronous <i>intrp_rdy</i> TL command shall be sent when software actions have completed and the operation can be retried. It is strongly recommended that an AFU wait for the <i>intrp_rdy</i> before retrying the command. However, using a retry back off timer is permitted to determine when to retry the command. Such an implementation shall examine <i>intrp_rdy</i> for the results and take action based on those results. <ul style="list-style-type: none"> <i>intrp_rdy</i> uses TL.vc.0. The implementation shall ensure that the <i>intrp_resp</i> carrying the response code of <i>intrp_pending</i> is added to the VC prior to the <i>intrp_rdy</i>.
'0101' - '0110'	Reserved.
'0111'	Reserved.
'1000'	Data error (dError). Used only in response to <i>intrp_req.d</i> . The received data was corrupted and not correctable, or might have been marked bad in the control flit associated with the data transfer, or might have been damaged in the host. The operation is aborted.
'1001'	Unsupported operand length. Used only in response to <i>intrp_req.d</i> . The operation specifies an operand length that is not supported by the device. A retry of the operation shall not be successful.
'1010'	Reserved.
'1011'	Bad object handle specification. The object handle is specified by the platform architecture. A retry of the operation shall not be successful.
'1100'	Reserved.
'1101'	Reserved.
'1110'	Failed. The operation has failed and cannot be recovered. This code point indicates that the state of the host due to the error occurrence does not allow a successful retry of the operation. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the <i>Resp_code</i> = Failed. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved.

Note: The errors specified by *Resp_code* do not include the fatal error conditions described in *Table 7-1* on page 104.

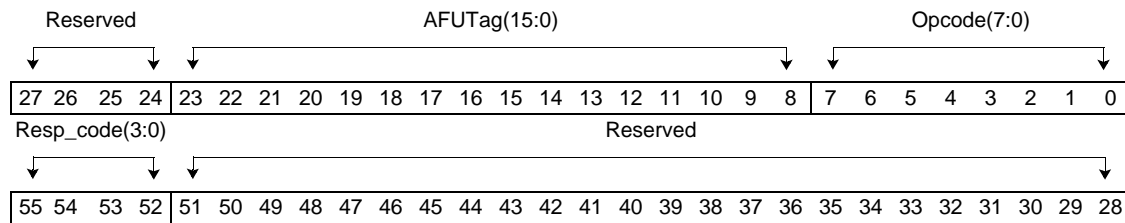
intrp_resp responds to the TLX commands found in *Table 2-15*. For each command only the *Resp_codes* indicated with a Y may be used. *Resp_code* indicated with an N shall not be used.

Table 2-15. *intrp_resp* *Resp_code* use by TLX command

TLX command	rty_req (2)	xlate_pending (4)	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
<i>intrp_req</i>	Y	Y	N	N	Y	Y
<i>intrp_req.d</i>	Y	Y	Y	Y	Y	Y

Approved

Wake Host Thread Response	wake_host_resp	'0001 0000'
message response	TL.vc.0	2



This packet is used in response to a **wake_host_thread** command. The operation in the host was either successful in waking the thread specified by the command or it was not. The Resp_code field and the reporting priority when multiple errors are detected is specified in *Table 2-16*.

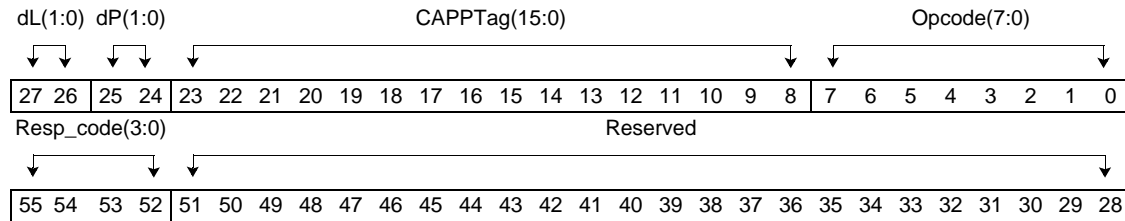
Table 2-16. The Resp_code specification for wake_host_resp

Resp_code encode	Description
'0000'	Thread found. Thread woken.
'0001'	Reserved.
'0010'	Retry request (rty_req). Indicates that the operation could not be completed at this time. The operation may be retried at a later time. This is a long <i>back-off event</i> .
'0011'	Reserved.
'0100'	Reserved.
'0101'	Thread not found. An interrupt is required to service the operation.
'0110'	Reserved.
'0111'	Reserved.
'1000' - 1010'	Reserved.
'1011'	Bad object handle specification. The object handle is specified by the platform architecture. A retry of the operation shall not be successful.
'1100'	Reserved.
'1101'	Reserved.
'1110'	Failed. The operation has failed and cannot be recovered. This code point indicates that the state of the host due to the error occurrence does not allow a successful retry of the operation. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the Resp_code = Failed. The specification of the error collection facility should be documented in the host's platform architecture.</p> </div>
'1111'	Reserved.

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

Approved

Memory read failure	mem_rd_fail	'0000 0010'
mem_response	TLX.vc.0	2



In response to a memory read command initiated by the host, the AFU is indicating that the read failed. The host determines which command to associate with the failure by using the CAPPTag provided with the command and returned with the response.

For **rd_mem**, the dLength (dL) and dPart (dP) fields specify how much of the read operation is being reported. A single response may cover the entire operation. For a single response, the dLength field must match the dLength specified by the command, and the dPart field must indicate a starting offset of 0. Multiple response packets may be received for a single memory read command. When the dLength field in the response does not match the full amount of data specified by the command, the dPart field is used to indicate the offset within the *naturally aligned data block* specified by the address in the memory read command. For multiple responses to a single command, the CAPPTag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TL may see the values of dPart returned in any order. When multiple responses are received for a memory read command, a combination of **mem_rd_response** and **mem_rd_fail** responses may be received. Taken together, all responses shall contain a combination of dLength and dPart to cover the command's dLength specification.

For the **pr_rd_mem** command and **config_read**, the dLength and dPart fields shall be specified as 64 bytes and offset at 0. A single response shall be returned. Violating this rule results in a *Bad response received* error event.

The Resp_code field indicates the type of failure being reported. The Resp_code field is specified in *Table 2-17*.

Table 2-17. The Resp_code specification for mem_rd_fail (Page 1 of 2)

Resp_code encode	Description
'0000' - '0001'	Reserved.
'0010'	Reserved.
'0011'	Reserved.
'0100' - '0111'	Reserved

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

Approved

Table 2-17. The *Resp_code* specification for *mem_rd_fail* (Page 2 of 2)

Resp_code encode	Description
'1000'	Data error (dError). The memory access completed. The data obtained by the AFU has been corrupted and is not correctable. Data is not sent to the host. This may be a recoverable error by retrying the operation. See the device documentation for additional information. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering note</p> <p>A dError condition may also be reported using mem_rd_response, and shall indicate that the data is bad using the bad data indication in the control flit as specified for the data carrier used.</p> </div>
'1001'	Unsupported operand length. The operation specifies an operand length that is not supported by the device. A retry of the operation shall not be successful.
'1010'	Reserved.
'1011'	Bad address specification. The address specified is not naturally aligned on a boundary specified by the operand length. A retry of the operation shall not be successful.
'1100' - '1101'	Reserved.
'1110'	Failed. The operation has failed and cannot be retried. This code point indicates that the state of the device due to the error occurrence does not allow a successful retry of the operation. This includes the following: <ul style="list-style-type: none"> The device and function number specified in the address of a config_read is not recognized by the AFU. config_read is issued with T=1. Any other failure detected by the AFU that is not included in any of the specified response codes. <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the Resp_code = Failed. The specification of the error collection facility should be documented in the device documentation.</p> </div>
'1111'	Reserved.

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

mem_rd_fail responds to the TL commands found in *Table 2-18*. For each command only the Resp_codes indicated with a Y may be used. Resp_code indicated with an N shall not be used.

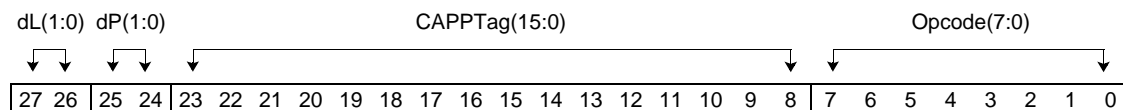
Table 2-18. *mem_rd_fail* Resp_code use by TL command

TL command	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
rd_mem	Y	Y ¹	Y	Y
pr_rd_mem	Y	Y ¹	Y	Y
config_read	Y	Y	Y	Y

1. May occur during MMIO space read access only.

Approved

Memory write response	mem_wr_response	'0000 0100'
mem_response	TLX.vc.0	1

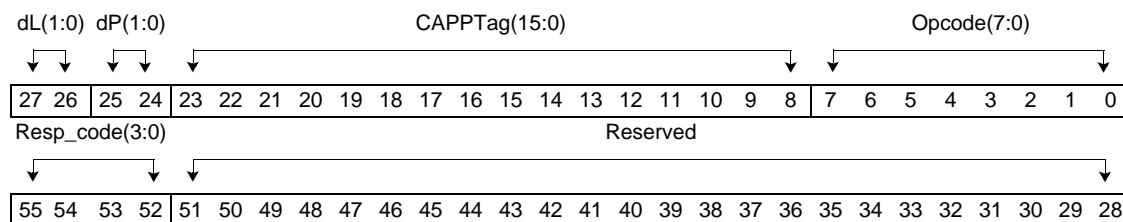


This packet is used in response to a command that writes to memory. This response is used to indicate the successful completion of all or a portion of the operation. Data specified by this response is global visible. That is, a subsequent read shall see the new data.

For **write_mem**, the dLength (dL) and dPart (dP) fields specify how much of the write operation is being reported. A single response may cover the entire operation. For a single response, the dLength must match the dLength specified by the command, and dPart must indicate a starting offset of 0. Multiple response packets may be received for a single memory write command. When the dLength field in the response does not match the full amount of data specified by the command, the dPart field is used to indicate the offset from the starting address specified in the memory write command. For multiple responses to a single command, the CAPPTag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TL may see the values of dPart returned in any order. When multiple responses are received for a memory write command, a combination of **mem_wr_response** and **mem_wr_fail** responses may be received. Taken together, all responses shall contain a combination of dLength and dPart to cover the command's dLength specification.

For the **pr_wr_mem**, **write_mem.be**, and **config_write** commands a single response shall be returned. The dLength and dPart fields shall be specified as 64 bytes and offset at 0.

Memory write failed	mem_wr_fail	'0000 0101'
mem_response	TLX.vc.0	2



In response to a **write_mem** command initiated by the host, the AFU is indicating that the write failed. The host determines which command to associate with the failure by using the CAPPTag provided with the command and returned with the response.

The dLength (dL) and dPart (dP) fields specify how much of the write operation is being reported. A single response may cover the entire operation. For a single response, the dLength must match the dLength specified by the command, and dPart must indicate a starting offset of 0. Multiple response packets may be received for a single memory write command. When the dLength field in the response does not match the full amount of data specified by the command, the dPart field is used to indicate the offset from the starting address specified in the memory write command. For multiple responses to a single command, the CAPPTag is unchanged. That is, the dLength may vary and the dPart shall vary when multiple responses are returned

Approved

for a single command. For multiple responses to a single command, there is no order requirement placed by the architecture. That is, the TL may see the values of dPart returned in any order. When multiple responses are received for a memory write command, a combination of **mem_wr_response** and **mem_wr_fail** responses may be received. Taken together, all responses shall contain a combination of dLength and dPart to cover the command's dLength specification.

For **config_write**, **pr_wr_mem** and **write_mem.be**, only one response shall be returned; dLength and dPart shall be specified as 64 bytes and offset at 0.

The Resp_code field indicates the type of failure being reported. The Resp_code field is specified in *Table 2-19*.

Table 2-19. The Resp_code specification for mem_wr_fail (Page 1 of 2)

Resp_code encode	Description
'0000' - '0001'	Reserved.
'0010'	Reserved.
'0011'	Reserved.
'0100' - '0111'	Reserved
'1000'	<p>Data error (dError). The AFU's operation has completed. The data sent by the TL was corrupted prior to the completion of the operation. Changes, if any to the memory location specified by the response are globally visible.</p> <ul style="list-style-type: none"> • Memory locations specified by the command's PA and length that correspond to system memory space shall contain SUE data. • Memory locations specified by the command's PA and length that correspond to either MMIO or configuration space may be unmodified, may contain undefined data, or may contain SUE data. <p>The corruption of the data might have occurred anywhere in the AFU's hardware or might have been detected by a bad data indication.</p> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering note</p> <p>If an implementation is unable to modify the memory location specified by the command's PA that correspond to system memory space to contain SUE data, the implementation shall not report a dError. The implementation shall report a Failed.</p> </div> <div style="border: 1px solid black; padding: 5px; margin: 10px 0;"> <p style="text-align: center;">Engineering note</p> <p>A dError condition may also be reported by the consumer of the data. That is, the reporting of the dError condition may be delayed until the data is consumed by a read operation. This requires that the actions taken when the error condition is detected either shall cause the memory location to contain SUE data or shall use an alternate method to report the data is invalid prior to or when it is consumed. When either of these methods are used, the AFU may response with mem_wr_response, instead of a mem_wr_fail.</p> </div>
'1001'	Unsupported operand length. The operation specifies an operand length that is not supported by the device. A retry of the operation shall not be successful.
'1010'	Reserved.
'1011'	Bad address specification. The address specified is not naturally aligned on a boundary specified by the operand length. A retry of the operation shall not be successful.
<p>Note: The errors specified by Resp_code do not include the fatal error conditions described in <i>Table 7-1</i> on page 104.</p>	

Approved

Table 2-19. The *Resp_code* specification for *mem_wr_fail* (Page 2 of 2)

Resp_code encode	Description
'1100 - '1101'	Reserved.
'1110'	<p>Failed. The operation has failed and cannot be retried. This code point indicates that the state of the device due to the error occurrence does not allow a successful retry of the operation. This includes the following:</p> <ul style="list-style-type: none"> • The device and function number specified in the address of a config_write is not recognized by the AFU. • config_write specified with T=1. • A dError event was detected and the implementation is unable to modify the memory locations specified by the command's PA and length that correspond to system memory space to contain SUE data. Changes, if any to the memory location specified by the response are globally visible. The memory location may be unmodified, or may contain undefined data. • Any other failure detected by the AFU that is not included in any of the specified response codes. The failure may cause the modification of the memory location specified by the command. Changes, if any to the memory location specified by the response are globally visible. The memory location may be unmodified, may contain undefined data, or may contain SUE data. <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p style="text-align: center;">Engineering Note</p> <p>It is strongly recommended that an implementation provide error collection facilities to indicate the reason for the Resp_code = Failed. The specification of the error collection facility should be documented in the device documentation.</p> </div>
'1111'	Reserved.

