Network Woking Group Internet Draft Expire in six months Danny Cohen Myricom Craig Lund Mercury Computers Tony Skjellum Mississippi State University Robert George Mississippi State University Thom McMahon Mississippi State University May 1998

The End-to-End (EEP) PacketWay Protocol for High-Performance Interconnection of Computer Clusters <<u>draft-ietf-pktway-protocol-eep-spec-03.txt</u>>

Status of this Memo

This document is an Internet-Draft. Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

To view the entire list of current Internet-Drafts, please check the "1id-abstracts.txt" listing contained in the Internet-Drafts Shadow Directories on ftp.is.co.za (Africa), ftp.nordu.net (Northern Europe), ftp.nis.garr.it (Southern Europe), munnari.oz.au (Pacific Rim), ftp.ietf.org (US East Coast), or ftp.isi.edu (US West Coast).

Table of Content:

<u>1</u> .	Introduction2
1a.	PktWay and IP2
1b.	General
1c.	The Level-2 Operation of PktWay5
<u>2</u> .	A note about the PktWay documents $\underline{6}$
<u>3</u> .	Notations
<u>4</u> .	PktWay EEP Messages
4a.	The PktWay Message Structure
4b.	The Optional Fields8
<u>5</u> .	Optional Sequence of L2RHs and Symbols9
5a.	L2 Routing Headers (L2RHs) <u>9</u>
5b.	Symbols

<u>6</u> .	EEP Header		. <u>12</u>
6a.	Version		. <u>12</u>
6b.	Priority		. <u>12</u>
6c.	Destination-Type		. <u>13</u>
6d.	Packet Type Extension		. <u>14</u>
6e.	Packet Type		. <u>14</u>
6f.	Endianness		. <u>15</u>
6g.	Padding Length		. <u>15</u>
6h.	Data Length		. <u>15</u>
6i.	Options flag		. <u>15</u>
6j.	Reserved		. <u>15</u>
6k.	Source Address		. <u>16</u>
<u>7</u> .	Optional Header Fields		. <u>16</u>
<u>8</u> .	Optional Data Block		. <u>17</u>
<u>9</u> .	Optional Trailer Fields		. <u>17</u>
<u>10</u> .	EEP Trailer		. <u>18</u>
<u>11</u> .	Appendix-A: Recommendation for PktWay Address Assign	nment	. <u>18</u>
<u>12</u> .	Appendix-B: Glossary		. <u>19</u>
<u>13</u> .	Appendix-C: Acronyms and Abbreviations		. <u>20</u>
<u>14</u> .	Appendix-D: PktWay at a Glance ("cheat-sheet")		. <u>22</u>
<u>15</u> .	Security Considerations		. <u>23</u>
<u>16</u> .	Editor's Address		. <u>23</u>
Cohen et	al	[Page	1]
Internet	-Draft PktWay End-to-End Protocol	October 19	997

1. Introduction

PktWay is an open family of specifications for inter-networking high performance SANs (System Area Networks) and high performance LANs (Local Area Networks) into computing clusters.

Most modern SANs have much in common, such as high data rates, low message latency and low bit error rates. Such SANs are often packet networks made of point-to-point links with flow control and source routing. Yet these SANs do not provide heterogeneous networking support, and are subsequently incapable of direct intercommunications with other SANs. PktWay's goal is to provide high performance "internetting" of such SANs and of high performance LANs.

The core PktWay protocol comprises the End-to-End Protocol (EEP) and the Router-to-Router-Protocol (RRP). This document specifies the EEP (End-to-End) protocol of PktWay. A companion document ("Specification for the Router-to-Router (RRP) PktWay Protocol") specifies the Router-to-Router protocol of PktWay.

Computing clusters and modern MPPs (Massively Parallel Processing systems) are sets of processors interconnected by high performance SANs. Examples are Intel's Paragon and ASCI-red, CRAY's T3D and T3E, and IBM's SP2 and SP3. Most modern SANs have much in common, such as high data rates, low message latency and low bit error rates. Such SANs are often packet networks made of point-to-point links with flow control and source routing.

Unfortunately, there is no efficient way to "internet" these SANs to allow each computing node to have high performance communication directly with any other computing node, in any other interconnected SAN. Hence, there is no way to interconnect such high performance SANs to form as efficient computing cluster as possible.

The objective of PktWay is to provide high performance communication among all the processors in a cluster of tightly coupled heterogeneous SANs. PktWay borrows heavily from the experience and wisdom of IP, with a few modifications needed for high performance. PktWay sacrifices generality and scalability to improve performance.

1a. PktWay and IP

IP is the general solution for "internetting" heterogeneous diverse networks, proven for over 25 years. However, IP was designed for the generality required for Wide Area Networks, without regard to the high performance requirements of tightly coupled systems. In addition, IP was designed to addresses "systems" rather than individual processors in MPPs (as PktWay does). For example, a 9,000 processor system is not expected to be assigned 9,000 IP addresses.

Cohen et al

[Page 2]

Internet-Draft PktWay End-to-End Protocol October 1997

PktWay is slightly below IP in the OSI Reference Model. It has many Level-3 features, like IP, but also can support IP as if PktWay was a Level-2 protocol. Hence, it is below IP. In addition, PktWay supports Level-2 optimizations (such as source routing).

Like IP, as a heterogeneous network layer, PktWay packets are transported by the native data-link layer of each SAN. As a result, PktWay packets are encapsulated with any native routing headers and trailers as required by the local network fabric.

Like IP, PktWay uses routers between its SANs. When an HR (half-router) receives a packet for a destination on its own SAN it forwards that packet directly to its destination. If the packet is for a destination out of this SAN, the HR forwards it to another HR which is en route to that destination.

Unlike IP, RRP defines the communication among the HRs both

intra-router and inter-router. In the IP environment only the intra-router communication is not defined, only the inter-router communication.

