

# PCI-E note

## 2.0

SYSTEM software — BIOS  
operating SYSTEM — Linux

PCI Express 版本	推出	Line 编码	原始传输率 <sup>[i]</sup>	带宽 (带宽) <sup>[i]</sup>				
				x1	x2	x4	x8	x16
1.0	2003	8b/10b	2.5 GT/s	250 MB/s	0.50 GB/s	1.0 GB/s	2.0 GB/s	4.0 GB/s
2.0	2007	8b/10b	5.0 GT/s	500 MB/s	1.0 GB/s	2.0 GB/s	4.0 GB/s	8.0 GB/s
3.0	2010	128b/130b	8.0 GT/s	984.6 MB/s	1.97 GB/s	3.94 GB/s	7.88 GB/s	15.8 GB/s
4.0	2017	128b/130b	16.0 GT/s	1969 MB/s	3.94 GB/s	7.88 GB/s	15.75 GB/s	31.5 GB/s
5.0 <sup>[5][6]</sup>	2019 <sup>[7][8]</sup>	NRZ 128b/130b	32.0 GT/s <sup>[ii]</sup>	3938 MB/s	7.88 GB/s	15.75 GB/s	31.51 GB/s	63.0 GB/s
6.0	2021	PAM4 & FEC 128b/130b	64.0 GT/s	7877 MB/s	15.75 GB/s	31.51 GB/s	63.02 GB/s	126.03 GB/s

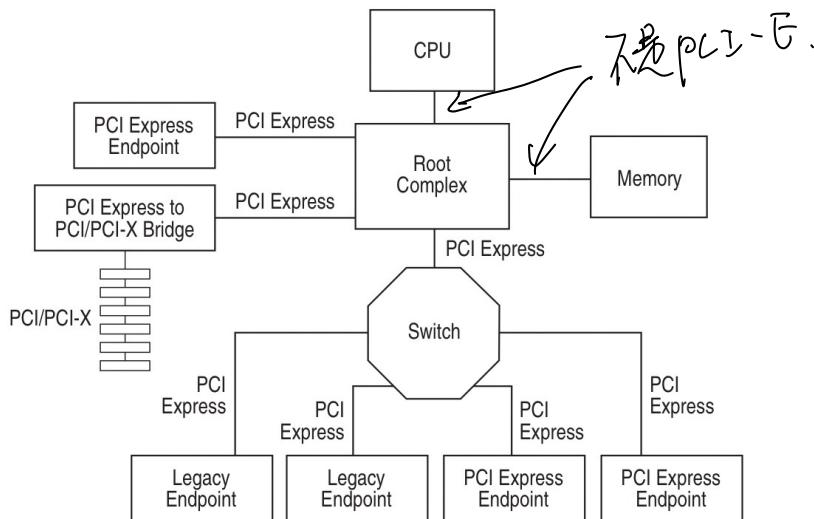
i. ^ 1.0.1.1 每条通道 (lane) 是全双工通道。

ii. ^ 出于技术可行性, 最初也考虑过25.0 GT/s

以PCIe 2.0为例, 每秒5GT (Gigatransfer) 原始数据传输率, 编码方式为8b/10b (每10个比特只有8个有效数据), 即有效带宽为4Gb/s = 500MByte/s。

### 1.3. PCI Express Fabric Topology

A fabric is composed of point-to-point Links that interconnect a set of components – an example fabric topology is shown in Figure 1-2. This figure illustrates a single fabric instance referred to as a hierarchy – composed of a Root Complex (RC), multiple Endpoints (I/O devices), a Switch, and a PCI Express to PCI/PCI-X Bridge, all interconnected via PCI Express Links.



OM13751A

Figure 1-2: Example Topology

1. RC > CPU 和 memory 是 I/O Port.
2. RL 有 T PCI-E Port, 因此是 -T domain
3. RL 可拆包.

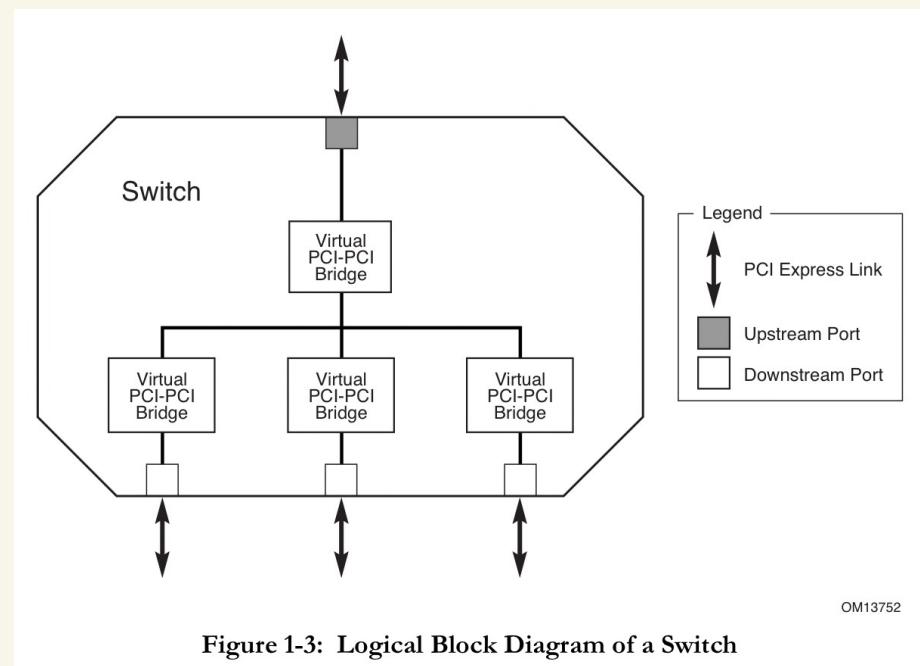
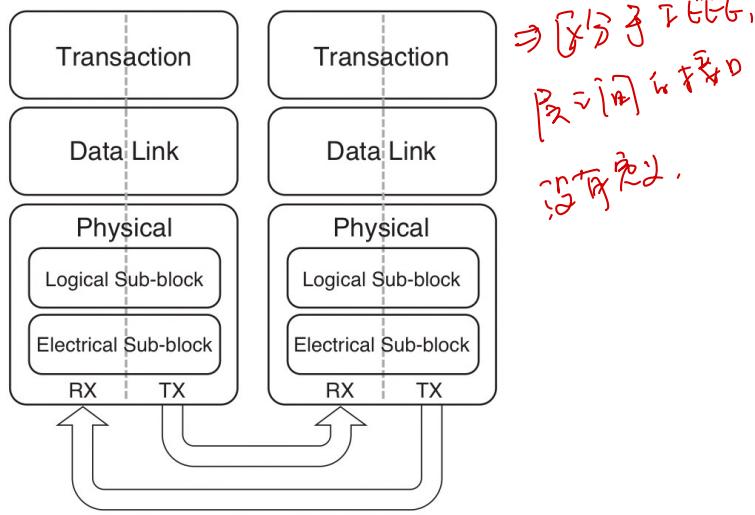