Note: The errors specified by Resp_code do not include the fatal error conditions described in *Table 7-1* on page 104.

mem_wr_fail responds to the TL commands found in *Table 2-20*. For each command only the Resp_codes indicated with a Y may be used. Resp_code indicated with an N shall not be used.

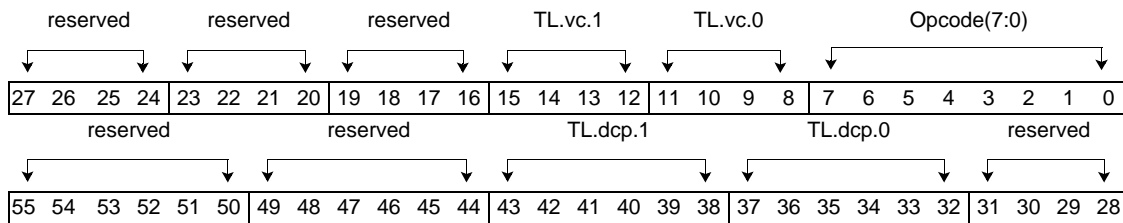
Table 2-20. *mem_wr_fail* Resp_code use by TL command

TL command	dError (8)	Unsupported operand length (9)	Bad address specification (11)	Failed (14)
write_mem	Y	Y ²	Y	Y
write_mem.be	Y	N	Y	Y ²
pr_wr_mem	Y	Y ¹	Y	Y
config_write	Y	Y	Y	Y

1. Unsupported operand length may occur only when target memory is defined as MMIO space and the command's specified pLength is not supported at the address specified.
2. May occur during MMIO address space write operation only.

Approved

Return TL credits	return_tl_credits	'0000 1000'
credit return	NA	2



This response packet is used by the TLX to return VC and DCP credits to the TL. There is no VC associated with this response, and credits are not required to service this response. Each TL.* field contains the number of credits being returned.

This response packet shall be placed only in slots 1 to 0 of any control flit using a template which specifies those slots as a 2-slot or larger location.

TL.vc.{0, 1} and TL.dcp.{0, 1} credits are returned. TL credits are for resources owned by the TLX that the TL consumes. The TLX controls the total number of credits for each of the VC and DCP it provisions the TL with.

3. Virtual channel and data credit pool specification

Commands and responses are assigned to virtual channels (VCs) to allow ordering and the specification of servicing service queues and virtual queues. Credits to use these VCs are managed by the destination (where the resources are consumed) and are released to the source (where the command or response originates) when the resource is available. Each VC has its own credit pool, which may be varied by the destination. VCs in each direction are numbered; the value assigned is used only for differentiation between VC. Each VC credit permits the sending of one command or response.

The following VC descriptions use the specific command or response or the classification of the command or response. See the command and response descriptions in *Section 2 TL and TLX command and response specifications* on page 30 for command and response VC classifications.

VCs are specified between the TL and TLX and between the TLX and the TL. Ordering is maintained within a VC and is assured between these endpoints only. VCs cannot block each other; blocking within a VC may occur due to ordering requirements. With the exception of *Command ordering* on page 26, any queuing or ordering occurring in the upper protocol layers (host bus and AFU) is not assured to be retained. Synchronization points are managed at the interfaces between the TL and host bus protocol stack and between the TLX and AFU protocol stack.

Data credit pool (DCP) credits are required when a command or response has immediate data; that is, data that is associated with the sending of the command or response. For example, a write command has immediate data while a read command does not. Commands and responses with immediate data shall obtain the necessary credits for both the VC and DCP assigned in an atomic fashion.

For example, sending 128 bytes requires atomically obtaining two DCP credits. See *Section 5.1.3 Data transport, order, and alignment* on page 97 for details on the relationship between DCP credits and *data carrier* use. If the credits are not available, the command or response cannot be serviced. That is, the command or response shall not be placed into a DL frame for transmission. See the description of *DCP* on page 18.

The order of data sent to the destination is the same as the order of the data's corresponding command or response that is sent to its destination. This allows the destination to use the arrival order of commands and responses with immediate data to associate the immediate data with its command or response.

Credits are released to the consumer of the credits by the owner of the resources that the credits represent. That is, TLX credits represent TL resources that are consumed by the TLX. TL credits represent TLX resources that are consumed by the TL. The responses used to return credits are **return_tlx_credits** and **return_tl_credits**. A compliant implementation shall:

- Provide a 16-bit counter to track credits available for use for each DCP and VC specified in the following sections.
- Provision a minimum of one and less than 64K credits for each VC specified in the following sections.
- Provision a minimum of four and less than 64K credits for each DCP specified in the following sections.

Developer Note

- The host is required to release credits only for the VC/DCP that will be used by the AFU. Releasing credits for VC/DCPs that are not going to be used might not be optimal because the released credits correspond to resources in the host that could have been used for actual command/data/response traffic from the AFU.
- The AFU is required to release credits only for the VC/DCP that will be used by the host. Releasing credits for VC/DCPs that are not going to be used might not be optimal because the released credits correspond to resources in the AFU that could have been used for actual command/data/response traffic from the host.

The credits required can be determined by knowing the AFU capabilities as described in *Section 1.2 Host operation modes* on page 25, the commands and responses used, and the associated VC and DCP.

Table 3-3 on page 87 specifies the assignment of VC and DCP to commands and responses.

Developer Note

The architecture permits an implementation to release a minimum-total of four DCP credits. Releasing a minimum-total of a single DCP credit was considered and discarded because this limits the data to a single 64-byte transfer. This was considered too restrictive and might not easily enable implementations to optimize for the data block size normally used by the implementation.

3.1 Virtual channel

Ordering shall be maintained between elements within a VC. Blocking between VCs shall not be permitted.

3.1.1 TLX command and response VC (TLX.vc)

The VC is directed from the AP to the host. Two VCs are specified {0,3}. VC credits are consumed by the TLX and are returned by the TL using **return_tlx_credits**.

Engineering Note

TLX VC credits represent resources in the TL used for processing TLX commands and responses. Each credit released by the TL represents a *unique* resource and shall not be shared with any other VC. Doing so would result in breaking the accounting rules implied by this specification. For example, if the TL used the same resource for two different VCs, the actual credit available for VCs using the shared resource would also be diminished, and the TLX would be unaware of this change.

The ability of the TL to return TLX VC credits shall not be dependent on any action taken by the TLX, including the return of TL credits.

3.1.2 TL command and response VC (TL.vc)

The VC is directed from the host to the AP. Two VCs are specified {0, 1}. VC credits are consumed by the TL and are returned by the TLX using **return_tl_credits**.

Engineering Note

TL VC credits represent resources in the TLX used for processing TL commands and responses. Each credit released by the TLX represents a *unique* resource and shall not be shared with any other VC. Doing so would result in breaking the accounting rules implied by this specification. For example, if the TLX used the same resource for two different VCs, the actual credit available for VCs using the shared resource would also be diminished, and the TL would be unaware of this change.

The ability of the TLX to return TL VC credits shall not be dependent on any action taken by the TL, including the return of TLX credits.

3.1.3 VC credit count specification

Table 3-1 specifies the maximum number of VC credits supported by an OpenCAPI-compliant device. A device is not required to release or support the maximum count. However, the consumer of the credit shall provide a counter that supports the architected maximum count.

Table 3-1. VC maximum credit count specification

VC	Maximum credit count
TLX.vc.0	64K-1
TLX.vc.3	64K-1
TL.vc.0	64K-1
TL.vc.1	64K-1

3.2 Data credit pool

3.2.1 TLX data DCP (TLX.dcp)

The data credit pool is used when moving immediate data from the AP to the host. Two DCPs are specified {0, 3}. DCP credits are consumed by the TLX and returned by the TL using **return_tlx_credits**.

Engineering Note

TLX DCP credits represent resources in the TL used for accepting data from the TLX. Each credit released by the TL represents a *unique* data resource and shall not be shared with any other DCP. Doing so would result in breaking the accounting rules implied by this specification. For example, if the TL used the same resource for two different DCPs, the actual credit available for DCPs using the shared resource would also be diminished, and the TLX would be unaware of this change.

See Section 5.1.3 *Data transport, order, and alignment* on page 97 for the association between a DCP credit and the immediate data associated with the command or response.

The ability of the TL to return TLX DCP credits shall not be dependent on any action taken by the TLX, including the return of TL credits.

3.2.2 TL data DCP (TL.dcp)

This data credit pool is used when moving immediate data from the host to the AP. Two DCPs are specified {0, 1}. DCP credits are consumed by the TL and are returned by using the TLX using **return_tl_credits**.

Engineering Note

TL DCP credits represent resources in the TLX used for accepting data from the TL. Each credit released by the TLX represents a *unique* data resource and shall not be shared with any other DCP. Doing so would result in breaking the accounting rules implied by this specification. For example, if the TLX used the same resource for two different DCPs, the actual credit available for DCPs using the shared resource would also be diminished, and the TL would be unaware of this change.

See *Section 5.1.3 Data transport, order, and alignment* on page 97 for the association between a DCP credit and the immediate data associated with the command or response.

The ability of the TLX to return TL DCP credits shall not be dependent on any action taken by the TL, including the return of TLX credits.

3.2.3 DCP credit count specification

Table 3-2 specifies the maximum number of DCP credits supported by an OpenCAPI-compliant device. A device is not required to release or support the maximum count. However, the consumer of the credit shall provide a counter that supports the architected maximum count.

Table 3-2. DCP maximum credit count specification

DCP	Maximum credit count
TLX.dcp.0	64K-1
TLX.dcp.3	64K-1
TL.dcp.0	64K-1
TL.dcp.1	64K-1

Table 3-3. Summary VC and DCP assignments (Page 1 of 2)

VC	Classification/ command	DCP	Comments	Command	Response
TL.vc.0	Touch response				X
TL.vc.0	DMA read (response)	TL.dcp.0	TL.dcp.0 is used only for read_response and is not used for read_failed .		X
TL.vc.0	Write response (OK and failed)				X
TL.vc.0	Interrupt and wake host thread responses		Message responses.		X
TL.vc.0	xlate_done		Asynchronous notification. Asynchronous command reporting the status of a previously AFU-initiated action (address translation touch).	X	
TL.vc.0	Interrupt ready		Asynchronous notification.	X	
TL.vc.1	MEM read both multiples of 64 bytes or partial commands			X	
TL.vc.1	MEM read		An MMIO read operation uses a pr_rd_mem CAPP command.	X	
TL.vc.1	MEM write	TL.dcp.1	An MMIO write operation uses a pr_w_mem CAPP command.	X	

Approved

Table 3-3. Summary VC and DCP assignments (Page 2 of 2)

VC	Classification/ command	DCP	Comments	Command	Response
TL.vc.1	Configuration register read			X	
TL.vc.1	Configuration register write	TL.dcp.1		X	
TLX.vc.0	MEM read response; for example, mem_response	TLX.dcp.0	mem_rd_fail does not use TLX.dcp.0.		X
TLX.vc.0	MEM write response				X
TLX.vc.3	Non-cacheable read and write operations - all forms			X	
TLX.vc.3	acTag management			X	
TLX.vc.3	Interrupts and wake host thread	TLX.dcp.3	Message commands. TLX.dcp.3 is used only for intrp_req.d commands.	X	
TLX.vc.3	xlate_touch - ATC prefetch			X	

3.3 TL Virtual channel and service queues

3.3.1 Host TLX command handling

Figure 3-1 on page 89 illustrates the steps a TLX command follows from the VC queue in the TL to its service queue. Commands are removed from the DL frame in slot order from a control flit and are loaded into the VC queue specified by the command. After the command reaches the head of the VC queue, it is examined.

1. Any error found in the command entry is noted. Errors found at this point in the flow shall not be reported until the command reaches the head of the service queue.
2. If the command is **assign_actag**, the command shall be removed from the VC queue and shall be serviced at this time. This service results in the update of the acTag table. Subsequent commands shall see the new state of the acTag table. The **assign_actag** command is removed from the head of the service queue.