Unlike IP, the PktWay routers do not have to pop each packet back to Level-3, and are capable of operating entirely at Level-2, if this operation is requested by the communicating hosts. This Level-2 operation is discussed later in this introduction.

Like IP, the PktWay protocol utilizes the native capabilities of its constituent SANs and routers. PktWay defines neither how each HR maps the network in the SAN to which it is attached, nor how each half-router constructs SAN-headers for each of its hosts. The PktWay protocol also does not define how error-checking is conducted by each SAN (e.g., CRC8, CRC32, CRC64, or anything else). Instead, PktWay assumes that these capabilities are native to each SAN, and defines only how these maps are exchanged, and how these error indications are carried from where they were detected, to the destination node.

Like IP, the PktWay protocol defines neither how routes are selected, nor what corrective actions should be taken in case of faults. Instead, PktWay provides the information needed by the host nodes for devising routes and detecting and circumventing faults.

Like in IP environments, when hosts are powered up they may contact their default half-routers to register themselves and to inquire about other hosts (by name or node capabilities). This registration could be used in support of dynamic discovery procedures. The half-routers may help nodes discover each other (like IP's DNS) and may provide routing alternatives, possibly with different characteristics (e.g., MTU, length, and cost). The PktWay protocol does not specify how to choose among them.

Cohen	et	al	

[Page 3]

Internet-Draft

PktWay End-to-End Protocol

October 1997

To sum it all up, PacketWay has learned many lessons from IP, but has been heavily optimized for high performance SANs, while IP is the protocol of choice for WANs.

1b. General

PktWay supports resource discovery, by name or capabilities.

PktWay's unit of data is 64-bit long (8 bytes). Hence, a PktWay packet is always a multiple of 8B quantities. PktWay provides hosts with padding as required.

PktWay iself is big-Endian 8B-word based. Hence, the terms "first bit" and "first byte" are equivalent to MSbit and MSByte.

PktWay handles the Little vs. Big-Endian issue for its payload by providing a field in the EEP header which defines the endianness and "the chunk-size" of the data in the payload (Data Block). The intent is that byte-swapping hardware, if any, could be used to invert the endianness of payloads with uniform data elements (e.g., all the data being 32-bit floating point). Although this approach does not address the problems of transporting general structures (e.g., a "struct" of C), it does allows the participation of smart memory cards as PktWay nodes, as well as supporting direct memory access (DMA) operations.

The PktWay protocol is designed to allow wormhole (or "cut-through") forwarding, in which a router can start forwarding packets after receiving the first four bytes only (that include the PktWay-protocol version, priority, and the destination-type) without waiting for information that may not be needed for the packet forwarding task. This is unlike IP routers that receive the sender address before receiving the destination address, even though the former is not always needed whereas the latter is.

PktWay's addresses are short (23 bits) because, unlike IP, PktWay is not designed for global operation. The amount of state that is stored in the half-routers per node (type, name, paths, capabilities, etc.) makes it impractical for scalability beyond a few tens (hundreds?) of thousands of nodes, over a (relatively) small number of SANs.

PktWay does not support SAR (Segmentation And Reassembly). Instead, it provides means for hosts to discover the minimum transmission unit (MTU) over several alternative paths to any other node. A PktWay packet must never exceed the minimum MTU along all the network hops from the source node to the destination node.

Cohen et al

[Page 4]

October 1997

Internet-Draft PktWay End-to-End Protocol

Several protocol extensions, which are layered on the core PktWay protocol, have been defined. These include dynamic resource and routing discovery, secure PktWay, and multicast PktWay. These protocol extensions will be described in documents to be provided later.

1c. The Level-2 Operation of PktWay

PktWay's goal is to move data from a source node, (on some arbitrary SAN) to a destination node, (either on the same SAN, or on another SAN). Sources and destinations can be physical entities, such as a processor or a smart memory board, or logical entities, such as a group of cooperating processes or a collection of threads. Sources, destinations, and routers are such nodes.

Within each PktWay configuration all nodes have unique 23-bit physical PktWay addresses. A system designer can assign these PktWay addresses manually. Alternatively, the optional PktWay Server Layer may provide a way to assign and discover addresses dynamically. Throughout this document "address" always means the 23-bit physical PktWay address.

To optimize for performance, PktWay has a data transfer mode that directly leverages the native message routing schemes used within each SAN. This mode uses a "Planned Transfer" paradigm. During the planning phase, a source node collects information on optimal routes to a destination, expressed in the various native formats of all the intervening SANs. A source node later uses this information for low latency transfers to that destination. In PktWay, the transfer phase of a Planned Transfer is called "L2-forwarding". The RRP document demonstrates the use of L2-forwarding.

PktWay also supports a more traditional data transfer mode that requires no planning. Such transfers specify the destinations by their addresses only. In PktWay, this more traditional approach is called "L3-forwarding".

PktWay packets may be routed by Level-2 (L2) forwarding, Level-3 (L3) forwarding, or a combination thereof.

In L3-forwarding (similar to IP forwarding), the L2-routing through each SAN is determined by an inter-SAN router upon entering that SAN. The router prefixes the packet with an L2 routing header (such as a source route) corresponding to the destination address specified in the packet directing the packet either to its destination or to an intermediate router. It is a task for that router to determine the L2-routing-header corresponding to the given PktWay-address.

Cohen et al

[Page 5]

Internet-Draft PktWay End-to-End Protocol October 1997

In L2-forwarding the source prefixes the packet with all the L2-routing headers needed along the entire path to the

destination. Each router has only to get the L2-routing-header from the leading L2RH (L2-Routing-Header record) that was provided by the source.

PktWay allows hosts to construct a source-route built entirely of Level-2 headers, allowing each SAN to exploit the full performance of its native interconnection fabric. These SAN-headers (equivalent to MAC-headers) are provided by the SANs that will use them, in their native format. PktWay does not define the format of the local routing envelope. Instead, it defines how the encapsulated PktWay packets should be passed between half-routers, leaving it up to the local network of each SAN to properly deliver the packet.