Figure 1-3: Logical Block Diagram of a Switch

1. PCI-E switch [多播基址] address-based routing.
2. switch 不可拆包
3. ingress port [向] IP 包 [向] round-robin & weighted round-robin

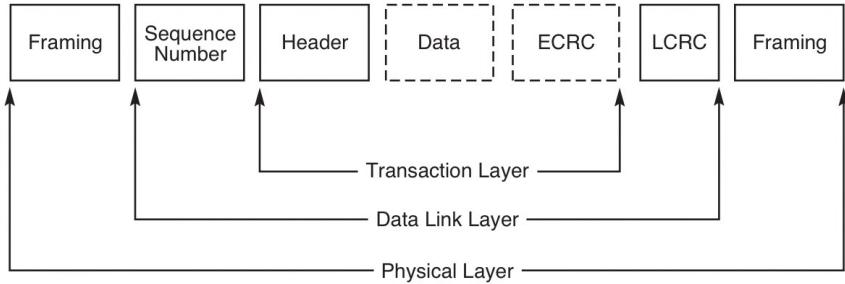
Note that this layering does not imply a particular PCI Express implementation.



OM13753

Figure 1-4: High-Level Layering Diagram

PCI EXPRESS BASE SPECIFICATION, REV. 2.0



OM13754

Figure 1-5: Packet Flow Through the Layers

## transation Layer:

1. communicate transactions, R&W

2. 每个包同一个标记

3. 4种协议类型

memory

I/O

configuration

adds message

加入了额外的  
头，用来管理信息。

源，电源管理信息。

5. flow-control

## Data Link Layer

1. 速率管理，纠错重传机制

2. IEEE 802.11 data protection code 算法

sequence number.

3. RX-检测，如果错，到重传队列中。

4. DUP 包顺序管理。

## physical layer

1. HW driver, input buffers:

2. 线缆，连接转换

3. PLL

4. interface inst.

At a high level, the key aspects of the Transaction Layer are:

- ❑ A pipelined full split-transaction protocol
- ❑ Mechanisms for differentiating the ordering and processing requirements of Transaction Layer Packets (TLPs)
- ❑ Credit-based flow control
- ❑ Optional support for data poisoning and end-to-end data integrity detection.

The Transaction Layer comprehends the following:

- ❑ TLP construction and processing
- ❑ Association of transaction-level mechanisms with device resources including:
  - Flow Control
  - Virtual Channel management
- ❑ Rules for ordering and management of TLPs
  - PCI/PCI-X compatible ordering
  - Including Traffic Class differentiation

This chapter specifies the behaviors associated with the Transaction Layer.

### **2.1.1. Address Spaces, Transaction Types, and Usage**

Transactions form the basis for information transfer between a Requester and Completer. Four address spaces are defined, and different Transaction types are defined, each with its own unique intended usage, as shown in Table 2-1.

**Table 2-1: Transaction Types for Different Address Spaces**

<b>Address Space</b>	<b>Transaction Types</b>	<b>Basic Usage</b>
Memory	Read Write	Transfer data to/from a memory-mapped location.
I/O	Read Write	Transfer data to/from an I/O-mapped location
Configuration	Read Write	Device Function configuration/setup
Message	Baseline (including Vendor-defined)	From event signaling mechanism to general purpose messaging

I/O space 是由 device 提供的，不折衷使用

- Read Request/Completion

- Write Request/Completion

I/O Transactions use a single address format:

- Short Address Format: 32-bit address

configure space

Configuration Transactions are used to access configuration registers of Functions within devices.

Configuration Transactions include the following types:

- Read Request/Completion

- Write Request/Completion

memory space

Memory Transactions include the following types:

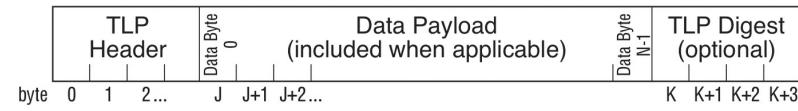
- Read Request/Completion

- Write Request

Memory Transactions use two different address formats:

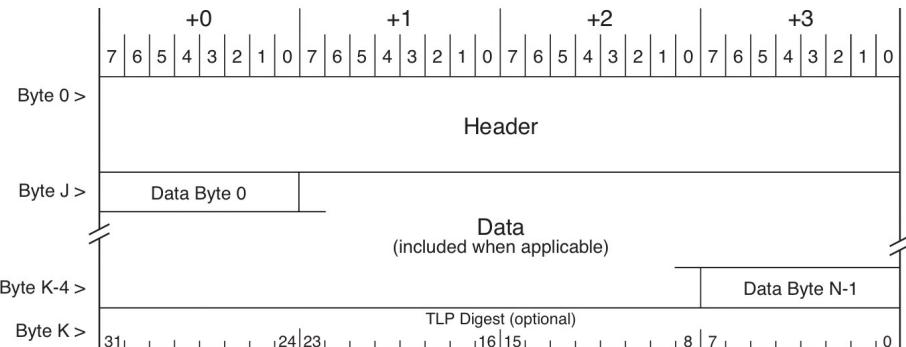
- Short Address Format: 32-bit address

- Long Address Format: 64-bit address



OM14547

Figure 2-2: Serial View of a TLP



OM13756

Figure 2-3: Generic TLP Format

byte / bit 序和internet不同. 字节内的大端. 小端是  
early address decode.