If the command is not **assign_actag**, the command shall be serviced as follows:

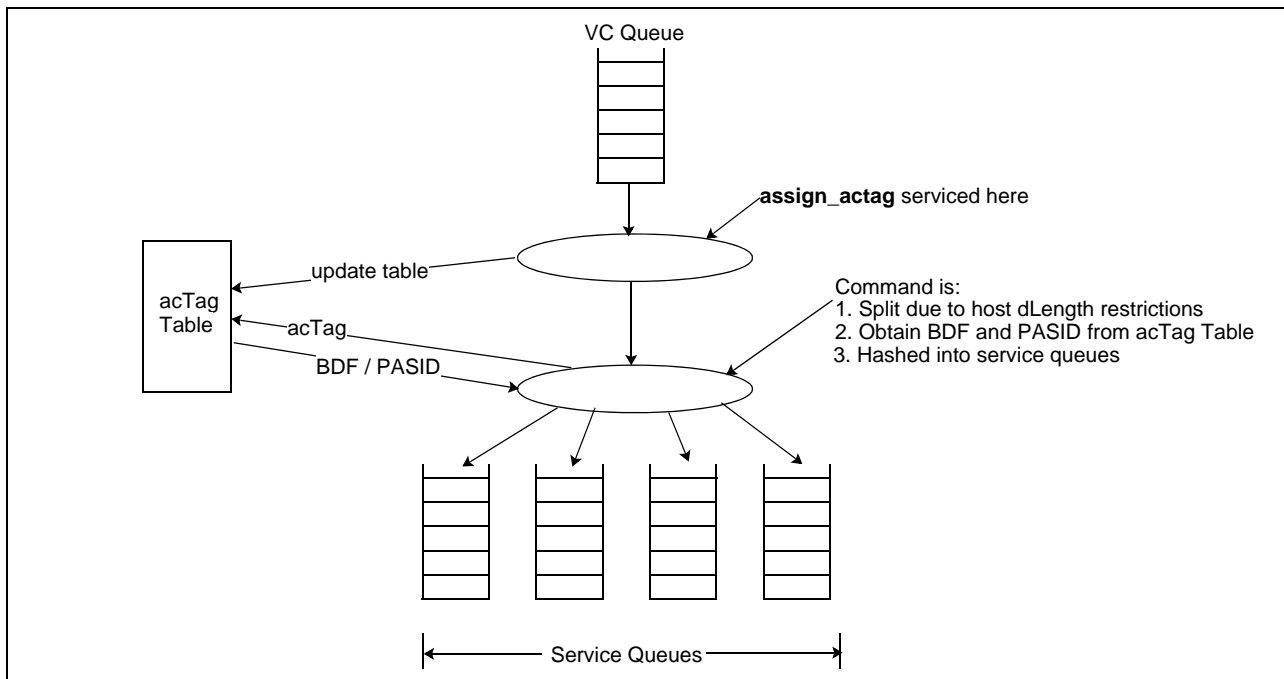
- a. Commands with a dLength specified that is larger than the host's maximum supported length for a single operation are split into commands with dLength set equal to or less than the host's maximum supported length. The address is adjusted for each command to account for the dLength serviced. The AFUtag, acTag, and stream_id are obtained from the original command. Split command entries are noted with the dPart value to be used when returning responses. Split commands are serviced based on increasing address order. All commands continue with the next step.
- b. Using the acTag found in the command entry, the BDF and PASID are obtained. An error detected at this step is fatal. See the description of *Bad BDF and PASID combination on page 105* for additional details.

Approved

- c. A VC- and implementation-specific hash is used to determine the service queue to add the command to. See the specification of the hash applied to the command in the definition of a *service queue* on page 21.

The command is then added to the service queue determined in step c and removed from the head of the VC queue.

Figure 3-1. TL command flow from the VC queue to the service queue TLX VC queues shown



As defined in *Terms* on page 17, the difference between a *virtual queue* and a *service queue* is that a *service queue* may contain multiple *virtual queues*.

The architectural model of a service queue specifies that the commands in the body of the service queue may receive the following services in the following order:

1. Error checking of the command.
2. Address translation as required by the command's specification.

Errors occurring in either of these services shall be noted in the service queue entry and shall not be reported until the command reaches the head of the service queue. The results of the address translation shall be noted in the service queue entry.

The architectural model of a service queue specifies that the head of the service queue shall receive the following services in the following order:

1. Error checking of the command. The check may occur here, or the error noted in the entry may be used to determine the error checked state of the command.
2. Address translation as required by the command's specification. The translation may occur here, or may have previously occurred. The results of the translation noted in the entry may be used at this time.
3. When the previous two services are completed, the command shall be removed from the head of the service queue.

Approved

Step 3 requires that a command at the head of the service queue be processed with the indicated precedence as follows:

1. Failed due to an error, which may be due to either an error found with the command or failed address translation attempt.

And

- That error shall result in either a response being returned to the requester or an error event being asserted,

Or

2. The operation is completed³ and a response shall be returned to the requester.

Or

3. Shall be dispatched to the host protocol layer.

Engineering Note

The architecture *requires* that the implementation of a service queue shall appear to behave as if the architectural model of a service queue is implemented.

3.3.2 Host TLX response handling

TLX responses are assigned to a VC and are removed from the DL frame in the same manner as TLX commands. Handling of TLX responses is simpler in that TLX.vc.0 which is used for TLX responses has no hash involved when selecting a service queue. Responses assigned to TLX.vc.0 are passed directly to a dedicated service queue and allowed to dispatch to the host interface.

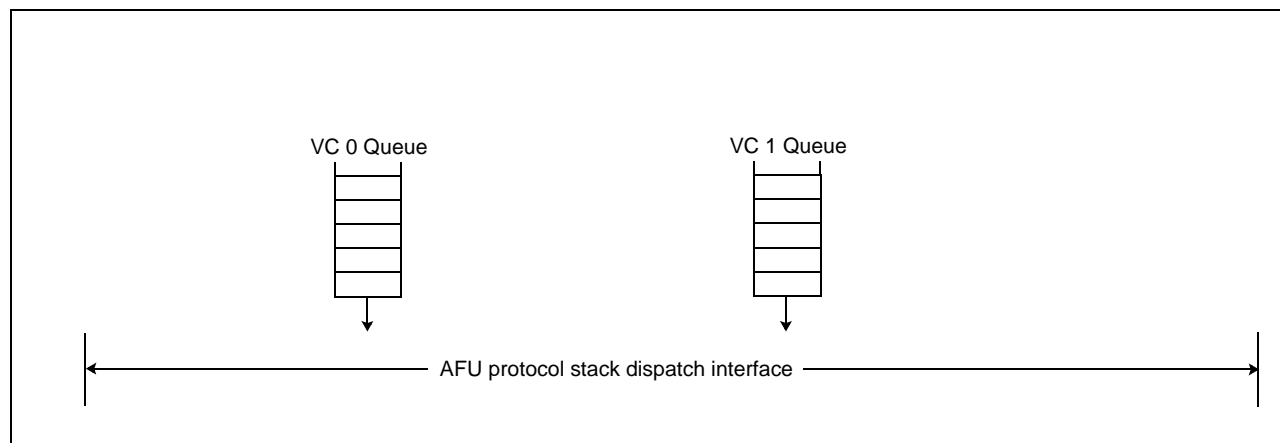
3.4 Device TL virtual channel queues

Figure 3-2 on page 91 shows the steps a TL VC queue entry follows from the time it is dequeued from the TL VC queue until it is dispatched to the AFU protocol stack interface. The TL VC queue entry can be a TL command- or response-packet. TL command- and response-packets are removed from the DLX frame in slot order from a control flit. They are loaded as queue entries into the TL VC queue specified by the command or response. After the queue entry reaches the head of the TL VC queue, it is examined.

1. Any error found in the queue entry shall be reported, and the entry shall be dequeued. If the TL VC queue entry is a TL command-packet, the errors shall be reported either by a response returned or by an error event. If the TL VC queue entry is a response-packet, errors shall be reported through error events such as *Bad response received on page 106*. Other error events are described in *Section 7.1* beginning on page 104.
2. With error checking completed, the entry shall be dequeued and sent to the AFU protocol dispatch interface. At the dispatch interface, arbitration between the VC queues is implementation dependent and beyond the scope of this specification.

3. For example, the command required only an address translation action, such as **xlate_touch**, or the address translation was not successful or is required to be retried.

Figure 3-2. TLX command and response flow from the VC to the AFU protocol stack TL VC queues shown



3.5 Virtual channel dependency rules

Commands and responses are assigned to virtual channels. AFU and host designs are cautioned not to implement a design where there is a potential for a dead lock when forward progress of one command is dependent on the forward progress of another.

The TL specification specifies natural VC dependencies. For example a TLX **rd_wnitc** using TLX.vc.3 is dependent on a TL **read_response** using TL.vc.0 to complete the operation. Deadlock conditions can occur when:

- The host cannot respond to a TLX command without issuing a TL command and the host design does not allow the TL command to be issued.
- The AFU cannot respond to a TL command without issuing a TLX command and the AFU host design does not allow the TLX command to be issued.

Since the implementation of the host and the AFU are beyond the scope of this specification, the designers of the host and the AFU shall ensure that such dead lock conditions are prohibited by the implementation.

Figure 3-3 illustrates the interdependencies between VC that occur in the existing specification. An implementation, of either a host or device, shall not introduce dependencies not shown in *Figure 3-3*.

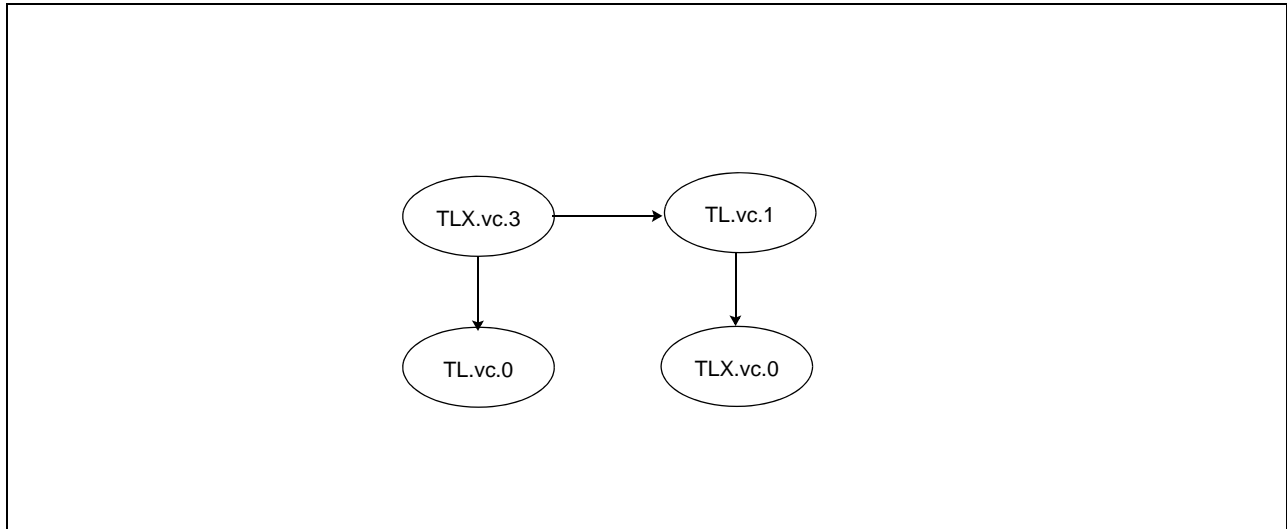
Graphically, the dependencies are shown as a line between two ovals. The ovals specify a virtual channel. For illustrative purposes, consider two ovals connected by an arrow. The tail of the arrow is connected to the oval indicating a TL- or TLX-packet using VC.a. The head of the arrow is connected to a TL- or TLX-packet using VC.b. Reading the graphic becomes:

Servicing an incoming packet using VC.a requires issuing a packet using VC.b.

The intent of the architecture is to require implementations to insure that servicing an incoming packet using VC.a does not require the use of the same resources that are required to issue a packet using VC.b. This intent might restrict implementation choices.

Approved

Figure 3-3. VC dependency graph



4. The acTag table

The section provides the architectural specification of the acTag table. Implementations may choose to implement this table using any method. However, the externally observable behavior of the table and contents of the table shall, at a minimum, comply to this architectural specification.

The acTag table shall be included in the host's implementation and shall provide a minimum of one entry.

4.1 acTag table contents

Each entry of the acTag table (acTag entry) shall contain the following fields:

- 1-bit entry valid. Architected states are defined as {valid, invalid}
- 16-bit BDF
- 20-bit PASID

An implementation may choose to add additional fields to the acTag entry.

4.2 acTag table access

The OpenCAPI device maintains a copy of each acTag entry to determine the acTag value used in TLX commands specified with an acTag field.

The acTag table is accessed using the acTag as follows:

- Read access uses the acTag provided in a TLX command specified with an acTag field as an index into the table to locate the acTag entry to be read. Reading an acTag entry returns the entry valid indication and, when valid, a BDF and PASID.
- Write access uses the acTag provided in the **assign_actag** command as an index into the table to locate the acTag entry to be updated. The command provides a BDF and PASID, which are loaded into the acTag entry. Successful completion of the write access sets the entry valid bit to the valid state.

4.2.1 Error cases when accessing the acTag table

Error	Action
acTag entry not valid	This is a fatal error. See "acTag specified in the command points to an invalid entry" in <i>Table 7-1</i> on page 104.
Address context not valid	This is determined either at the time the acTag entry is created or when the acTag entry is used to obtain the address context. This is a fatal error. See "Bad BDF and PASID combination" in <i>Table 7-1</i> on page 104.

4.3 acTag entry management

The entries of the acTag table are managed by the attached OpenCAPI device. The OpenCAPI device shall maintain its own mapping between an acTag and the contents of the acTag entry held in the host's acTag table.

Approved

During configuration, the host writes to the OpenCAPI device's configuration space to indicate the maximum size of the acTag, which indirectly specifies the size of the host's acTag table. An acTag size of 0 indicates a single entry acTag table, and the only valid acTag value is '0'. The OpenCAPI device may use any value within the allowed range to specify an acTag entry and subsequently refer to that entry using a TLX command specifying the acTag.