If hosts so prefer, they can address their destinations either by any arbitrary name, a PktWay physical address (which is handled like the Level-3 IP-address), or by concatenating a sequence of Level-2 SAN-headers. Although the generation of a sequence of L2 Routing Headers requires more effort to construct initially, PktWay source routing results in considerably lower network latencies, as the packets are allowed to cut-through route through the intervening SAN networks .

2. A note about the PacketWay Documents

The PacketWay protocol is defined by a series of documents:

- * EEP (End-to-End Protocol)
- * RRP-1 (basic Router-to-Router Protocol)
- * RRP-2 (dynamic inter-SAN routing)
- * PktWay enumerations

Each of these documents should include the same "PacketWay at a Glance (Cheat-Sheet)", this note, and the Notations page. They should include also (as appendices) a copy of the PacketWay glossary of terms and its acronyms and abbreviations list.

The EEP and the RRP documents will be published first as Internet-Drafts and later as Proposed-Standards, Draft-Standards, and Standards.

The Enumeration Document will be first published as an "Informational-RFC" and later will be maintained by IANA.

The enumeration document may be attached to the EEP/RRP documents, as a matter of convenience. The enumeration is NOT a part of the PktWay standard, just as RFC0739 (the original "Assigned Numbers" RFC) is not a part of RFC0791, that defines IP.

Cohen et al

Similarly, the EEP-document has "Appendix-A: A Recommendation for PktWay Address Assignment" which is a recommendation only and NOT a part of the PktWay standard, just as IP-address-assignment is not a part of <u>RFC0791</u>, that defines IP.

The appendices are brought for clearance and convenience. They are not a part of the PktWay specification.

Information about the PktWay activity may be found in the URL: http://www.erc.msstate.edu/PktWay/

Notations

The shorter "PktWay" is used for "PacketWay".

8B means "8-byte" (64 bits).

0x indicates hexadecimal values, e.g., 0x0100 is 2^8=256(decimal).

Ob indicates binary values, e.g., Ob0100 is 4(decimal).

- xxxx indicate a field that is discarded without any checking (e.g., padding).
- [fff] indicates that fff is an optional field, when appropriate.
- [exp] in equations, is the integral part, rounded down, of `exp`. e.g., [23/8]=2.

All length fields do not include themselves, and therefore may be zero.

Lengths are specified either (a) by byte count, implying that some padding bytes may follow to fill 8B-words, or (b) by 8B-word count and PL, the number of trailing padding bytes (with PL between 0 and 7).

<u>4</u>. PktWay EEP Messages

4a. The Pktway Message Structure

PktWay messages have 6 components, including 4 optional ones:

[1]: [Optional Sequence of L2-Routing-Headers and Symbols] [2]: EEP Header (16 bytes) (PH)

- [3]: [Optional Header fields] (OH)
- [4]: [Optional, Most likely: Data Block] (DB)

[5]: [Optional Trailer fields] (OT)
[6]: EEP Trailer (8 bytes) (TAIL)

Cohen et al

[Page 7]

Internet-Draft PktWay End-to-End Protocol October 1997

4b. The Optional Fields

- [1]: as explained later, if the 9th+10th bits of a messages are Ob10 then the message starts with an L2RH, but if the 9th through the 12th bits of a message are Ob1111 then this message starts with a "symbol". The other values of these 4 bits indicate the lack of L2RH and symbols and that the message begins with the EEP-header.
- [3]: if the h-bit in the EEP header [2] is 1 then there are optional header (OH) fields. The sequence of these OH fields is terminated with an OH field marked as being the last one (with C=1).
- [4]: if DL>0, in the EEP header, zero then a Data Block (DB) is included in this message.
- [5]: the optional header fields, [3], may indicate that some optional trailer fields are present after the DB, [4]. The order and the formats of the trailer fields are defined by the optional header fields.

It is expected that most messages will have Data Blocks (DB), and that most messages will not have Optional Header fields (OH), nor Optional Trailer fields (OT).

Leading L2RHs and symbols [1] are consumed by the HRs before reaching the destination which receives only the other components, [2] through [6]. These parts, [2] to [6], constitute the End-to-End Protocol of PktWay.

TAIL, the EEP trailer, [6] may be modified along the way to the destination, unlike [2], [3], [4] and [5], which arrive exactly as sent by the source.

Each PktWay packet may be first L2-forwarded (zero or more times) before being L3-forwarded (zero or more times).

Although PktWay headers and trailers are always in Big Endian order, the byte order of the Data Block is not defined by PktWay.

Since all the elements of PktWay (L2RHs, EEP-headers, optional fields, data, and EEP-trailers) are always multiples of 8B-words,

it is recommended that PktWay headers (and data) be aligned on 8B-boundaries in the nodes' memory.

Cohen et al

[Page 8]

Internet-Draft PktWay End-to-End Protocol October 1997

5. Optional Sequence of L2RHs and Symbols [1]

PktWay messages may start with a mix of L2RHs and symbols.

A PktWay source may specify native routes, by placing the native routes before the PktWay Header. The native routes (for all SANs and LANs beyond the initial one) must appear within a sequence of PktWay L2-Routing-Header records (L2RH).

In certain situations symbols may be included among the L2RHs. These symbols are used for conveying information to the routers that handle the messages, such as about encryption. A symbol does not specify its destination and is processed (and consumed) by the entity that encounters it.

In L2-forwarding each intermediate HR consumes an L2RH and the preceeding symbols (if any). When a packet reaches its destination all of [1] (the Optional Sequence of L2RHs and Symbols) should be consumed.