帧格式

1. 由 requests to complete trans [块]

{  
读回读  
I/O ack  
响应字}

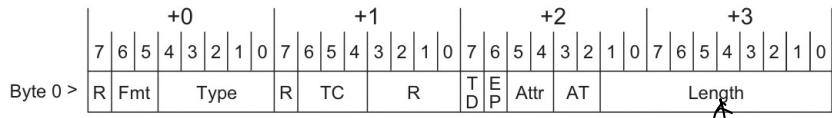


Figure 2-4: Fields Present in All TLPs

Table 2-2: Fmt[1:0] Field Values

Fmt[1:0]	Corresponding TLP Format
00b	3 DW header, no data
01b	4 DW header, no data
10b	3 DW header, with data
11b	4 DW header, with data

Table 2-3: Fmt[1:0] and Type[4:0] Field Encodings

TLP Type	Fmt [1:0] <sup>2</sup> (b)	Type [4:0] (b)	Description
MRd	00	0 0000	Memory Read Request
	01		
MRdLk	00	0 0001	Memory Read Request-Locked
	01		
MWr	10	0 0000	Memory Write Request
	11		
IORD	00	0 0010	I/O Read Request
IOWR	10	0 0010	I/O Write Request
CfgRd0	00	0 0100	Configuration Read Type 0
CfgWr0	10	0 0100	Configuration Write Type 0
CfgRd1	00	0 0101	Configuration Read Type 1
CfgWr1	10	0 0101	Configuration Write Type 1
TCfgRd	00	1 1011	Deprecated TLP Type <sup>3</sup>
TCfgWr	10	1 1011	Deprecated TLP Type <sup>3</sup>
Msg	01	1 0r <sub>2</sub> r <sub>1</sub> r <sub>0</sub>	Message Request – The sub-field r[2:0] specifies the Message routing mechanism (see Table 2-12).
MsgD	11	1 0r <sub>2</sub> r <sub>1</sub> r <sub>0</sub>	Message Request with data payload – The sub-field r[2:0] specifies the Message routing mechanism (see Table 2-12).
Cpl	00	0 1010	Completion without Data – Used for I/O and Configuration Write Completions and Read Completions (I/O, Configuration, or Memory) with Completion Status other than Successful Completion.
CplD	10	0 1010	Completion with Data – Used for Memory, I/O, and Configuration Read Completions.
CplLk	00	0 1011	Completion for Locked Memory Read without Data – Used only in error case.
CplDLk	10	0 1011	Completion for Locked Memory Read – otherwise like CplD.
			All encodings not shown above are Reserved.

Table 2-12: Message Routing

r[2:0] (b)	Description	Bytes 8 Through 15 <sup>5</sup>
000	Routed to Root Complex	Reserved
001	Routed by Address <sup>6</sup>	Address
010	Routed by ID	See Section 2.2.4
011	Broadcast from Root Complex	Reserved
100	Local – Terminate at Receiver	Reserved
101	Gathered and routed to Root Complex <sup>7</sup>	Reserved
110-111	Reserved - Terminate at Receiver	Reserved

Y-axis: 1000-1000  
X-axis: 0-1000

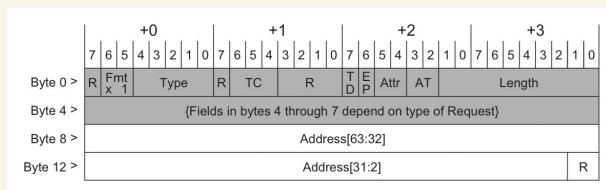
Table 2-4: Length[9:0] Field Encoding

Length[9:0]	Corresponding TLP Data Payload Size
00 0000 0001b	1 DW
00 0000 0010b	2 DW
...	...
11 1111 1111b	1023 DW
00 0000 0000b	1024 DW

读写 64/128 字节包提高系统吞吐量

# Router-type

## 1. address-based



← } memory req  
} I/O req

Figure 2-5: 64-bit Address Routing

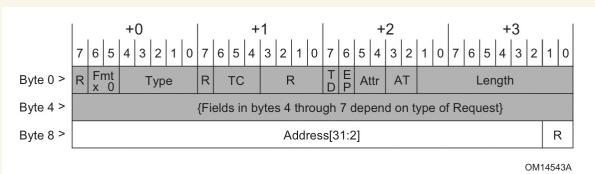


Figure 2-6: 32-bit Address Routing

## 2. ID-based

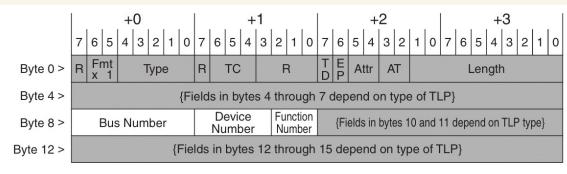
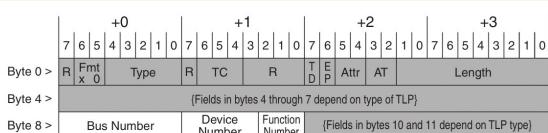