The OpenCAPI device learns its bus number from the address specified in a **config_write**, T=0 command. Device and function numbers are assigned by the OpenCAPI device and discovered by the host during configuration. See the specification of **config_write** for the format of the address field that contains the bus, device and function numbers.

The OpenCAPI device is configured with one or more PASIDs during initialization and operation.

After an acTag entry is set to a valid state, it is set to an invalid state only by the host upon detection of a link failure that requires resetting the OpenCAPI interface and the OpenCAPI device.

Engineering note

In a host implementation, a configuration register would be useful to vary the size of the acTag table, which allows stress testing of the design. A proposed specification of the configuration register follows:

The register contains a value N where the size of the acTag table is 2^N and the acTag range is limited, as described previously, to a range of $0..2^{N-1}$.

Permitted access methods for this configuration register, as well as the register contents, are specified by the platform architecture.

5. DL frame format

The TL and TLX contain framer and parser functions that work with the DL frame format. The DL frame format is specified as a set of 64-byte flits. There are two types of flits:

- Control flits. The control flit contains TL command/response content and DL content. The DL content contains several DL-generated subfields including the CRC that covers the control flit and any *preceding* data flits. There are fields in the DL content that are generated by the TL. For more information, see *Section 5.1.1 DL content* on page 96.
- Data flits. There are 0 to 8 data flits between each control flit.

Table 5-1. DL frame format showing CRC and “bad data flit” coverage

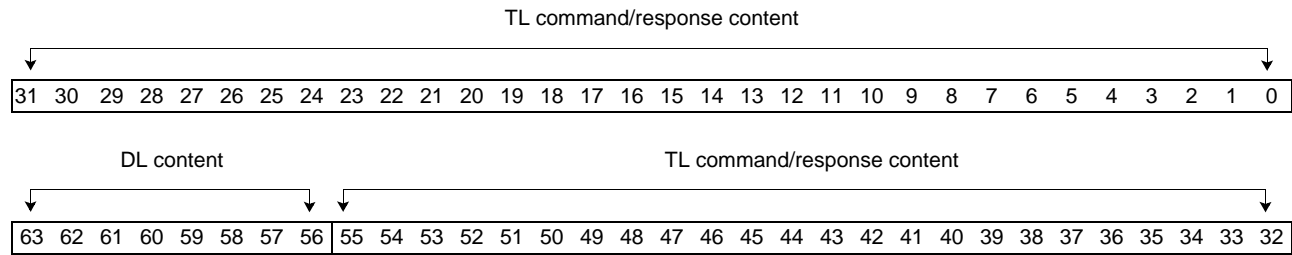
← Bytes(63:0) →	
DL content	TL command/response content
Data flit 0	
Data flit 1	
Data flit 2	
Data flit 3	
Data flit 4	
Data flit 5	
Data flit 6	
Data flit 7	
DL content	TL command/response content
Data flit 0	
Data flit 1	
DL content	TL command/response content
DL content	TL command/response content

Table 5-1 uses color to illustrate the coverage of the CRC found in the DL content of the control flit. The CRC covers the control flit that it is contained in and all *previous*, if any, data flits. The DL content found in the control flit also contains “bad data flit indicators” for the previous data flits. A control flit specified by the TL shall always follow data flits as shown in Table 5-1. The DL may insert DL-idle-control flits when the TL interface to the DL is idle.

The transmit order in Table 5-1 is from right to left and top to bottom. That is, data flits are transmitted in increasing address order. Control flits follow this convention.

Approved

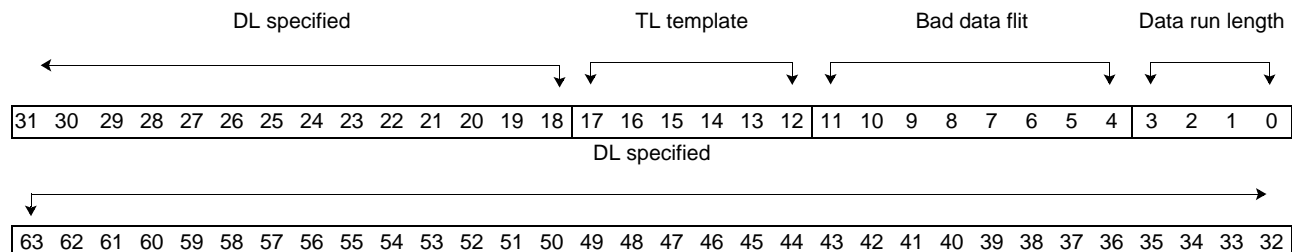
5.1 DL frame control flit (64 bytes)



Bytes	Field name	Description
55:0	TL command/response content	This field contains information provided by the TL. This 448-bit field (447:0) is comprised of sixteen 28-bit slots. One or more slots comprise either a null entry, TL command packet, or TL response packet. See <i>Section 5.1.2 TL command/response content</i> on page 97 for slot layout information.
63:56	DL content	This field contains information added by both the TL and the DL layer. See <i>Section 5.1.1 DL content</i> for the specification of this field.

5.1.1 DL content

This field is bytes 63:56 of the DL frame control flit.

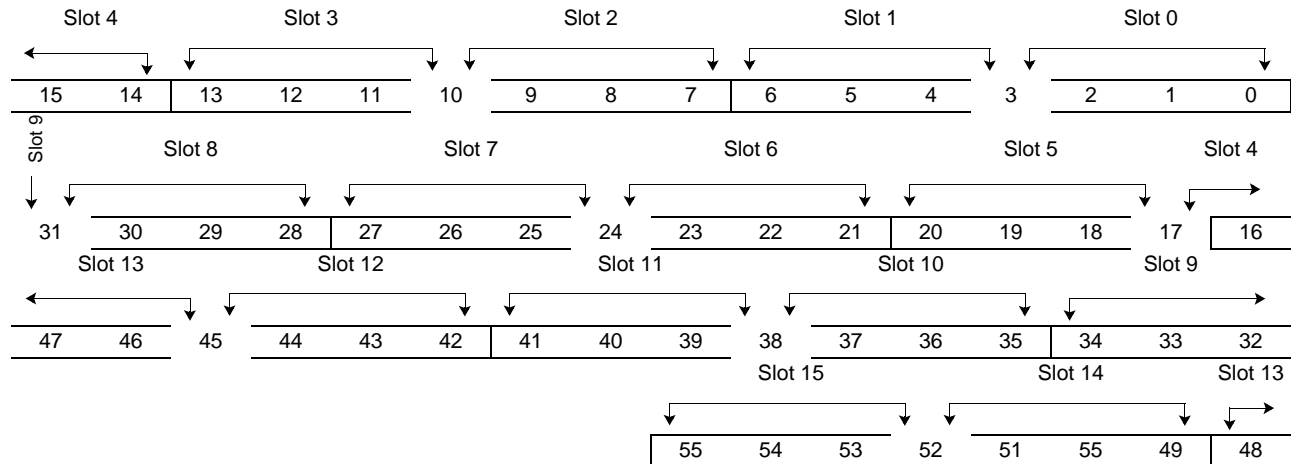


Bit	Field name	Description
3:0	Data run length	This 4-bit field indicates the number of data flits until the next control flit. A value of 0 indicates that the next flit is a control flit. Valid values are {0...8}.
11:4	Bad data flit indication	This 8-bit field indicates that data flits received prior to this control flit contain bad data and shall not be used without being marked as bad (for example, mark as SUE). Each bit corresponds to one data flit (for example, bit 0 corresponds to data flit 0, which is the first data flit following the previous control flit). See <i>Table 5-1</i> on page 95 to match up the data flit being reported to the bit in this field. 11 Data flit 7 is in error. 10 Data flit 6 is in error. 9 Data flit 5 is in error. 8 Data flit 4 is in error. 7 Data flit 3 is in error. 6 Data flit 2 is in error. 5 Data flit 1 is in error. 4 Data flit 0 is in error.
17:12	TL template	This 6-bit field specifies the locations of opcodes found in the 448-bit TL command and response content field. See <i>Section 6 TL and TLX template specifications</i> on page 99 for the specification of this field.
63:18	DL specified	The specification for this 46-bit field is found in the OpenCAPI DL specification . This field is expected to contain the CRC, ACK, and ACK count information.

Approved

5.1.2 TL command/response content

Slots are packed into the DL control flit as shown in the following figure. Each slot occupies 3.5 bytes (28 bits). The slot number and control flit byte are shown. The number of slots used by a command or response can be found in the specification of the command or response.



5.1.3 Data transport, order, and alignment

Data is transported between the TL and TLX using *data carriers*, which are specified as 64-byte data flits. To transport data, one or more DCP credits shall be obtained atomically when obtaining the VC credit required to send a command or response. A DCP credit is associated with 64 bytes of data.

One DCP credit is required for each data carrier used to transport the data. For example, sending 128 bytes requires two DCP credits when the data is sent in two 64-byte data flits.

A data carrier shall be associated with a single command or response. Multiple data carriers may be associated with a single command or response.

Within each control flit, there is an order to the commands and responses as shown in *Section 5.1.2 TL command/response content* on page 97. Commands and responses loaded into lower numbered slots are ordered before commands and responses loaded into higher numbered slots. Data shall be loaded into data carriers by the TL in the same order as the commands and responses are specified in the control flits.

For each command and response associated with data, there is a dLength field and a dPart field specified. These are used to pull the data out of the data carriers and associate the data with the command or response. Each data carrier is examined using the length information and the data is associated with a single command or response. (Additional association can be made using the AFU or CAPP tag provided to locate the machine associated with the command or response.) In some cases (for example, **dma_pr_w**), the dLength field is implicit and is specified as 64 bytes (one data flit).

Data shall be address aligned within a 64-byte data carrier. That is, this architecture treats a data carrier as if it were a memory-mapped naturally aligned data block, where each byte of data is loaded into the data carrier based on the address of the byte being loaded. When a command or response with immediate data uses multiple data carriers, the data shall be loaded in increasing address order. That is, offset 0 from the address specified by the command or response, and adjusted by the dPart field, shall be loaded first and the remaining data is loaded in increasing address order.

Approved

When the data is not associated with an address, the data shall be placed starting at byte 0 of the data carrier and increasing byte locations until all the data has been loaded into the data carrier. A command or response with this type of data shall use only a single data carrier. Additional restrictions might be found in the command and response description.

Engineering note

There are naturally occurring cases when all bytes of a data carrier are not fully specified by the command or response associated with the data. For example, the data associated with a **dma_pr_w** does not specify all bytes within a 64-byte data flit. That is, there are bytes within the data carrier that are defined by the write operation based on pLength and starting address, and there are bytes that are undefined by the architecture.

It is strongly recommended that the architecturally undefined data bytes found within a data carrier do not contain information associated with any application other than the application associated with the command or response and is limited to the permissions granted to the application.

The method used to ensure that the contents of undefined data locations within a data carrier are not from a different process is determined by the implementation. Suggested methods include, but are not limited to the following:

- Set all undefined byte locations to zero.
- Set all undefined byte locations to a fixed or random non-zero pattern of bits
- Replicate the content of defined byte locations to undefined byte locations.

The architecture does not provide architectural conformance statements regarding the contents of undefined byte locations within a data carrier other than the conformance statement specified by the definition of the architectural term “*undefined*”.

A 64-byte data flit may be used with any type of data that is specified for:

- Any command

One to four data flits may be used to provide data associated with an address. The number of data flits is dependent on the number of bytes specified by the command or response associated with the immediate data. When multiple data flits are used, the data flits shall be loaded in increasing address order based on the dLength and dPart specified in the command or response.

6. TL and TLX template specifications

This specification defines the allowed placement formats of TL/TLX packets within the DL/DLX frame's control flit. The allowed formats are captured in four capability descriptions defined in *Table 6-1* on page 100.

A TL template field is specified within the DL content of a DL packet's control flit. The DL content of the control flit is shown in *Section 5.1 DL frame control flit (64 bytes)* on page 96. The DLX frame has the same format as the DL frame.

The TL architecture specifies all template capability descriptions and specifies a number for each specification that is used in the TL template field. When transmitting a packet:

- The TL shall place the TL transmit template number used to form the control flit of the DL frame.
- The TLX shall place the TLX transmit template number used to form the control flit of the DLX frame.

The template capabilities specify the legal locations of one or more TL/TLX command or response packets' starting slot⁴ as well as the contiguous number of slots used by the packet. Unused control flit slots are reserved. That is, unused control flit slots shall be set to an all zero state when transmitted and shall not be examined on receipt for any purpose other than CRC checking.

The template restricts the maximum length of TL/TLX command or response packets placed in the control flit. The template specification permits a smaller packet to be placed into a larger specification footprint. For example, a six-slot template specification may be filled with a 4-, 2-, or 1-slot packet. The state of the unused slots is undefined.

Developer Note

Allowing smaller TL packets to occupy larger packet-specified locations helps to reduce the number of template specifications.

A template may further restrict the TL/TLX command or response packet placed into packet locations. These restrictions include and are not limited to the following:

- Command
- Response
- VC used
- DCP requirement (present or absent)

Table 6-1 on page 100 defines the template capabilities.

4. Bits 27:0 of the TL/TLX packet specification.

Approved

Table 6-1. Template capability definitions

Capability	Definition	Specified by
TLX receive template	Specifies the templates that the OpenCAPI device supports when receiving DL frames.	OpenCAPI device
TL transmit template	Specifies the templates that the host supports when transmitting a DL frame to the OpenCAPI device.	Host
TL receive template	Specifies the templates that the host supports when receiving DLX frames.	Host
TLX transmit template	Specifies the templates that the OpenCAPI device supports when transmitting a DLX frame to the host.	OpenCAPI device

The intersection of the TL transmit template capability and the TLX receive template capability shall not be a null set. All OpenCAPI devices shall support TLX receive template x'00'.