5a. L2 Routing Headers (L2RHs) Records

The contents of the L2RH are totally SAN dependent, with the exception of the first 2 bytes that distinguish this record from an EEP-header and also provide the Length (0 < L < 64) indicating the number of routing bytes of that L2RH (not including these 2 bytes).

This distinction (between L2RHs and EEP-headers) is necessary for routers that L2-forward packets starting with L2RHs, but L3-forward packets starting with EEP-headers. Similarly, hosts expect packets to start with EEP-headers (with optionally preceeding symbols), and may discard packets that start with L2RHs.

It's up to each SAN to provide padding, as needed, to fill the L2RH words.

Each L2RH is defined by the entity that will process it. In addition to routing information per se, it may also include demuxing information such as a local message-type. For example, over Myrinet the L2RH should end with 0x0300 which is the Myrinet-type assigned to PktWay (and possibly some padding, too).

The L2RH must contain enough information to allow a router to create any necessary local routing headers and trailers. Although the low-level network implementation is beyond the scope of this document, the native source routing format must be documented in sufficient detail to allow for heterogeneous network interoperability.

Cohen et al

[Page 9]

Internet-Draft PktWay End-to-End Protocol October 1997

When a PktWay message is encapsulated inside any native SAN message (Paragon or Myrinet, for example), it's up to that SAN to distinguish between it and its own native packets. This is not a PktWay issue. For example, Myrinet uses its Message-Type to recognize PktWay messages.

PktWay-Routers on boundaries between SANs L2-forward packets starting with L2RH or L3-forward packets starting with EEP-headers. L2RH are distinguished from EEP-headers by the value of the first two bits of the Destination-Type field.

5a1. L2RH FORMAT:

Each L2RH is in the format:

+----+ |vv000000|10LLLLLL| SR01 | SR02 |.....|.....|.....| xxxx | +----+

 $\wedge \wedge$

The first 2 bits are vv=0b00 for the working version of the protocol. They may have other values for experimental versions.

The next 6 bits should be all zeroes.

The next two bits must be 0b10 to indicate that this is an L2RH record. This 0b10 was chosen to be consistent with the 0b10 of PktWay-addresses, as described in [2] below.

The next 6 bits are the byte count (L) of the routing information that starts in the next byte and is followed by as many padding bytes as needed to fill to the next 8B-boundary.

L does not include itself, hence it could be between 0 and 63. However, since this record contains some routing bytes, L is greater than 0. The total number of 8B-words in the L2RH is [(L+9)/8] where the square brackets indicate the integer part, rounded down, of the quantity within. Therefore, the number of padding bytes is PL=8*[(L+9)/8]-2-L.

5a2. L2RH EXAMPLES:

An L2RH with an SR with 5 routing bytes:

т	0b10	L=5	#1	т	#2	Ŧ	#3	т	#4	т	#5	- k	badding	т
vv0000	00 1000	00101	SR01		SR02		SR03		SR04		SR05		xxxx	 -+-
	^^	<		rou	ting i	nfo	rmatio	n -			->			- 1

Cohen et al

[Page 10]

An L2RH with an SR with 13 routing bytes:

	0b10	9 L=1	3	#1		#2		#3		#4		#5		#6	
+	+- 00 10	000110	-+- 1 -+-	SR01	-+-	SR02	-+-	SR03	-+-	SR04	-+-	SR05	-+-	SR06	+ · +
SR07		SR08		SR09		SR10		SR11		SR12		SR13		xxxx	
#7		#8	- + -	#9	- + -	#10	- + -	#11	- + -	#12	- + -	#13		paddinç))

5b. Symbol Records

5b1. Symbol Format:

Each symbol is in the format:

+----+ |vv000000|1111ssss|sssssss|sssssss| Length | data |.....| +----+ ^^^^<---- Symbol-Type --->

The 5th byte is the byte-count (L) of the data for this field that starts in the next byte, and is padded with as many padding bytes as needed to fill 8B-words.

The length (L) does not include itself, hence it is between 0 and 255. The total number of 8B-words in the symbol L2RH is [(L+12)/8] where the square brackets indicate the integer part, rounded down, of the quantity within. Therefore, the number of padding bytes is PL=8*[(L+12)/8]-2-L.

Symbols may be mixed among the L2RHs, before the EEP-header.

The values of the Symbol-Type field are defined in the PktWay Enumeration document.

5b2. Symbol Example:

A symbol with 9 data bytes.

 Cohen et al

[Page 11]

6. EEP Header [2]

The EEP (aka PH) has 16 bytes.

2	6	24			16		16	
+-+- V	Р	++ Destination	-++- i-Type	Туре	+ -Extensio	+ n	Packet-Type	+
+-+-	++ PL	Data-Length>=0	(8B-words) h	+	0 +	Source	-Address	+ +
+	4	3 2	25	1	7		24	Т

These fields are described below:

		Bytes.	bits
a.	Version	(V)	0.2
b.	Priority	(P)	0.6
с.	Destination-Type	(DT)	3.0
d.	Packet Type Extension	(TE)	2.0
e.	Packet Type	(PT)	2.0
f.	Endianness	(E)	0.4
g.	Padding Length	(PL)	0.3
h.	Data Length	(DL)	3.1
i.	Options flag	(h)	0.1
j.	Reserved	(RZ)	0.7
k.	Source Address	(SA)	3.0

6a. Version (V) 2 bits

This field is static. Its 2 bits are 0b00 for the working version of the protocol. These bits should have other values for co-existing experimental versions.

6b. Priority (P) unsigned integer, 6 bits

It is anticipated that some SANs, especially those working in real time, will want to implement priorities. This field supports such usage.

All ones is the highest priority, and all zeroes the lowest. Ideally, packets with higher priority should gain access to contested resources before packets with lower priority. Implementations may ignore the Priority field.

Cohen et al

[Page 12]

6c. Destination-Type (DT) 24 bits

The purpose of this field is to specify the header type, as well as the destination of the packet, when applicable.