Figure 2-7: ID Routing with 4 DW Header

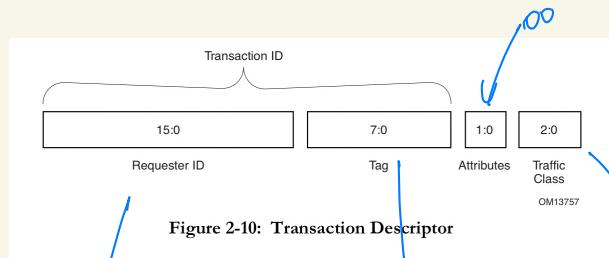


OM14541A

Figure 2-8: ID Routing with 3 DW Header

报文由4字节组成，所以无源时TE会向RTE-EN发送。来表示

字母读用作 flush, 保证之前的字母都生效。



render to ZD

$\rightarrow$  bus number  
 $\rightarrow$  device number  
 $\rightarrow$  function number

## ③ - 時間因变量 -

1. TAG大小决定可暂存同一类型 TEP；数量。
2. configuration space 中 capability to : TAG  
→ TEP size 为 8 bit  $\Rightarrow$
3. 强制优先级为 3 因为 T Non-posted 消息 = Cmpl

大部分右端模式 (BYTE4 - BYTE8)

OM13764A

Figure 2-13: Request Header Format for 64-bit Addressing of Memory

固定格式 for I/O

	+0	+1	+2	+3	
Byte 0 >	R Fmt x 0	Type	R TC 0 0 0	Reserved	
Byte 4 >	Requester ID		Tag	Last DW BE 0 0 0 0 1st DW BE	
Byte 8 >	Address[31:2]				R

OM13765A

Figure 2-15: Request Header Format for I/O Transactions

	+0	+1	+2	+3
Byte 0 >	R Fmt x 0	Type	R TC 0 0 0	Reserved
Byte 4 >	Requester ID		Tag	Last DW BE 0 0 0 0 1st DW BE
Byte 8 >	Bus Number	Device Number	Function Number	Reserved Ext. Reg. Number Register Number R

OM13766A

Figure 2-16: Request Header Format for Configuration Transactions

MSI / MSI-X 报文格式类似于 memory req  
地址由 memory 提供

	+0	+1	+2	+3
Byte 0 >	R Fmt x 1	Type	R TC	Reserved
Byte 4 >	Requester ID		Tag	Message Code
Byte 8 >	{Except as noted, bytes 8 through 15 are reserved.}			
Byte 12 >				

OM14539B

Figure 2-17: Message Request Header

message code + type [>:0]  
不固定 long version

所有讀請求和 Non-posted 写回需要 Completion.

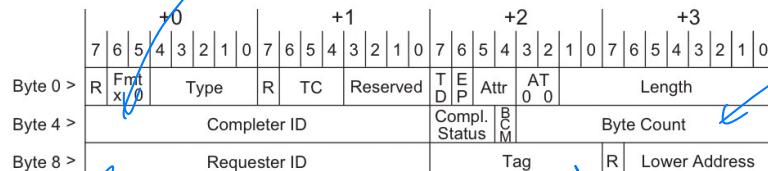


Figure 2-19: Completion Header Format

Table 2-21: Completion Status Field Values

Completion Status[2:0] Field Value (b)	Completion Status
000	Successful Completion (SC)
001	Unsupported Request (UR)
010	Configuration Request Retry Status (CRS)
100	Completer Abort (CA)
all others	Reserved

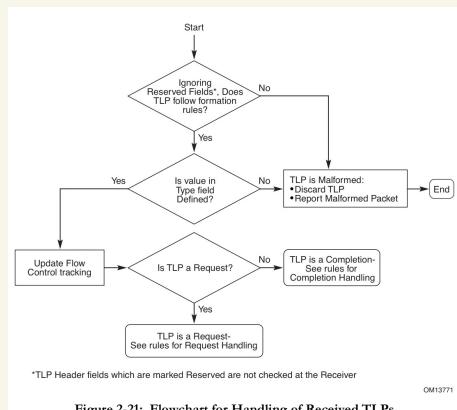


Figure 2-21: Flowchart for Handling of Received TLPs

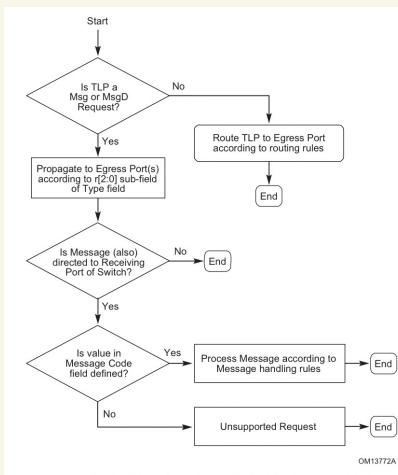


Figure 2-22: Flowchart for Switch Handling of TLPs

当 Read REQ 长度大于 RCB CPL  
 且 byte 成 128B， 来回为 128B 时， 有以下  
 往回于 completion

Table 2-22: Calculating Byte Count from Length and Byte Enables

1 <sup>st</sup> DW BE[3:0] (b)	Last DW BE[3:0] (b)	Total Byte Count
1xx1	0000 <sup>11</sup>	4
01x1	0000	3
1x10	0000	3
0011	0000	2
0110	0000	2
1100	0000	2
0001	0000	1
0010	0000	1
0100	0000	1
1000	0000	1
0000	0000	1
xxx1	1xxx	Length <sup>12</sup> * 4
xxx1	01xx	(Length * 4) - 1
xxx1	001x	(Length * 4) - 2
xxx1	0001	(Length * 4) - 3
xx10	1xxx	(Length * 4) - 1
xx10	01xx	(Length * 4) - 2
xx10	001x	(Length * 4) - 3
xx10	0001	(Length * 4) - 4
x100	1xxx	(Length * 4) - 2
x100	01xx	(Length * 4) - 3
x100	001x	(Length * 4) - 4
x100	0001	(Length * 4) - 5
1000	1xxx	(Length * 4) - 3
1000	01xx	(Length * 4) - 4
1000	001x	(Length * 4) - 5
1000	0001	(Length * 4) - 6