The intersection of the TLX transmit template capability and the TL receive template capability shall not be a null set. All OpenCAPI hosts shall support TL receive template x'00'.

See the host's platform architecture for additional information about how the receive and transmit capabilities are resolved and what is stored into the OpenCAPI device's configuration space.

6.1 TLX receive and TL transmit template capability specification

Table 6-2. TLX receive/TL transmit template

Slot #	x'00'	x'01'	x'02'	x'03'
0	return_tlx_credits^a	4-slot TL packet	2-slot TL packet	4-slot TL packet
1			2-slot TL packet	
2	reserved			
3				
4	6-slot TL packet	4-slot TL packet	2-slot TL packet	6-slot TL packet
5			2-slot TL packet	
6				
7		4-slot TL packet	2-slot TL packet	
8				
9				
10	reserved	4-slot TL packet	2-slot TL packet	6-slot TL packet
11			2-slot TL packet	
12		4-slot TL packet		
13				
14			2-slot TL packet	
15				

a. Template x'0' slots 0 and 1 shall contain either a **nop** or **return_tlx_credits**

Approved

6.2 TL receive and TLX transmit template capability specification

Table 6-3. TL receive/TLX transmit template

Slot #	x'00'	x'01'	x'02'	x'03'
0	return_tl_credits ^a	4-slot TLX packet	2-slot TLX packet	4-slot TLX packet
1			2-slot TLX packet	
2	reserved			
3				
4	6-slot TLX packet	4-slot TLX packet	2-slot TLX packet	6-slot TLX packet
5			2-slot TLX packet	
6				
7		4-slot TLX packet	2-slot TLX packet	
8			2-slot TLX packet	
9				
10	reserved	4-slot TLX packet	2-slot TLX packet	6-slot TLX packet
11			2-slot TLX packet	
12		4-slot TLX packet		
13			2-slot TLX packet	
14				
15			2-slot TLX packet	

a. Template x'00' slots 0 and 1 shall contain either a **nop** or **return_tl_credits**.

6.3 Control-flit rate capability

Each receive and transmit template capability specification has an associated control-flit rate capability. When a template is used to format a control flit, the template's associated control flit rate capability specifies when the next control flit may be sent.

The control-flit rate capability controls the DL flit spacing between control flits. The configuration space for this capability provides 4 bits.

- A value of x'0' indicates that a control flit may be sent in the following *flit-cycle*.
- A value of x'1' indicates that the cycle following the control flit shall contain either a null control flit or a data flit.

In general, a value of "n" indicates that there shall be a gap of "n" 64-byte flits before the next control flit can be sent. During the gap, either null control flits or data flits shall be inserted.

A null control flit is defined as using template x'00'. The 6-slot packet contains a 1-slot null command, and the remaining five slots are undefined. A return credit response found in slots 0 and 1 may be used to return credits.

Engineering note

At the start of initialization of an OpenCAPI device, the flit rate capability is unknown since the configuration space has not yet been examined. Template x'00' is used to issue **config_read** commands to determine the device's capabilities.

Until the device's capabilities are determined, template 0 shall be used and the control flit rate capability shall be assumed to be x'F'.

7. Error detection

This section identifies errors and classes of errors detected by the TL and TLX. Error notifications and the collection of error signatures are specified.

- The host or device may provide:
 - Additional error signature information for error events defined by this architecture

Specification of these implementation-specific error signature extensions might be found in either the host's platform architecture, the host's user's guide, or the manufacturer's documentation provided with the OpenCAPI device.

Actions taken by hypervisors, operating systems, firmware, or device drivers when error notifications are asserted are beyond the scope of this architecture.

Error events are specified in *Table 9-1* on page 213 using the following format:

Error event name	Description of error event	
	Error signature:	The minimum set of information captured by a compliant design. Engineering Note The error signatures specified in <i>Table 9-1</i> on page 213 shall be accessible and may be obtained for diagnostic use by an examination of hardware facilities and may require additional host- or device-specific software manipulation.
	Error Class	Specifies the class and architecture conformance requirements.

Error classes are specified as follows:

- *Correctable error events* are error conditions that the hardware can recover without any loss of function, state, or data.
- *Fatal error events* are error conditions that result in the unrecoverable loss of function, state, or data. Continued use of the link and attached device might not be safe. A reset might be required to return the link and device to a safe operational state. The architecture does not place conformance requirements on an implementation once a fatal error event has occurred. That is, continued use of the link may occur and the results of operations after a fatal error are undefined.
- *Non-fatal error events* are error conditions that affect the operation of a single transaction. The link is considered to be operationally safe. The results of the transaction might not be as intended. That is, the results of the operation are undefined. Devices associated with the error might require a reset. Devices not associated with the error are not affected by the error. Continued use of the link may occur and the results of operations after a non-fatal error shall conform to the requirements of the architecture.

Error events are assigned error types and may be assigned an error subtype. These assignments are found in bold text in the description of the error event.

Error events are assigned an error class and a conformance requirement.

- *Required error events* are demanded by the architecture and shall be included in any architecture conformance testing. All error events specified as required shall be included in both the TL and TLX implementation unless otherwise specified.

Approved

- *Optional error events* are not required by the architecture. Careful reading of the description of the error events assigned as optional is strongly recommended since detection may be required due to architecturally required actions to be taken. Conformance testing may indicate the presence or absence of the hardware's capability to detect and report the error event.

7.1 Error events

The following error events are specified

acTag specified in a command is outside the configured specification set	acTag specified in the command points to an invalid entry	
Age out specified for xlate_touch.n	Bad BDF and PASID combination	Bad data flit indication error
Bad data received	Bad opcode and template combination	Bad response received
Bad template x'00' format		Control flit overrun
Illegal return credit command location	log2_page_size specification in xlate_touch is bad.	
PA specified is out of bounds	Reserved field not transmitted as 0	Reserved field value used
Reserved opcode used	Returned credit overflows credit counter	Unexpected data carrier
Unsupported template format		

Table 7-1. Error event specification (Page 1 of 6)

Error event	Description				
acTag specified in a command is outside the configured specification set	On receipt of a TLX command packet containing an acTag field, it is determined that the acTag is specified outside the configured specification set. <ul style="list-style-type: none"> • The operation is aborted without changes to the machine state, and a malformed packet error type 2 event is asserted. 				
	<table border="0" style="width: 100%;"> <tr> <td style="width: 30%;">Error signature:</td> <td>opcode(7:0), AFUtag(15:0)</td> </tr> <tr> <td></td> <td style="border: 1px solid black; padding: 5px;"> <p style="text-align: center;">Engineering note</p> <p>When the TLX command is assign_actag, the AFUtag in the error signature is reserved.</p> </td> </tr> </table>	Error signature:	opcode(7:0), AFUtag(15:0)		<p style="text-align: center;">Engineering note</p> <p>When the TLX command is assign_actag, the AFUtag in the error signature is reserved.</p>
	Error signature:	opcode(7:0), AFUtag(15:0)			
	<p style="text-align: center;">Engineering note</p> <p>When the TLX command is assign_actag, the AFUtag in the error signature is reserved.</p>				
Error Class:	Fatal/Required (TL only)				
acTag specified in the command points to an invalid entry	On receipt of a TLX command packet containing an acTag field, the entry in the acTag table is examined and found to be marked invalid. The operation is aborted without changes to the machine state, and an address context error type 0 event is asserted.				
	Error signature:	opcode(7:0), AFUtag(15:0), acTag(11:0)			
	Error Class:	Fatal/Required (TL only)			
Age out specified for xlate_touch.n	An xlate_touch.n is specified with a command flag of "age out". This is a nonsensical combination. The dot-n directive is ignored, and the operation proceeds to completion. <ul style="list-style-type: none"> • An xlate_touch error type 1 event may be asserted. 				
	Error signature:	AFUtag(15:0)			
	Error Class:	Non-fatal/Optional (TL only)			

Approved

Table 7-1. Error event specification (Page 2 of 6)

Error event	Description																								
Bad BDF and PASID combination	<p>On receipt of a TLX command packet containing an acTag field, the entry in the acTag table is found to be valid. The BDF and PASID specified are not valid for use. The operation is aborted without changes to the machine state, and an address context error type 1 event is asserted.</p> <table border="1" data-bbox="613 457 1456 541"> <tr> <td>Error signature:</td> <td>opcode(7:0), AFUtag(15:0), acTag(11:0)</td> </tr> <tr> <td>Error Class:</td> <td>Fatal/Required (TL only)</td> </tr> </table>	Error signature:	opcode(7:0), AFUtag(15:0), acTag(11:0)	Error Class:	Fatal/Required (TL only)																				
Error signature:	opcode(7:0), AFUtag(15:0), acTag(11:0)																								
Error Class:	Fatal/Required (TL only)																								
Bad data flit indication error	<p>On the receipt of a control flit, the bad data flit field indicates that a data flit is bad and is located beyond the scope of the control flit. That is, the control flit indicates that n data flits follow, and the bad data flit field indicates that data flit n or above is in error. In the following valid combinations, an 'x' indicates that the field may take on either a 0 or 1 state.</p> <table border="1" data-bbox="659 674 1013 1100"> <thead> <tr> <th>Data run length</th> <th>Bad data flit</th> </tr> </thead> <tbody> <tr><td>x'0'</td><td>'0000 0000'</td></tr> <tr><td>x'1'</td><td>'0000 000x'</td></tr> <tr><td>x'2'</td><td>'0000 00xx'</td></tr> <tr><td>x'3'</td><td>'0000 0xxx'</td></tr> <tr><td>x'4'</td><td>'0000 xxxx'</td></tr> <tr><td>x'5'</td><td>'000x xxxx'</td></tr> <tr><td>x'6'</td><td>'00xx xxxx'</td></tr> <tr><td>x'7'</td><td>'0xxx xxxx'</td></tr> <tr><td>x'8'</td><td>'xxxx xxxx'</td></tr> </tbody> </table> <ul style="list-style-type: none"> A malformed control flit error type 3 event is asserted. <p>The bad data flit information beyond the scope of the control flit is ignored.</p> <table border="1" data-bbox="613 1178 1456 1310"> <tr> <td>Error signature:</td> <td>The bad data flit and data run length fields are captured. These are found in the DL content of the control flit found in bits 11:0 as shown in <i>Section 5.1.1 DL content</i> on page 96.</td> </tr> <tr> <td>Error Class:</td> <td>Non-fatal/Optional (TL and TLX)</td> </tr> </table>	Data run length	Bad data flit	x'0'	'0000 0000'	x'1'	'0000 000x'	x'2'	'0000 00xx'	x'3'	'0000 0xxx'	x'4'	'0000 xxxx'	x'5'	'000x xxxx'	x'6'	'00xx xxxx'	x'7'	'0xxx xxxx'	x'8'	'xxxx xxxx'	Error signature:	The bad data flit and data run length fields are captured. These are found in the DL content of the control flit found in bits 11:0 as shown in <i>Section 5.1.1 DL content</i> on page 96.	Error Class:	Non-fatal/Optional (TL and TLX)
Data run length	Bad data flit																								
x'0'	'0000 0000'																								
x'1'	'0000 000x'																								
x'2'	'0000 00xx'																								
x'3'	'0000 0xxx'																								
x'4'	'0000 xxxx'																								
x'5'	'000x xxxx'																								
x'6'	'00xx xxxx'																								
x'7'	'0xxx xxxx'																								
x'8'	'xxxx xxxx'																								
Error signature:	The bad data flit and data run length fields are captured. These are found in the DL content of the control flit found in bits 11:0 as shown in <i>Section 5.1.1 DL content</i> on page 96.																								
Error Class:	Non-fatal/Optional (TL and TLX)																								
Bad data received	<p>On processing data flits, a data flit is marked bad by the bad data flit indication field found in the control flit as described in <i>Section 5.1.1 DL content</i> on page 164.</p> <ul style="list-style-type: none"> A bad data flit error event may be asserted. The command or response associated with the bad data is provided with the data and the bad data indication. The bad data indication shall be propagated to the final destination of the data. It is strongly recommended that the error propagation be implemented in a fashion that allows for error isolation; that is, first error incidence reporting. <table border="1" data-bbox="613 1503 1456 1703"> <tr> <td>Error signature:</td> <td>Data is associated with: TLX response packet: opcode(7:0), CAPPTag(15:0). TLX command packet: opcode(7:0), acTag(11:0), AFUtag(15:0). TL response packet: opcode(7:0), AFUtag(15:0). TL command packet: opcode(7:0), CAPPTag(15:0)</td> </tr> <tr> <td>Error Class:</td> <td>Non-fatal/Optional (TL and TLX)</td> </tr> </table>	Error signature:	Data is associated with: TLX response packet: opcode(7:0), CAPPTag(15:0). TLX command packet: opcode(7:0), acTag(11:0), AFUtag(15:0). TL response packet: opcode(7:0), AFUtag(15:0). TL command packet: opcode(7:0), CAPPTag(15:0)	Error Class:	Non-fatal/Optional (TL and TLX)																				
Error signature:	Data is associated with: TLX response packet: opcode(7:0), CAPPTag(15:0). TLX command packet: opcode(7:0), acTag(11:0), AFUtag(15:0). TL response packet: opcode(7:0), AFUtag(15:0). TL command packet: opcode(7:0), CAPPTag(15:0)																								
Error Class:	Non-fatal/Optional (TL and TLX)																								

Approved

Table 7-1. Error event specification (Page 3 of 6)