This field may specify:

- * A physical PktWay address (of 23 bits);
- * An L2-Routing-Header (L2RH) of a variable length;
- * A logical address (of 20 bits); or
- * A symbol (of 20 bits).

In addition, it is anticipated that additional types will be needed in the future.

A variant of Huffman coding is used to accommodate all these methods for the Destination-Type field. This is done by assigning the MSbit of 0 to physical addresses, 2 MSbits of 0b10 to L2RH, 3 MSbits of 0b110 to future needs, 4 MSbits of 0b1110 to logical addresses, and 4 MSbits of Ob1111 to symbols.

This assignment is summarized in the following table:

MSbits | Method -----0xxx | Physical 10xx | L2RH 110x | Reserved 1110 | Logical 1111 | Symbol

A single C-style 16-way switch can dispatch quickly the protocol processor to the right handler required for any of the methods used to specify the destination.

The Physical addresses are unique within each instance of PktWay. Nodes should have addresses assigned to them. The method of assigning unique addresses within each PktWay is not specified here.

Examples of potentially addressable PktWay nodes include: groups of cooperating processes, an entire MPP, or each of an MPP's many processors or processes.

The Ob10xx was chosen for L2RH to be consistent with the Ob10 indication of L2RHs, as described earlier in this document.

"Logical Addresses" (e.g., for broadcast and for multicast groups)

are also in this address space. The destination-Type is a "Logical Address" if its 4 MSbits are set to 0b1110.

Cohen et al

[Page 13]

- A few Physical-addresses are reserved:
- 0x000000 Undefined address (illegal where an address is expected, but is allowed in the SA field)
- 0x7FFFFE ("Hey-You!") This address could be used at power up to address nodes or routers, over point-to-point links. ("If you receive it, it's for you.")
- 0x7FFFFF (Broadcast) This address is reserved for broadcast operations which may be added in later versions. ("If you receive it, it's for you.")

6d. Type Extension (TE) 2 bytes

An extension of the following PT field.

Logically, the TE should be after the PT. However, the PT is 8B-word aligned, easier to process than the TE which is 2B-aligned, but not 8B-aligned. Since the PT is more frequently used than the TE, it was assigned to the better aligned field.

6e. Packet Type (PT) 2 bytes

The PT field provides the information needed for efficient de-multiplexing of multiple protocol layers. Whereas traditional protocol layering requires several stages of sequential de-multiplexing, PktWay provides enough information to support a single combined de-multiplexing operation (such as in support of zero copy TCP). Thus, the PT field may indicate, for example, that the data blocks contain IP, SNMP, ATM, Ethernet, or other layered protocols.

PT values to support popular parallel programming APIs such as MPI have been defined. The PktWay Enumeration document defines several values for this PT field.

The PT field value of "RRP" indicates that message contains commands used in the PktWay Router-to-Router Protocol (RRP).

Some PTs will also use the 2 byte Type Extension (TE) field which precedes the PT for passing PT-specific parameters, such as implementation specific de-multiplexing information.

RRP messages (as described in the PktWay RRP document) use the TE field to distinguish among the various RRP-messages.

Cohen et al

[Page 14]

Special Packet Types

RRP - PktWay's Router/Router protocol (see the RRP document).

ERR - Error reporting packet, usually sent to the Source Address (SA, see below) in response to a PktWay message that could not be properly handled, such as "Destination Unknown." The TE indicates the nature of the error (e.g., UNK) as defined in the PktWay Enumeration document.

6f. Endianness (E), 4 bits

If the SAN interface of the receiving-node detects Endianness that is different than its own and if the entire Data Block (DB) consists of N-byte fields, then it may activate byte-swapping hardware for N-byte fields, saving much work for the receiving node.

The first bit (MSbit) of E, 'e' indicates whether the DB is in Big-Endian order (e=0) or in Little-Endian order (e=1). The next 3 bits could control hardware byte swapping, if any, which assumes that all the data consists of words of the same length.

The meaning associated with the values of the 3 LSbits of this field are defined in the PktWay enumeration document.

6g. Pad Length (PL) unsigned integer, 3 bits

The number of padding bytes that were added at the end of the DB (i.e., from the end of the data to the end of the DB). PL can be between 0 and 7.

6h. Data Length (DL) unsigned integer, 25 bits

Length, in 8B-words, of the data block, not including the L2RHs, EEP-header, OH, OT, and TAIL, including any optional padding. Hence, the net length of the Data Block is 8*DL-PL bytes. The minimum is zero, and the maximum length is $(2^{25-1})^{*8}$ bytes = -2^{28} = 256 MBytes.

6i. Optional Header-Field Flag (h) 1 bit

This bit is set to 1 if there are one (or more) optional header (OH) fields following the standard 16-byte EEP-header.

6j. Reserved (RZ) 7 bits

This field is reserved for future use. Applications should neither use it, nor count on others not to use it. It should be always set to zero (0b000000).

Cohen et al

[Page 15]

6k. Source Address (SA) 24 bit

This field contains the physical address of the packet's original source in the same format as the DT. However, unlike the DT, the SA must be a physical address.

Filling in this field is optional. A value of zero means that the SA is not specified.

Routers may use this field to identify the sender to which error messages may be returned.

7. Optional Header Fields (OH) [3]

A PktWay-message has Optional Header fields (OH) following the EEP-header, if the Option-Flag (h) is set to 1 in the EEP-header.

Each OH is in the format:

++-		+	+	+	+	++	
tttttttt LLLLLLL	data						
++-		+	+	+	+	++	

The first byte indicates the optional header field type (OH-TYPE).

The first bit, T, of the first byte indicates the processing of this OH-TYPE:

T=0: Optional (may drop this field if this OH-TYPE is unknown) T=1: Mandatory (should not process this message if this OH-TYPE is unknown)

The second bit, C, of the first byte indicates whether there are more header fields (i.e., whether this is the last field of this message).