Table 2-24: Ordering Rules Summary Table

Row Pass Column?		Posted Request	Non-Posted Request		Completion	
		Memory Write or Message Request (Col 2)	Read Request (Col 3)	I/O or Configuration Write Request (Col 4)	Read Completion (Col 5)	I/O or Configuration Write Completion (Col 6)
Posted Request	Memory Write or Message Request (Row A)	a) No b) Y/N	Yes	Yes	a) Y/N b) Yes	a) Y/N b) Yes
Non-Posted Request	Read Request (Row B)	No	Y/N	Y/N	Y/N	Y/N
	I/O or Configuration Write Request (Row C)	No	Y/N	Y/N	Y/N	Y/N
Completion	Read Completion (Row D)	a) No b) Y/N	Yes	Yes	a) Y/N b) No	Y/N
	I/O or Configuration Write Completion (Row E)	Y/N	Yes	Yes	Y/N	Y/N



## IMPLEMENTATION NOTE

### Large Memory Reads vs. Multiple Smaller Memory Reads

Note that the rule associated with entry D5b in Table 2-24 ensures that for a single Memory Read Request serviced with multiple Completions, the Completions will be returned in address order. However, the rule associated with entry D5a permits that different Completions associated with distinct Memory Read Requests may be returned in a different order than the issue order for the Requests. For example, if a device issues a single Memory Read Request for 256 bytes from location 1000h, and the Request is returned using two Completions (see Section 2.3.1.1) of 128 bytes each, it is guaranteed that the two Completions will return in the following order:

1<sup>st</sup> Completion returned: Data from 1000h to 107Fh.

2<sup>nd</sup> Completion returned: Data from 1080h to 10FFh.

However, if the device issues two Memory Read Requests for 128 bytes each, first to location 1000h, then to location 1080h, the two Completions may return in either order:

1<sup>st</sup> Completion returned: Data from 1000h to 107Fh.

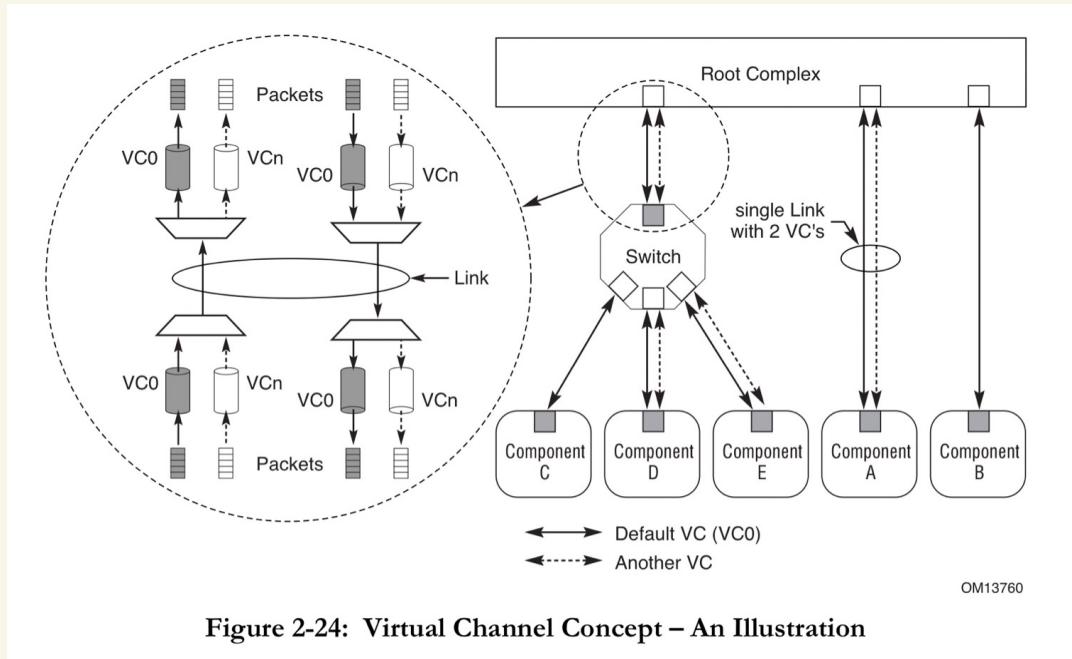
2<sup>nd</sup> Completion returned: Data from 1080h to 10FFh.

— or —

1<sup>st</sup> Completion returned: Data from 1080h to 10FFh.

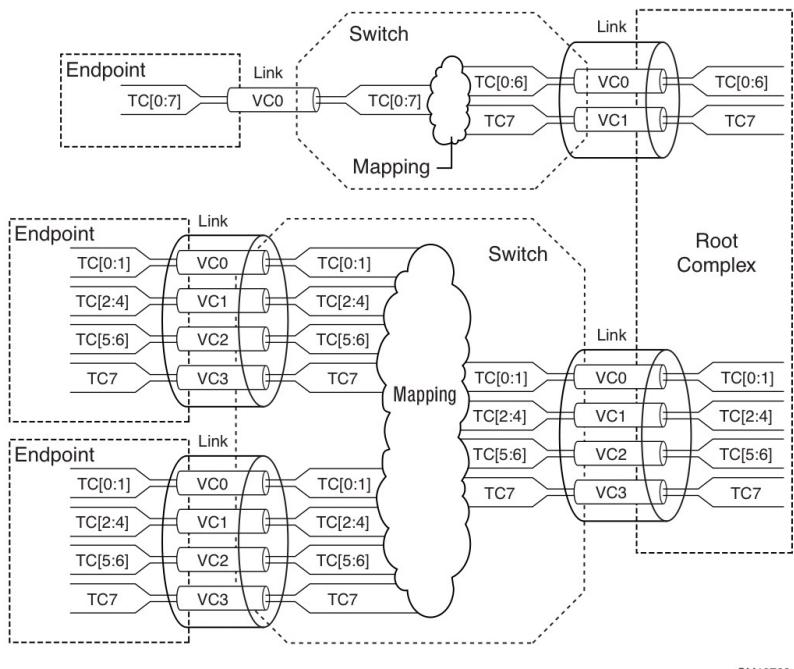
2<sup>nd</sup> Completion returned: Data from 1000h to 107Fh.