Error event	Description				
Bad opcode and template combination	<p>The format of the control flit specified by the template is in error. Opcodes specified indicate packet sizes larger than allowed by the template found in the control flit. This error is detected in both the TL and TLX.</p> <p>All commands or responses identified in the control flit are aborted and do not cause any machine state changes, and a malformed control flit error type 2 event is asserted.</p> <table border="1" data-bbox="613 478 1458 674"> <tr> <td data-bbox="613 478 802 632">Error signature:</td> <td data-bbox="802 478 1458 632"> <ul style="list-style-type: none"> • Template (5:0) found in the control flit. • The opcode (7:0) found in the control flit where it was determined that the template packet size rules were violated. • The slot location, a 4-bit field, where it was determined that the template packet size rules were violated. </td> </tr> <tr> <td data-bbox="613 632 802 674">Error Class:</td> <td data-bbox="802 632 1458 674">Fatal/Required. (TL and TLX)</td> </tr> </table>	Error signature:	<ul style="list-style-type: none"> • Template (5:0) found in the control flit. • The opcode (7:0) found in the control flit where it was determined that the template packet size rules were violated. • The slot location, a 4-bit field, where it was determined that the template packet size rules were violated. 	Error Class:	Fatal/Required. (TL and TLX)
Error signature:	<ul style="list-style-type: none"> • Template (5:0) found in the control flit. • The opcode (7:0) found in the control flit where it was determined that the template packet size rules were violated. • The slot location, a 4-bit field, where it was determined that the template packet size rules were violated. 				
Error Class:	Fatal/Required. (TL and TLX)				
Bad response received	<p>TL: On receipt of a TLX response packet, the CAPPTag is first examined to determine if a prior command has been issued using the CAPPTag found in the response packet. It is reported as a bad response received variant 0 if the CAPPTag has not been used.</p> <p>If the response packet is not a variant 0, the response opcode is checked to determine if it is a valid response for the command opcode used. It is reported as a bad response received variant 1 if it is not.</p> <hr/> <p>TLX: On receipt of a TL response packet, the AFUTag is first examined to determine if a prior command has been issued using the AFUTag found in the response packet. It is reported as a bad response receive variant 0 if the AFUTag has not been used.</p> <p>If the response packet is not a variant 0, the response opcode is checked to determine if it is a valid response for the command opcode used. It is reported as a bad response received variant 1 if it is not.</p> <hr/> <ul style="list-style-type: none"> • The actions specified by the completion of the command due to a correct response are aborted, and a malformed packet error type 5 event is asserted. The state machine representing the command source is left in an undefined state. The undefined state of the machine is bounded, that is, the state is known to the implementation and the actions taken by the state machine in this architecturally undefined state is predictable by the implementation. <table border="1" data-bbox="613 1207 1458 1480"> <tr> <td data-bbox="613 1207 802 1444">Error signature:</td> <td data-bbox="802 1207 1458 1444"> <p>TL: CAPPTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>TLX: AFUTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>In the error signature, the variant(0) field is set to the variant error type.</p> <ul style="list-style-type: none"> • variant(0) = '0' when the error is variant 0. • variant(0) = '1' when the error is variant 1. </td> </tr> <tr> <td data-bbox="613 1444 802 1480">Error Class:</td> <td data-bbox="802 1444 1458 1480">Fatal/Required (TL and TLX)</td> </tr> </table>	Error signature:	<p>TL: CAPPTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>TLX: AFUTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>In the error signature, the variant(0) field is set to the variant error type.</p> <ul style="list-style-type: none"> • variant(0) = '0' when the error is variant 0. • variant(0) = '1' when the error is variant 1. 	Error Class:	Fatal/Required (TL and TLX)
Error signature:	<p>TL: CAPPTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>TLX: AFUTag(15:0), variant(0). For a variant 1, the command opcode(7:0) is also provided.</p> <p>In the error signature, the variant(0) field is set to the variant error type.</p> <ul style="list-style-type: none"> • variant(0) = '0' when the error is variant 0. • variant(0) = '1' when the error is variant 1. 				
Error Class:	Fatal/Required (TL and TLX)				
Bad template x'00' format	<p>The format of the control flit, specified as using the x'00' template, does not match the template x'00' format. Slot 0 does not contain either a nop, or return_tl_x_credits (detected by the TLX), or return_tl_credits (detected by the TL), opcode.</p> <ul style="list-style-type: none"> • Any commands or responses found in the control flit are aborted and do not cause any machine state changes. A malformed control flit error type 0 event is asserted. <table border="1" data-bbox="613 1648 1458 1728"> <tr> <td data-bbox="613 1648 802 1690">Error signature:</td> <td data-bbox="802 1648 1458 1690">Slot 0 (27:0) contents</td> </tr> <tr> <td data-bbox="613 1690 802 1728">Error Class:</td> <td data-bbox="802 1690 1458 1728">Fatal/Required (TL and TLX)</td> </tr> </table>	Error signature:	Slot 0 (27:0) contents	Error Class:	Fatal/Required (TL and TLX)
Error signature:	Slot 0 (27:0) contents				
Error Class:	Fatal/Required (TL and TLX)				

Approved

Table 7-1. Error event specification (Page 4 of 6)

Error event	Description
Control flit overrun	The destination is unable to accept a subsequent control flit. A possible cause is a violation of the control flit rate capability for the prior control flit's template. <ul style="list-style-type: none"> The incoming control flit is discarded. The machine state is unchanged. A control flit overrun error event is asserted.
	Error signature: None
	Error Class: Fatal/Required (TL and TLX)
Illegal return credit command location	Regardless of the template used, return_tl_x_credits and return_tl_credits shall be found only in slots 1:0. <ul style="list-style-type: none"> The credit return, as specified by the command, may occur. A malformed packet error type 3 event shall be asserted.
	Error signature: Template (5:0) found in the control flit; slot where return credit opcode was found (3:0).
	Error Class: Fatal/Required (TL and TLX)
log ₂ _page_size specification in xlate_touch is bad.	This error is detected when an xlate_touch is specified with a cmd_flag specification of age-out, and the page size specified by the ATC entry found does not match the page size specified by the command's log ₂ _page_size field. <ul style="list-style-type: none"> See <i>Figure 2-1 Address translation sequence: xlate_touch</i> on page 111 for actions taken when there is a mismatch. An xlate_touch error type 0 event may be asserted.
	Error signature: acTag(11:0)
	Error Class: Non-fatal/Optional (TL only)
PA specified is out of bounds	A command specifies a PA that is determined to be out of bounds for the AFU _M . For TL commands, the host has specified a PA that is outside the AFU _{M1} PA range. <ul style="list-style-type: none"> The operation is aborted without changes to the machine state. A PA specification error event is asserted.
	Error signature: opcode(7:0), PA(63:0)
	Error Class: Fatal/Required (TLX only)
Reserved field not transmitted as 0	On the receipt of a command or response packet, it is determined that a field specified by the architecture as reserved does not contain 0. <ul style="list-style-type: none"> The packet is used. That is, the operation specified by the command or response occurs normally. This is an architecture conformance violation. A malformed packet error type 4 event is asserted.
	Error signature: None.
	Error Class: Non-fatal/Optional (TL and TLX)

Approved

Table 7-1. Error event specification (Page 5 of 6)

Error event	Description																																								
Reserved field value used	<p>On the receipt of a command or response packet, it is determined that a field specification contains an architecturally reserved value. If multiple fields contain reserved values, only one field is reported.</p> <ul style="list-style-type: none"> The operation is aborted and the machine state is unchanged. A malformed packet error type 1 event is asserted. <p>The following fields have reserved values. Detection occurs at the receiver of the command or response packet.:</p> <ul style="list-style-type: none"> <i>cmd_flag</i>. Detected by TL and TLX. <i>dLength</i>. Detected by TL and TLX. <i>dPart</i> detected by TL and TLX. <i>pLength</i>. Detected by TL and TLX. <i>Resp_code</i>. Detected by TL and TLX. 																																								
	<p>Error signature: opcode(7:0), starting (LSb) field offset within the packet. For example, the acTag in a rd_wnitc TLX command packet has a field offset of 24.</p>																																								
	<p>Error Class: Fatal/Required (TL and TLX)</p>																																								
Reserved opcode used	<p>On the receipt of a command or response packet, the opcode field is examined and found to be a value reserved by the architecture.</p> <ul style="list-style-type: none"> The packet is dropped, the machine state is unchanged. A malformed packet error type 0 event is asserted. 																																								
	<p>Error signature: opcode(7:0)</p>																																								
	<p>Error Class: Fatal/Required (TL and TLX)</p>																																								
Returned credit overflows credit counter	<p>On processing of a return_tl_credits or return_tl_x_credits response packet, it is determined that the addition of the credits specified by the response will cause the counter to overflow.</p> <ul style="list-style-type: none"> The counter may increment and is allowed to saturate. That is, the counter shall not wrap. A credit return error event is asserted. 																																								
	<p>Error signature: opcode(7:0), specification of the counter or counters associated with the error using the following format:</p> <table style="margin-left: auto; margin-right: auto; border-collapse: collapse;"> <tr> <td style="border: 1px solid black; padding: 2px;">TL:</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">dcp.1</td> <td style="border: 1px solid black; padding: 2px;">dcp.0</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">vc.1</td> <td style="border: 1px solid black; padding: 2px;">vc.0</td> </tr> <tr> <td style="border: 1px solid black; padding: 2px;">TLX:</td> <td style="border: 1px solid black; padding: 2px;">dcp.3</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">dcp.0</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">vc.3</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">R</td> <td style="border: 1px solid black; padding: 2px;">vc.0</td> </tr> <tr> <td></td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> <td style="text-align: center;">↓</td> </tr> <tr> <td></td> <td style="border: 1px solid black; padding: 2px;">8</td> <td style="border: 1px solid black; padding: 2px;">7</td> <td style="border: 1px solid black; padding: 2px;">6</td> <td style="border: 1px solid black; padding: 2px;">5</td> <td style="border: 1px solid black; padding: 2px;">4</td> <td style="border: 1px solid black; padding: 2px;">3</td> <td style="border: 1px solid black; padding: 2px;">2</td> <td style="border: 1px solid black; padding: 2px;">1</td> <td style="border: 1px solid black; padding: 2px;">0</td> </tr> </table>	TL:	R	R	dcp.1	dcp.0	R	R	R	vc.1	vc.0	TLX:	dcp.3	R	R	dcp.0	R	vc.3	R	R	vc.0		↓	↓	↓	↓	↓	↓	↓	↓	↓		8	7	6	5	4	3	2	1	0
	TL:	R	R	dcp.1	dcp.0	R	R	R	vc.1	vc.0																															
TLX:	dcp.3	R	R	dcp.0	R	vc.3	R	R	vc.0																																
	↓	↓	↓	↓	↓	↓	↓	↓	↓																																
	8	7	6	5	4	3	2	1	0																																
<p>Error Class: Fatal/Required (TL and TLX)</p>																																									

Approved

Table 7-1. Error event specification (Page 6 of 6)

Error event	Description																								
Unexpected data carrier	<p>The destination has detected that the accumulated number of data carriers has exceeded the amount of data expected.</p> <p>The expected data count is determined by an examination of the commands and responses received that specify immediate data. The ordering requirements specifying commands or responses with immediate data are received before the data is found in <i>Section 5.1.3 Data transport, order, and alignment</i> on page 97.</p> <hr/> <p>TL: TLX packets received by the TL with immediate data specify the amount of data expected from the TLX. TLX packets that specify immediate data:</p> <table border="0" data-bbox="625 661 1356 772"> <tr> <td>mem_rd_response</td> <td>dma_w</td> <td>dma_w.n</td> <td>dma_w.be</td> </tr> <tr> <td>dma_w.be.n</td> <td>dma_pr_w</td> <td>dma_pr_w.n</td> <td>amo_rw</td> </tr> <tr> <td>amo_rw.n</td> <td>amo_w</td> <td>amo_w.n</td> <td>intrp_req.d</td> </tr> </table> <hr/> <p>TLX: TL packets received by the TLX with immediate data specify the amount of data expected by the TL. TL packets that specify immediate data:</p> <table border="0" data-bbox="625 982 1356 1050"> <tr> <td>read_response</td> <td>write_mem</td> <td>write_mem.be</td> <td>pr_wr_mem</td> </tr> <tr> <td colspan="4">config_write</td> </tr> </table> <ul style="list-style-type: none"> • The unexpected data carrier's contents may be discarded. An unexpected data carrier error event is asserted. <div style="border: 1px solid black; padding: 5px; margin-top: 10px;"> <p align="center">Developer Note</p> <p>Expected data can be counted by determining the number of DCP credits required to send the data as specified by the command or response. Counting DCP credits can be split out by data credit pool number or aggregated. Counts are incremented as expected-data-information, as described above, is observed. Counts are decremented as the data is received. That is, the DCP count of expected data is decremented as the data carriers arrive. When the count drops below zero, an unexpected data carrier error is detected. Other implementation specific methods may also be used.</p> </div> <table border="1" data-bbox="617 1432 1445 1514"> <tr> <td>Error signature:</td> <td>None</td> </tr> <tr> <td>Error Class:</td> <td>Fatal/Required (TL and TLX)</td> </tr> </table>	mem_rd_response	dma_w	dma_w.n	dma_w.be	dma_w.be.n	dma_pr_w	dma_pr_w.n	amo_rw	amo_rw.n	amo_w	amo_w.n	intrp_req.d	read_response	write_mem	write_mem.be	pr_wr_mem	config_write				Error signature:	None	Error Class:	Fatal/Required (TL and TLX)
mem_rd_response	dma_w	dma_w.n	dma_w.be																						
dma_w.be.n	dma_pr_w	dma_pr_w.n	amo_rw																						
amo_rw.n	amo_w	amo_w.n	intrp_req.d																						
read_response	write_mem	write_mem.be	pr_wr_mem																						
config_write																									
Error signature:	None																								
Error Class:	Fatal/Required (TL and TLX)																								
Unsupported template format	<p>An unsupported template is specified in a received control flit.</p> <ul style="list-style-type: none"> • Any commands or responses identified in the control flit are aborted. Changes to the machine state are undefined. A malformed control flit error type 1 event is asserted. <table border="1" data-bbox="617 1633 1445 1713"> <tr> <td>Error signature:</td> <td>Template (5:0) found in the control flit.</td> </tr> <tr> <td>Error Class:</td> <td>Fatal/Required (TL and TLX)</td> </tr> </table>	Error signature:	Template (5:0) found in the control flit.	Error Class:	Fatal/Required (TL and TLX)																				
Error signature:	Template (5:0) found in the control flit.																								
Error Class:	Fatal/Required (TL and TLX)																								

Approved

8. OpenCAPI profiles

A device shall use an OpenCAPI profile to specify groups of commands, responses, and templates supported by the device. Each row in the following tables identifies an architectural feature, and each column indicates the content of a profile. The full specification of a profile is comprised of the same column from all tables in this section.