C=0: More Optional Header fields follow C=1: End of Optional Header fields group (i.e., this is the last OH)

The other 6 bits of this byte, tttttt, define application-specific OH-TYPEs.

The second byte is the byte-count (L) of the data for this field that starts in the next byte, and is padded with as many padding bytes as needed to fill 8B-words.

The length (L) does not include itself, hence it is between 0 and 255. The total number of 8B-words in the symbol L2RH is [(L+9)/8] where the square brackets indicate the integer part, rounded down, of the quantity within. Therefore, the number of padding bytes is PL=8*[(L+9)/8]-2-L.

Cohen et al

[Page 16]

Internet-Draft

Example: An Optional Header Field (OH) with a mandatory OH-TYPE and 4 data bytes:

L=4 #1 #2 #3 #4 padding padding |1xtttttt|00000100| data01 | data02 | data03 | data04 | xxxx | xxxx | <----> value ---->

8. Optional Data Block (DB) [4]

The DB is free for applications to use in any way. Routers must not modify this field.

The DB has DL 8B-words, including optional padding (at the end) of PL bytes. Hence, the number of data bytes is 8*DL-PL. Both DL and PL are specified in the EEP-header.

The maximum length of the DB is $8*(2^2-1)B = -256$ MByte.

9. Optional Trailer Fields (OT) [5]

A PktWay-message has Optional Trailer fields (OT) if so indicated in an Optional Header field, e.g., an OH field may indicate that a CRC64 is in the OT.

An OT may have just the data for an OH defined above (following the EEP header), or be a stand alone, self-defined field in the same format as OH.

The OT-fields are in the order defined by the OHs. For example, if an OH-field indicating that a CRC32 is in the OT, is followed by another OH-fields indicating that a CRC64 is in the OT, then the OT with the CRC32 should be followed by the OT with the CRC64. Self defined OT fields must follow OTs defined by the OHs.

Cohen et al

[Page 17]

Internet-Draft

10. EEP Trailer (TAIL) [6]

The TAIL consists of only the Error Indication (EI) field which is a single 8B-word.

Routers may start forwarding packets toward their destinations before detecting transmission errors (such as in wormhole routing). The EI field provides such routers with a means to append an error indication to the end of a packet.

An all zero EI value means that no error was indicated. Any non-zero EI value indicates one or more errors.

The packet source will usually initialize the EI field to all zeros. However, as an alternative example, a memory board may create a packet with a non zero EI field (EI=1) that indicates that a parity error was detected by the memory board.

Each router does an arithmetic left shift, on the EI field by one bit unless its MSbit is 1. Routers that detect transmission errors also set the LSbit (after the shift) to 1.

This provides the ability to identify which routers have indicated errors (if the route is known).

11. Appendix-A: A Recommendation for PktWay Address Assignment

This section of the EEP document is a recommendation only, and not a part of the PktWay standard.

Unlike IP addresses, physical PktWay addresses are not globally unique, but must be locally unique within each PktWay configuration. Hence, when SANs that were developed independently are interconnected to form a PktWay, conflicting physical addresses may occur.

It is recommended not to attempt to assure local uniqueness of physical addresses by subdividing the global address space (hence, attempting to achieve global uniqueness).

Instead, it is recommended that every SAN would have local PktWay addresses, between 1 and the number of its local nodes, and also have a global "bias" to be added to all the addresses in that SAN. Hence, by proper setting of the biases of interconnected SANs, the local uniqueness of PktWay addresses is achieved.

The coordination of these biases is left (at least now) for manual (static) out-of-band coordination.

The use of such biases simplifies the mapping of physical addresses to their SANs.

Cohen et al

[Page 18]

Internet-Draft PktWay End-to-End Protocol October 1997

<u>12</u>. Appendix-B: Glossary

Address:	A unique designation of a node (actually an interface
	to that node) or a SAN.
Buddy-HR:	HRs are "buddies" if they are on the same SAN.
Cut-Thru:	See wormhole.
Destination:	The node to which a packet is intended
Dynamic-Routing	: Routing according to dynamic information
,	(i.e., acquired at run time, rather than pre-set).
Endianness:	The property of being Big-Endian or Little-Endian
	(transmission order, etc.)
Ethertype:	A 16-bit value designating the type of Level-3
	packets carried by a Level-2 communication system.
HR:	Half-Router, the part of a router that handles one
	network only.
12-Forwarding:	Forwarding based on Level-2 (i.e., data-link laver
22 i oi nai a±iigi	of the ISORM) information e.g. the native technique
	of each SAN or LAN Also called "source routing "
13-Eorwarding	Forwarding based on end-to-end
Lo i oi wai uiing.	(Level-3 i e network layer of the ISORM) addresses
	Also called "destination routing "
Man	The topology of a network
Mappor:	A node on a SAN/IAN that has the man and an PT
паррег	for that notwork. It is expected that the mapper
	dynamically undates the man and the PT
Multi homed Ned	a, A node with more than one network interface, where
Multi-Homed Nou	e. A node with more than one network interface, where
Nodou	Whatever can cond and receive packate
Noue.	
Nodo otructuro.	(e.g., a computer, an MPP, a software process, etc.)
Node structure:	A C-Struct (or equivalent) containing values for some
Dlannad Transfe	attributes of a node.
Planneu Transfe	r: Transfer of information, occurs after an initial
	phase in which the sender decides which Level-2 route
	to use for that transfer.
RUVF:	The "Received From" set includes all the physical
	adduces a through which on DT use discomingted
	addresses through which an RT was disseminated,
	addresses through which an RT was disseminated, starting with the address of the mapper that created
	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT.
Re-direct-messa	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be
Re-direct-messa	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address.
Re-direct-messa Router:	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device
Re-direct-messa Router: Security Contex	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device t: A relationship between 2 (or more) nodes that
Re-direct-messa Router: Security Contex	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device t: A relationship between 2 (or more) nodes that defines how the nodes utilize security services to
Re-direct-messa Router: Security Contex	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device t: A relationship between 2 (or more) nodes that defines how the nodes utilize security services to communicate securely.
Re-direct-messa Router: Security Contex Source:	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device t: A relationship between 2 (or more) nodes that defines how the nodes utilize security services to communicate securely. The node that created a packet.
Re-direct-messa Router: Security Contex Source: Source-Route:	addresses through which an RT was disseminated, starting with the address of the mapper that created that RT. ge: A message that tells nodes which HR should be used in order to get to a certain remote address. The inter-SAN communication device t: A relationship between 2 (or more) nodes that defines how the nodes utilize security services to communicate securely. The node that created a packet. A Level-2 route that is chosen for a packet by its source.