VC (Virtual channel) — 同 TC (Traffic Class) 之区别不同  
 to traffic, 隔离同一台设备但连接瓶颈  
 TCO 硬关联 VCD



每个 VC 通过自己的队列 queue, buffer 和控制逻辑.

TC → VL 是一对一关系, 而 T  
 port 有多个是 mapping.



OM13762

Figure 2-26: An Example of TC/VC Configurations

- flow control 交易层的流控制
1. 使用FCP, 是DCP的一种。
  2. 使用credit为基底，4DW为单位。
  3. 不同的VL有单独的flow control。
  4. 有6种不同类型分类

Table 2-26: Flow Control Credit Types

Credit Type	Applies to This Type of TLP Information
PH	Posted Request headers
PD	Posted Request Data payload
NPH	Non-Posted Request headers
NPD	Non-Posted Request Data payload
CPLH	Completion headers
CPLD	Completion Data payload

Table 2-27: TLP Flow Control Credit Consumption

TLP	Credit Consumed <sup>14</sup>
Memory, I/O, Configuration Read Request	1 NPH unit
Memory Write Request	1 PH + n PD units <sup>15</sup>
I/O, Configuration Write Request	1 NPH + 1 NPD Note: size of data written is never more than 1 (aligned) DW
Message Requests without data	1 PH unit
Message Requests with data	1 PH + n PD units
Memory Read Completion	1 CPLH + n CPLD units
I/O, Configuration Read Completions	1 CPLH unit + 1 CPLD unit
I/O, Configuration Write Completions	1 CPLH unit

to DATA Link Layer. 有 32 位 CRC, 2P + CRC.

to Transaction Layer. 有 32 位 CRC, 2P + ECRC.

{ ECRC. 位于发送端, 为链路接收到的.

LURL, 位于发送端, 双向发生会重新计算.

EP 用于 error forwarding. → { read completion DATA,  
write DATA.

DL - down by transaction layer 通过所有子层到 DL.

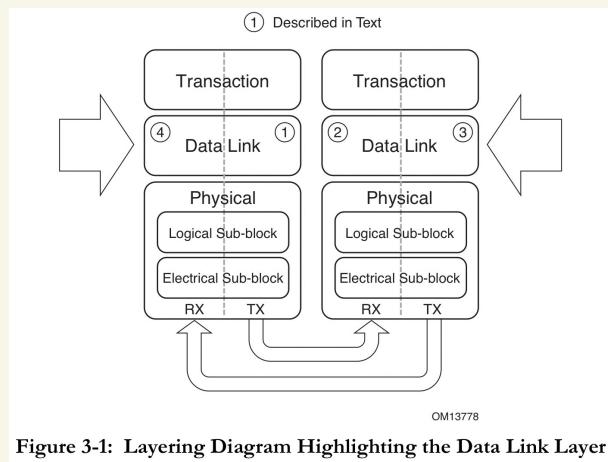


Figure 3-1: Layering Diagram Highlighting the Data Link Layer

DATA Link Layer 提供 TCP 及至的可靠化。

1. 接收 TCP 下达至物理层

2. 接受物理层上送至 TL.

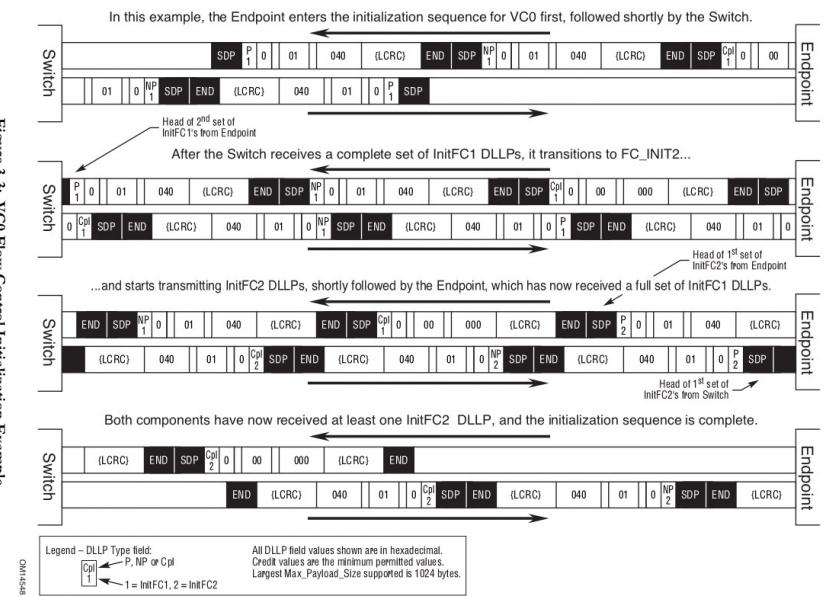
3. 错误检测及重传.

a. TCP 序列号

b. 支持 TCP 外延首尾

c. 基于免赔技术规范

d. 指挥 Link 层



DLLPs : Ack, Nak, InitFC1 & FC2 - Update Fc  
flow control  
 for power manager

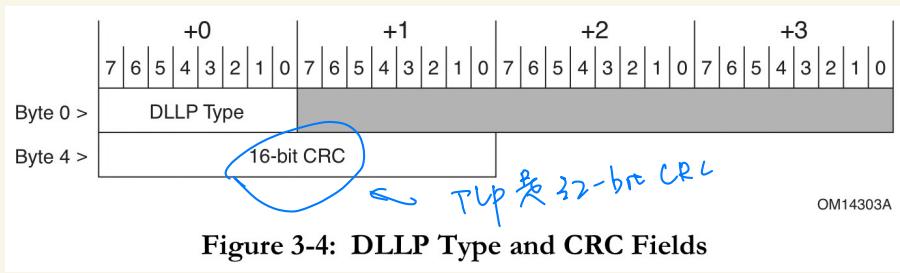


Figure 3-4: DLLP Type and CRC Fields

Table 3-1: DLLP Type Encodings

Encodings	DLLP Type
0000 0000	Ack
0001 0000	Nak
0010 0000	PM_Enter_L1
0010 0001	PM_Enter_L23
0010 0011	PM_Active_State_Request_L1
0010 0100	PM_Request_Ack
0011 0000	Vendor Specific – Not used in normal operation
0100 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC1-P (v[2:0] specifies Virtual Channel)
0101 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC1-NP
0110 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC1-Cpl
1100 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC2-P
1101 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC2-NP
1110 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	InitFC2-Cpl
1000 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	UpdateFC-P
1001 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	UpdateFC-NP
1010 0V <sub>2</sub> V <sub>1</sub> V <sub>0</sub>	UpdateFC-Cpl
All other encodings	Reserved