Within each profile a feature is marked using the notation found in *Table 8-1*. The specification of the compliance notation is taken from the view of the consumer of the command or response packet. That is, a command that is specified as mandatory requires that the consumer of the command shall process and execute the command per the architecture specification. Since the architecture specifies any responses or errors that are reported, those features become mandatory as well.

Table 8-1. Feature compliance requirement notation

Support requirement notation	Description
	(blank / empty) No conformance requirement, no recommendation guidance provided.
M	Mandatory
M.cx	Mandatory when the AFU is C1. Otherwise it is unsupported (U).
M.ir	Mandatory when the AFU issues any form of intrp_req . Otherwise it is unsupported (U).
M.mx	Mandatory when the AFU is M1. Otherwise it is unsupported (U).
M.wht	Mandatory when the TLX issues wake_host_thread . Otherwise it is unsupported (U).
M.xt	Mandatory when xlite_touch is issued by the TLX. Otherwise it is unsupported (U).
O	Optional. This feature may be supported. Conformance evaluation to determine the presence of the feature and when present shall test its architectural compliance.
O.E	Compliance specification for endianness support. Either big or little endian data formats used in atomic* class commands may be supported. An implementation shall support one of the formats.
U	Unsupported. This feature is not included in conformance evaluation. Use of a command or response noted as unsupported may result in a fatal error event.

Approved

Table 8-2 specifies the compliance requirements for the TLX and AFU accepting and processing a command from the TL and host.

Table 8-2. Profile specifications for TL commands

TL command	Device interface class, OpenCAPI 3.0
config_read	M
config_write	M
intrap_rdy	M.ir
nop	M
pr_rd_mem	M,mx
pr_wr_mem	M.mx
rd_mem	M.mx
write_mem	M.mx
write_mem.be	M.mx
xlate_done	M.cx

Table 8-3 specifies the compliance requirements for the TL and host accepting and processing a command from the TLX and AFU.

Table 8-3. Profile specifications for TLX commands (Page 1 of 2)

TLX command	Device interface class, OpenCAPI 3.0
amo_rd	M
amo_rd.n	M
amo_rw	M
amo_rw.n	M
amo_w	M
amo_w.n	M
assign_actag	M
dma_pr_w	M
dma_pr_w.n	M
dma_w	M
dma_w.be	M
dma_w.be.n	M
dma_w.n	M
intrap_req	M
intrap_req.d	M
nop	M
pr_rd_wnitc	M
pr_rd_wnitc.n	M
rd_wnitc	M

Approved

Table 8-3. Profile specifications for TLX commands (Page 2 of 2)

TLX command	Device interface class, OpenCAPI 3.0
rd_wntc.n	M
wake_host_thread	M
xlate_touch	M
xlate_touch.n	M

Table 8-4 specifies the compliance requirements for the TLX and AFU accepting and processing a response from the TL and host. Note that in most cases, TL responses are due to TLX commands issued to the host.

Table 8-4. Profile specifications for TL responses

TL response	Device interface class, OpenCAPI 3.0
intrap_resp	M.ir
nop	M
read_failed	M.cx
read_response	M.cx
return_tlx_credits	M
touch_resp	M.xt
wake_host_resp	M.wht
write_failed	M.cx
write_response	M.cx

Table 8-5 specifies the compliance requirements for the TL and the host accepting and processing a response from the TLX and AFU. Note that in most cases, TLX responses are due to TL commands issued to the TLX.

Table 8-5. Profile specifications for TLX responses

TLX response	Device interface class, OpenCAPI 3.0
mem_rd_fail	M
mem_rd_response	M
mem_wr_fail	M
mem_wr_response	M
nop	M
return_tl_credits	M

Table 8-6 and *Table 8-7* add template capability specifications to profiles. As discussed in *Section 6 TL and TLX template specifications* on page 99, the host's platform architecture provides additional information about how receive and transmit capabilities are resolved between the host and the attached OpenCAPI device. Other than the mandatory support for transmitting and receiving template x'00' template control flits, the profile specifications for templates provides guidance as to the recommended templates an implementation should support and is not a conformance requirement. The templates marked as Optional are recommended.

Approved

See *Section 6.1 TLX receive and TL transmit template capability specification* on page 100 and *Section 6.2 TL receive and TLX transmit template capability specification* on page 101. All host and devices shall support template x'00'.

Table 8-6 specifies the requirements and recommendations for the

- TLX to accept a control flit using the specified template.
- TL to transmit a control flit using the specified template.

Table 8-6. Profile specifications for TLX receive/TL transmit templates

TLX receive/TL transmit template	Device interface class, OpenCAPI 3.0
x'00'	M
x'01'	O
x'02'	O
x'03'	O

Table 8-7 specifies the requirements and recommendations for the

- TL to accept a control flit using the specified template.
- TLX to transmit a control flit using the specified template.

Table 8-7. Profile specifications for TL receive/TLX transmit templates

TL receive/TLX transmit template	Device interface class, OpenCAPI 3.0
x'00'	M
x'01'	O
x'02'	O
x'03'	O

Table 8-8 specifies the operation modes recommended to be supported by the host and the OpenCAPI device for different interface classifications and is not a conformance requirement. The operation modes marked as Optional are recommended. The definition of the host operation modes are found in *Section 1.2 Host operation modes* on page 25.

Table 8-8. Profile specifications host operation modes

AFU type	Device interface class, OpenCAPI 3.0
AFU _{C0}	O
AFU _{C1}	O
AFU _{M0}	O
AFU _{M1}	O

Table 8-9 adds page size support specification to profiles. The profile specification for page size provides guidance as to the recommended page sizes supported by the host's ATC and is not a conformance requirement. Address translation and ATC are discussed in *Section 1.6* on page 28. Page sizes marked as Optional are recommended.

Approved

Table 8-9. Profile specifications supported page size

page size	Device interface class, OpenCAPI 3.0
4K	O
64K	O

Table 8-10 specifies the compliance requirements for the TLX and AFU accepting and processing commands and responses from the TL and host with the dLength specified.

Table 8-10. Profile specifications supported dLength by TLX

dLength specification	Device interface class, OpenCAPI 3.0
64	M
128	M
256	O

Table 8-11 specifies the compliance requirements for the TL and host accepting and processing commands and responses from the TLX and AFU with dLength specified.

Table 8-11. Profile specifications supported dLength by TL

dLength specification	Device interface class, OpenCAPI 3.0
64	M
128	M
256	M

Table 8-12 specifies the compliance requirements for the TL and host accepting and processing atomic.* class commands based on the endianness of the data. The *E* field in the command specifies the endianness of the data. See Table 8-3 for the compliance requirements for the TL and host accepting and processing these commands.

Table 8-12. Profile specifications support of endianness data format by the TL

TLX atomic* class command	Device interface class, OpenCAPI 3.0	
	E=0	E=1
amo_rd	O.E	O.E
amo_rd.n	O.E	O.E
amo_rw	O.E	O.E
amo_rw.n	O.E	O.E

Appendix A. AP (TLX) command transaction diagrams

This section contains figures that illustrate AP command flows and TLX and TL interaction.

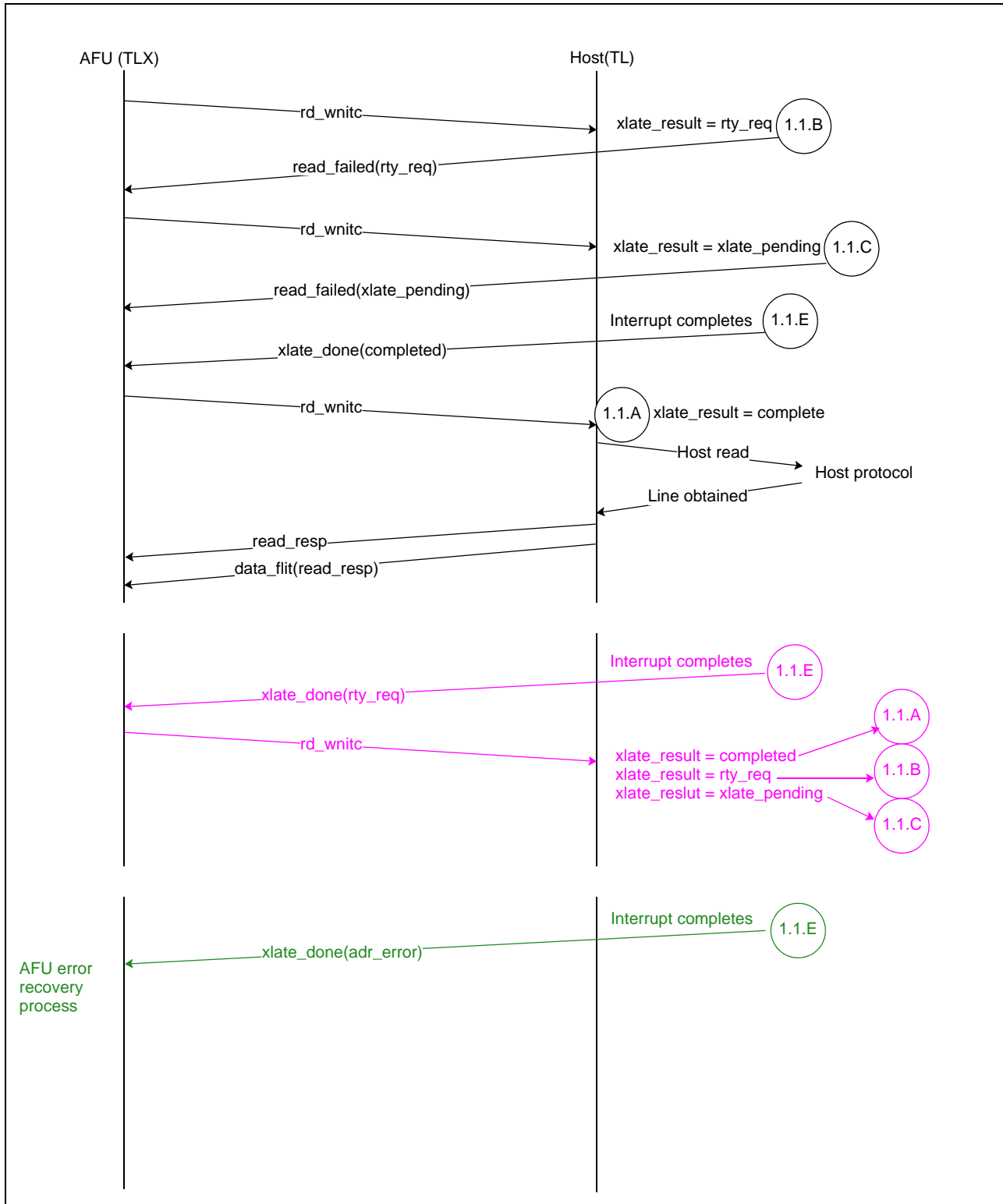
Rules:

1. Commands received at the TL are not serviced until all data, if any, specified by the AP command has arrived.

Approved

A.1 AFU read with no intent to cache; 128 bytes

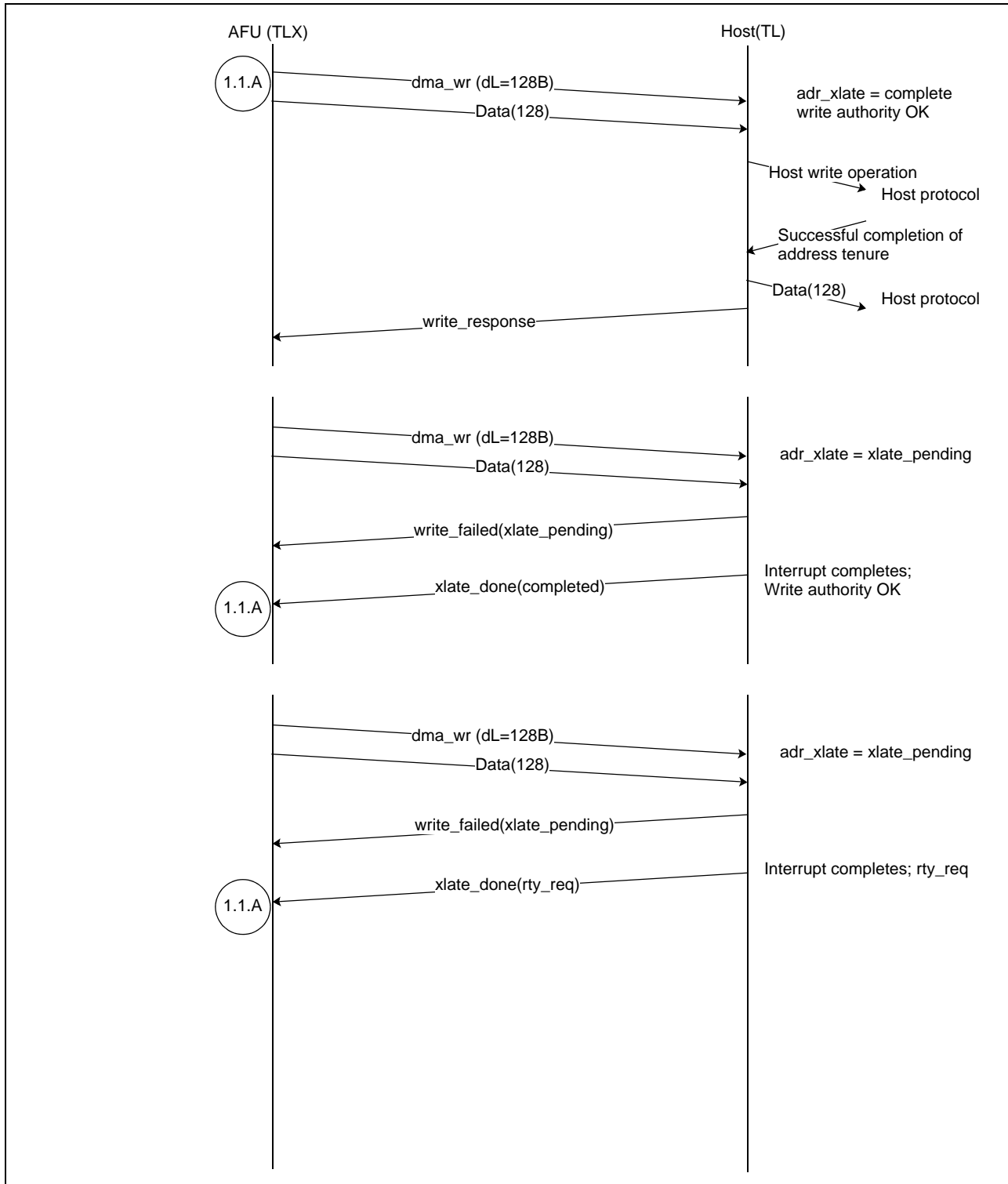
Figure A-1. TLX and TL interaction: *rd_wnitc*



Approved

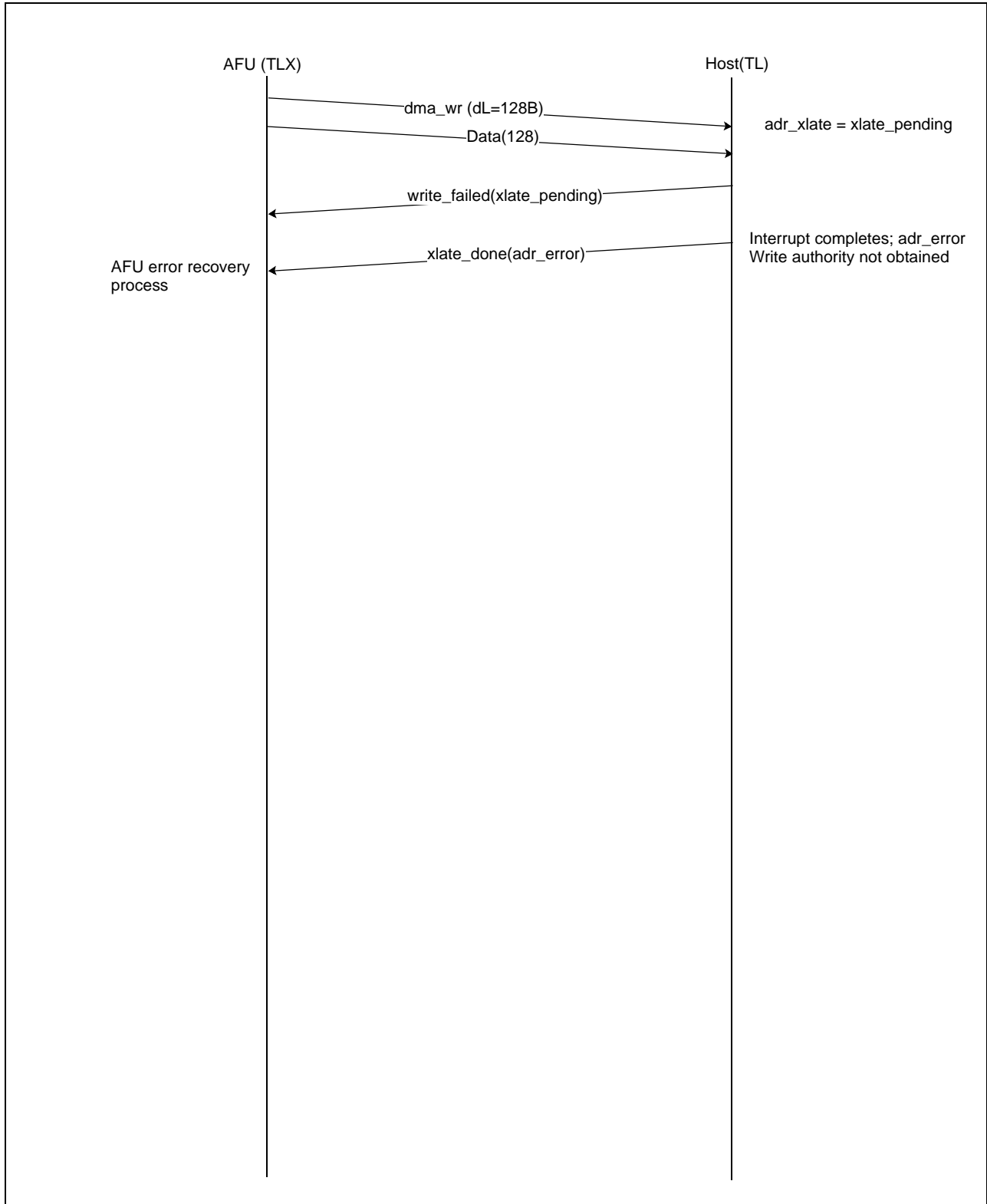
A.2 AFU DMA write; 128 bytes

Figure A-2. TLX and TL interaction: **dma_w** (Page 1 of 2)



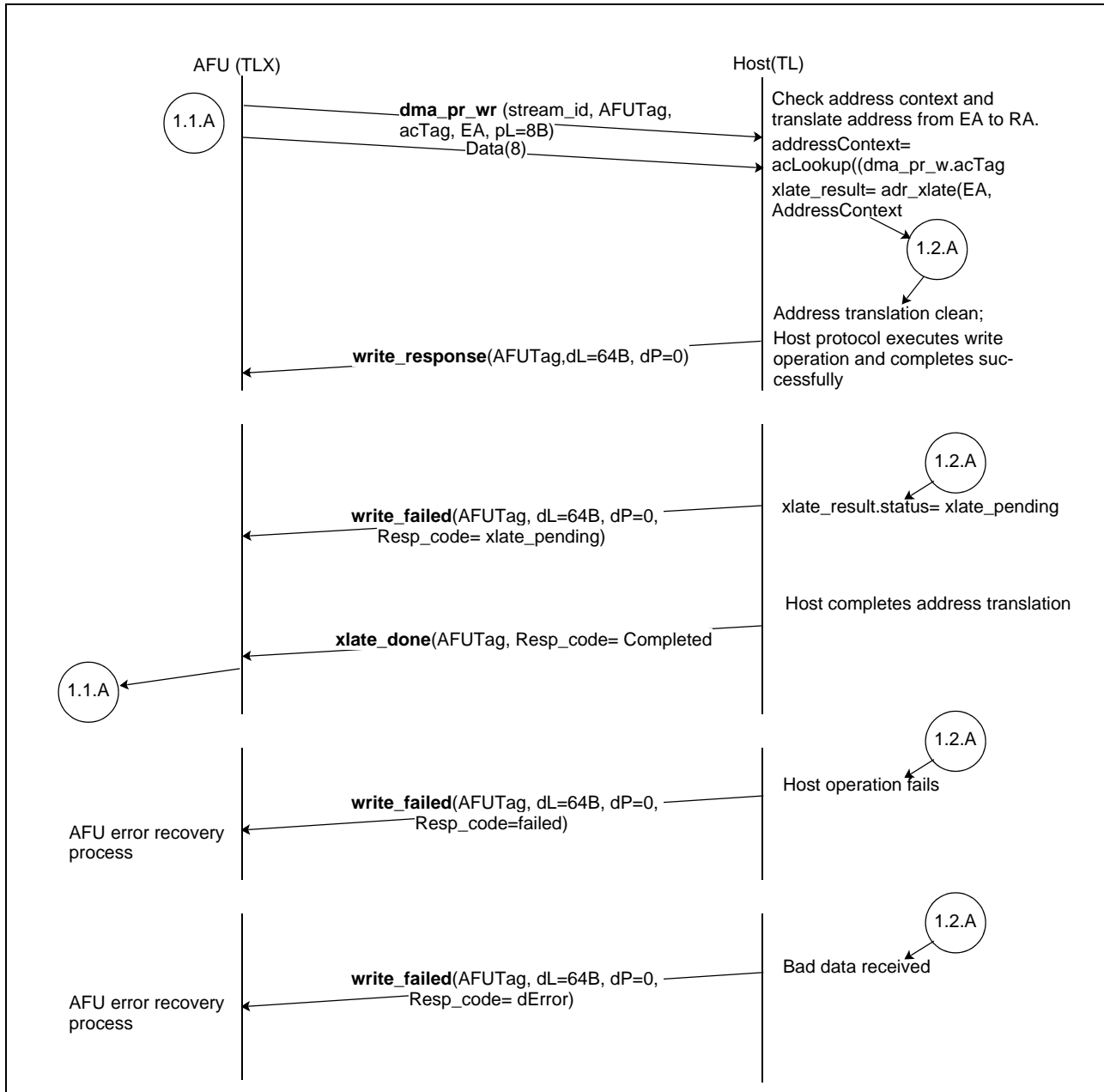
Approved

Figure A-2. TLX and TL interaction: **dma_w** (Page 2 of 2)



A.3 AFU DMA partial write; 8 bytes

Figure A-3. TL and TLX interaction: *dma_pr_w*



Appendix B. CAPP (TL) command transaction diagrams

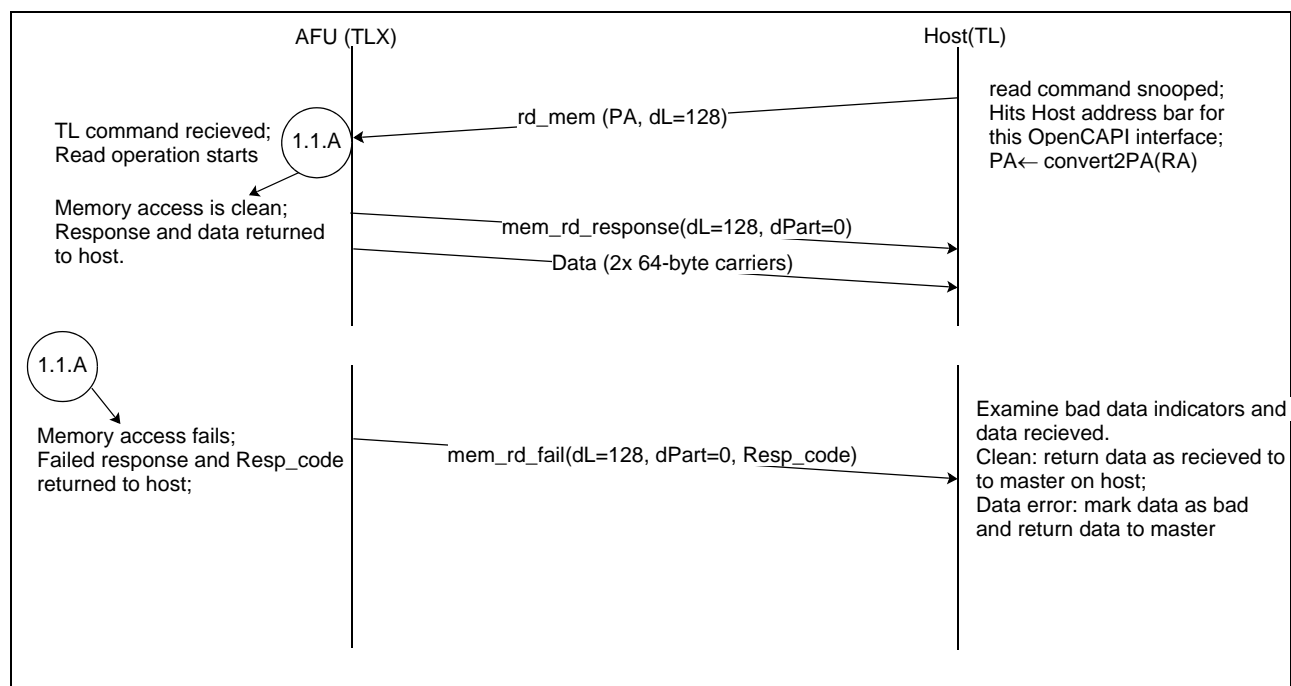
This section contains figures that illustrate CAPP command flows.

Rules:

1. Commands received at the TLX are not serviced until all data, if any, specified by the CAPP command has arrived.

B.1 CAPP memory read; 128 bytes

Figure B-1. TL and TLX transaction: *rd_mem*



Approved

B.2 CAPP memory write; 128 bytes

Figure B-2. TL and TLX transaction: *write_mem*