interleaving with the L2RHs.

Cohen et al

[Page 19]

Internet-Draft	PktWay End-to-End Protocol	October 1997
Twin-HR:	Two HRs are twins if they both are parts inter-SAN router.	of the same
Wormhole-routing	g: (aka cut-thru routing) forwarding packe switches as soon as possible, without sto entire packet in the switch (unlike Stop	ets out of oring that -and-forward)
Zero-copy TCP:	A TCP system that copies data directly be user area and the network device, bypass:	etween the ing OS copies

<u>13</u>. Appendix-C: Acronyms and Abbreviations

OXNNNNThe hexadecimal number NNNN (e.g., 0x0100 is 256-decimal)8B8 byte (64 bits) entityADDRThe Address-record of RRPAPINApplication/Program InterfaceATAddress TypeATMAsynchronous Transmission ModeBByte (e.g., 48)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVRL2An RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 for the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	ODNNNN	The binary number NNNN (e.g., 0b0100 is 4-decimal)
8B8 byte (64 bits) entityADDRThe Address-record of RRPAPInApplication/Program InterfaceATAddress TypeATMAsynchronous Transmission ModeBByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MShit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIGMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocol <t< td=""><td>0×NNNN</td><td>The hexadecimal number NNNN (e.g., 0x0100 is 256-decimal)</td></t<>	0×NNNN	The hexadecimal number NNNN (e.g., 0x0100 is 256-decimal)
ADDRThe Address-record of RRPAPInApplication/Program InterfaceATAddress TypeATMAsynchronous Transmission ModeBByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVRTAn RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDInternet protocolINFOAn RRP message providing information about nodesIPThe Internet protocolINFOLevel-2 for the ISORM (Link)L22RLevel-2 Routing HeaderL2SRSource Route	8B	8 byte (64 bits) entity
APInApplication/Program InterfaceATAddress TypeATMAsynchronous Transmission ModeBByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLevel-2 of the ISORM (Link)L22RLevel-2 Routing HeaderL2SRSource Route	ADDR	The Address-record of RRP
ATAddress TypeATMAsynchronous Transmission ModeBByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRT0An RRP message asking which HR to use for a given destinationIDIdentificationIGMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLevel-2 of the ISORM (Link)L22RLevel-2 Routing HeaderL2SRSource Route	APIn	Application/Program Interface
ATMAsynchronous Transmission ModeBByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message providing information about nodesIDThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	AT	Address Type
BByte (e.g., 4B)bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIGMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing Header	ATM	Asynchronous Transmission Mode
bbit (e.g., 32b)BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	В	Byte (e.g., 4B)
BCByte Count (of parameters)BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message providing information about nodesIPThe Internet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	b	bit (e.g., 32b)
BERBit Error RateCAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDIdentificationIDInternet group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	BC	Byte Count (of parameters)
CAPAThe CAPAbility-record of RRPCCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDIdentificationISORMThe IsO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	BER	Bit Error Rate
CCCapability CodeCSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDIdentificationISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	CAPA	The CAPAbility-record of RRP
CSRCommon Source-RouteDADestination AddressDBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIGMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	CC	Capability Code
 DA Destination Address DB Data Block DL Data Length (in 8B words) DSP Digital Signal Processor DT Destination-Type e The MSbit of E E The Endianness field (in the EEP header) EEP End/End Protocol EI Error Indication GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination I Identification I Identification I Identification I INFO An RRP message providing information about nodes IP The Internet protocol I ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR 	CSR	Common Source-Route
DBData BlockDLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationISMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	DA	Destination Address
DLData Length (in 8B words)DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRT0An RRP message asking which HR to use for a given destinationIDIdentificationIDInternet Group Management ProtocolINF0An RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	DB	Data Block
DSPDigital Signal ProcessorDTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDIdentificationISORMInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	DL	Data Length (in 8B words)
DTDestination-TypeeThe MSbit of EEThe Endianness field (in the EEP header)EEPEnd/End ProtocolEIError IndicationGPGeneral PurposeGVL2An RRP message, requesting L2 route to a given destinationGVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIDInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	DSP	Digital Signal Processor
 e The MSbit of E E The Endianness field (in the EEP header) EEP End/End Protocol EI Error Indication GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	DT	Destination-Type
 E The Endianness field (in the EEP header) EEP End/End Protocol EI Error Indication GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	е	The MSbit of E
 EEP End/End Protocol EI Error Indication GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRT0 An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INF0 An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	E	The Endianness field (in the EEP header)
 EI Error Indication GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	EEP	End/End Protocol
 GP General Purpose GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	EI	Error Indication
 GVL2 An RRP message, requesting L2 route to a given destination GVRT An RRP message asking an HR to give its routing tables h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	GP	General Purpose
GVRTAn RRP message asking an HR to give its routing tableshOptional header fields flagHRHalf RouterHRTOAn RRP message asking which HR to use for a given destinationIDIdentificationIGMPInternet Group Management ProtocolINFOAn RRP message providing information about nodesIPThe Internet protocolISORMThe ISO Reference ModelLLength field (exclusive of itself)L2Level-2 of the ISORM (Link)L2RHLevel-2 Routing HeaderL2SRSource Route	GVL2	An RRP message, requesting L2 route to a given destination
 h Optional header fields flag HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	GVRT	An RRP message asking an HR to give its routing tables
 HR Half Router HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	h	Optional header fields flag
 HRTO An RRP message asking which HR to use for a given destination ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route 	HR	Half Router
<pre>ID Identification IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route</pre>	HRT0	An RRP message asking which HR to use for a given destination
<pre>IGMP Internet Group Management Protocol INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route</pre>	ID	Identification
<pre>INFO An RRP message providing information about nodes IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route</pre>	IGMP	Internet Group Management Protocol
<pre>IP The Internet protocol ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route</pre>	INFO	An RRP message providing information about nodes
<pre>ISORM The ISO Reference Model L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route</pre>	IP	The Internet protocol
L Length field (exclusive of itself) L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route	ISORM	The ISO Reference Model
L2 Level-2 of the ISORM (Link) L2RH Level-2 Routing Header L2SR Source Route	L	Length field (exclusive of itself)
L2RH Level-2 Routing Header L2SR Source Route	L2	Level-2 of the ISORM (Link)
L2SR Source Route	L2RH	Level-2 Routing Header
	L2SR	Source Route