→ TLP 之類似於 TCP 協議  
 → PS 類似於 TCP 的序號  
 → power management  
 P: post  
 NP: no-post  
 Cpl: completion

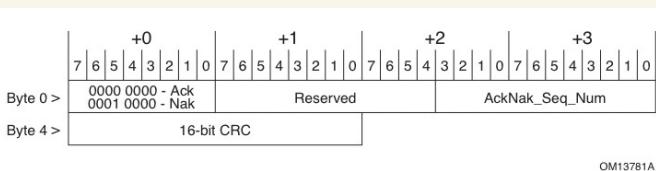


Figure 3-5: Data Link Layer Packet Format for Ack and Nak

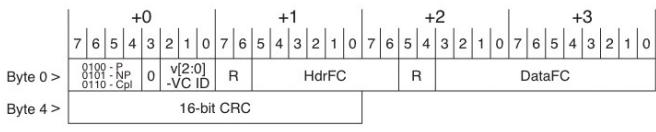


Figure 3-6: Data Link Layer Packet Format for InitFC1

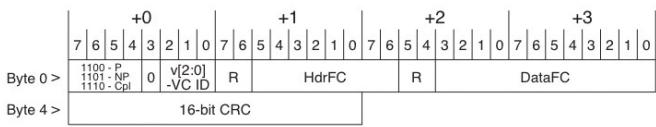


Figure 3-7: Data Link Layer Packet Format for InitFC2

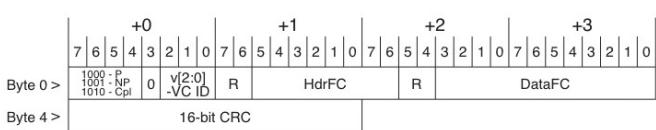


Figure 3-8: Data Link Layer Packet Format for UpdateFC

将 TLP 下发至 DATA Link Layer 为操作，会加上 Sequence Number 和 32bit 的 LCRC [Link CRC]  
 $\text{seq} := (\text{seq} + 1) \bmod 4096$

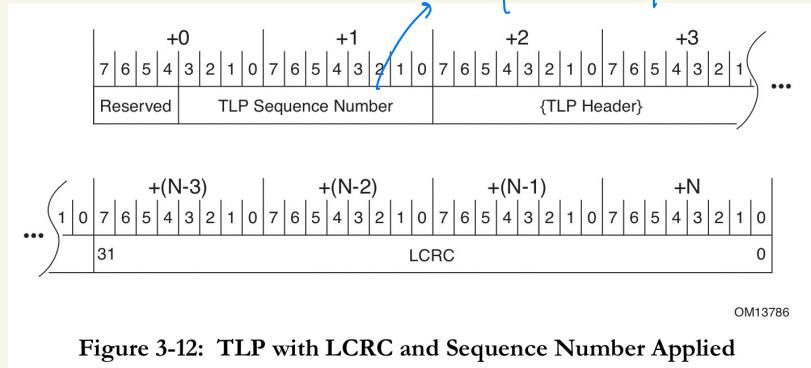


Figure 3-12: TLP with LCRC and Sequence Number Applied

通常来说，发送方会将每一个 TLP 放入 Replay buffer 中以备查，直到其接收到来自接收方的 ACK DLLP，确认以该 TLP 已成功的接收，才会删除这个 TLP。如果接收方发现 TLP 已成功的接收，才会删除这个 TLP。如果接收方接收到 TLP 后未接收到 ACK DLLP，则会向发送方发送 Nak DLLP。之后发送方会从 Replay buffer 中取出数据，重新发送该 TLP。

transmission 有优先级。

- 1) Completion of any transmission (TLP or DLLP) currently in progress (highest priority)
- 2) Nak DLLP transmissions
- 3) Ack DLLP transmissions scheduled for transmission as soon as possible due to:  
receipt of a duplicate TLP —OR—  
expiration of the Ack latency timer (see Section 3.5.3.1)
- 4) FC DLLP transmissions required to satisfy Section 2.6
- 5) Retry Buffer re-transmissions
- 6) TLPs from the Transaction Layer

- 7) FC DLLP transmissions other than those required to satisfy Section 2.6
- 8) All other DLLP transmissions (lowest priority)

解法:

1.  $T_{expected\_req} < T_{expected\_ack}$ .
2.  $\text{ACKNAK-LATENCY-TIMER}$ .
3.  $T_{expected\_req} + T_{expected\_ack} > T_{expected\_TLP}$ .
4.  $(ACKNAK, ACK)$ .
  - a.  $T_{expected\_req} < T_{expected\_TLP}$  &  $T_{expected\_ack} > T_{expected\_TLP}$  (duplicate).