L3 Level-3 of the ISORM (Network) LA Logical Address LADR The Logical-addresses-record of RRP

Cohen et al

[Page 20]

Local Area Network LAN LRT Local Routing Table LSbit Least Significant bit LSbyte Least Significant byte Message Authentication Code / Media Access Control MAC MPI Message Passing Interface MPP Massively Parallel Processing system MSbit Most Significant bit MSbyte Most Significant byte MSU Mississippi State University MTU Maximum Transmission Unit MTUR The MTU-record of RRP M/C Multicast The name-record of RRP NAME NFS Network File Server Optional Header field OH OH-TYPE The Type of an Optional Header field 0Т Optional Trailer field Ρ The Priority field PAD Padding After Data PBD Padding Before Data PCI The Peripheral Component Interconnect "standard" PH PacketWay Header Padding Length (always in bytes) PL PPP The Point-to-Point Protocol Programmable ROM (Read-Only-Memory) PROM ΡT Packet Type (2B) PVM Parallel Virtual Machine PW The Myrinet Packet Type assigned to PktWay (PW=0x0300) Quality (of a path) Q RCVF Received-From list, or the Received-From record of RRP RDRC A re-direct message of RRP RH Routing Header Record ID RID Record Length (in 8B-words) RL RRP Router/Router Protocol RT-hd RT (Routing Table) header Routing Table RT An RRP message proving a Routing Table RTBL RTHD The Routing-Table-Header record of RRP RTyp RRP's Record Type The Reserved field (in the EEP header) RZ SA Source Address System Area Network SAN SAN-ID The 24-bit PktWay-address of a SAN Segmentation and Reassembly SAR Serial Number SN SAN-ID SNID SNMP Simple Network Management Protocol

PktWay End-to-End Protocol

October 1997

Internet-Draft

SR	Source	Route	(always	at	Level-2)	
----	--------	-------	---------	----	----------	--

- SRQR The Source-Route-and-Q-record of RRP
- ST Symbol Type

Cohen et al

[Page 21]

PktWay End-to-End Protocol October 1997 Internet-Draft TAIL PacketWay EEP Trailer ΤE Type Extension (2B) TELL An RRP message requesting information about nodes partially specified Unknown UNK Version V An RRP message asking its recipient to identify itself WRU? XRT External Routing Table xxxx A padding byte

<u>14</u>. Appendix-4: PktWay at a Glance (aka "The Cheat-Sheet")

2	6	type		24			1	L6	16			
+-+ V	Р	P Destination-Type					ype-E	Exten	sion	Packet-Type		
E	E PL	Dat	a-Le	ngth (8B	-words)	 h +-+-	RZ	0 +	So	urce-Addre	ess	
4	3			25		1	7	1		23	,	
	type = 0xxx Physical Address 10xx L2RH 110x Reserved 1110 Logical Address 1111 Symbols											
L2R 2	211: 211: 211: 211: 211: 211: 211: 211:	2	6	8	8		8	8	8	8	8	
+ V	Р	+ 10LLL	+ LLL.	SR01	+ SR02 +	+ 		+ 		+4	++ 	
		Len	igth									
Sym 2	bol: 6	4	6	8	8	Ŧ	8	+	3	8	8	
V	Р	11111s	sss	SSSSSSSS	ssssssss	Le	ength	da	ata			
+		+·	+ (Symbol	+> Туре>	+		+		+	+	
0pt 2	10na1 2 6	Header 8	:	8	8	т	8	т Т	8	8	8	
TC	ttttt	: LLLLL	LLL	data	+ 							
т:	0=opti	Lonal,	1=ma	ndatory;	C: 0=mo	re 0)H-fi€	elds	follo	w, 1=last	OH-field	
RRP	Recor 8	-d: 8	+	8	8	+	8	;+	8	8	8	
	RТур	PL	.	R	L							

+----+ RRP-messages: GVL2, L2SR, RDRC, TELL, INFO, HRTO, WRU, GVRT, RTBL; RTyp: ADDR, NAME, CAPA, LADR, SRQR, MTUR, RCVF, RTHD;

Cohen et al

[Page 22]

<u>15</u>. Security Considerations

This RFC raises no security issues.

<u>16</u>. Editor' Address

Danny Cohen Myricom, Inc. 325 N. Santa Anita Ave Arcadia, CA 91006

Phone: 626-821-5555 Fax: 626-821-5316 Email: Cohen@myri.com Cohen et al

[Page 23]