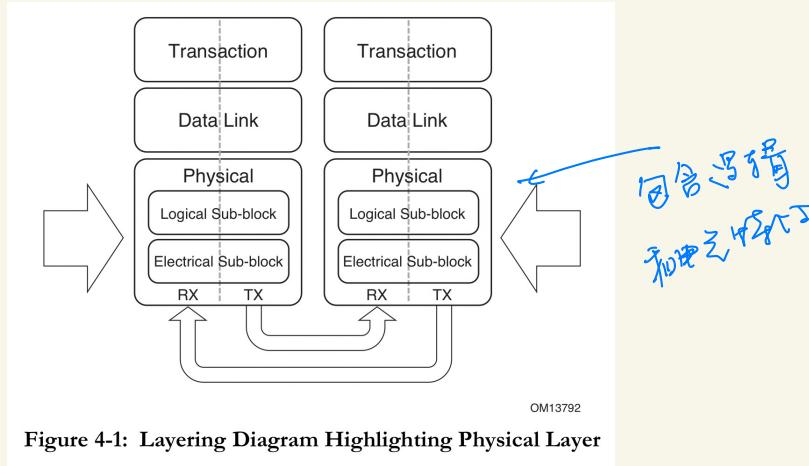


Figure 4-1: Layering Diagram Highlighting Physical Layer

gen2 - 8B/10B encoding.  
由 Lane 0 到 symbol 下发

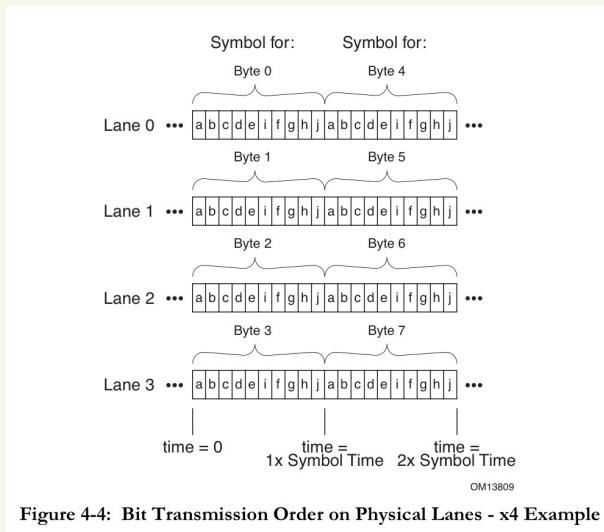


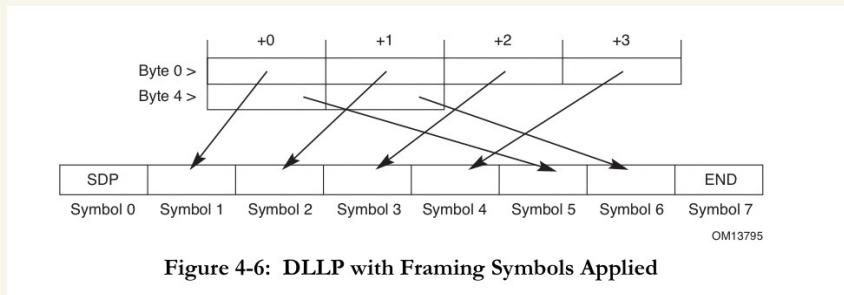
Figure 4-4: Bit Transmission Order on Physical Lanes - x4 Example

8B/10B 中有些格刷满 symbol. 可以区分不同种类的 frame  
或用于 link management.  
当链路空闲时，进入 logic IDLE 状态，此时持续且  
通知，因此添加 8B/10B

Table 4-1: Special Symbols

Encoding	Symbol	Name	Description
K28.5	COM	Comma	Used for Lane and Link initialization and management
K27.7	STP	Start TLP	Marks the start of a Transaction Layer Packet
K28.2	SDP	Start DLLP	Marks the start of a Data Link Layer Packet
K29.7	END	End	Marks the end of a Transaction Layer Packet or a Data Link Layer Packet
K30.7	EDB	End Bad	Marks the end of a nullified TLP
K23.7	PAD	Pad	Used in Framing and Link Width and Lane ordering negotiations
K28.0	SKP	Skip	Used for compensating for different bit rates for two communicating Ports
K28.1	FTS	Fast Training Sequence	Used within an Ordered Set to exit from L0s to L0
K28.3	IDL	Idle	Used in the Electrical Idle Ordered Set (EIOS)
K28.4			Reserved
K28.6			Reserved
K28.7	EIE	Electrical Idle Exit	Used in 2.5 GT/s Used in the Electrical Idle Exit Ordered Set (EIEOS) and sent prior to sending FTS at speeds other than 2.5 GT/s

和 IEEE 802.10B 中的 8B/10B一样存在 disparity。而且由于 lane 交错，  
TX 方向的问题，RX 方向的 symbol locking 从第一个 symbol 开始。  
信号 symbol → Lane 分发  
order set : 每个 lane 都会发  
TLP & DLLP : per Lane 分发  
By Lane 时 STP / SDP 必须在 Lane  $\text{mod } 4 = 0$  上  
END & EDB 在 Lane width - 1 上时，填补上 PAD



扰码，如果仅限于扰码，则其在 TX 端 8B/10B 的后 RX 端在

8B/10B 是什么解码

有关规定为 side-stream，且只对 ~~data symbols~~ special symbols 有效，COM 符号不受此限  
AT5.

disable scrambling is implementation specific and beyond the scope of this specification.

The data scrambling rules are the following:

- The COM Symbol initializes the LFSR.
- The LFSR value is advanced eight serial shifts for each Symbol except the SKP.
- All data Symbols (D codes) except those within a Training Sequence Ordered Sets (e.g., TS1, TS2) and the Compliance Pattern (see Section 4.2.8) are scrambled.
- All special Symbols (K codes) are not scrambled.
- The initialized value of an LFSR seed (D0-D15) is FFFFh. Immediately after a COM exits the Transmit LFSR, the LFSR on the Transmit side is initialized. Every time a COM enters the Receive LFSR on any Lane of that Link, the LFSR on the Receive side is initialized.
- Scrambling can only be disabled at the end of Configuration (see Section 4.2.6.3).
- Scrambling does not apply to a loopback slave.
- Scrambling is always enabled in Detect by default.

Link init 和 training.

The following are discovered and determined during the training process:

- Link width
- Link data rate
- Lane reversal
- polarity inversion.

Training does:

- Link data rate negotiation.
- Bit lock per Lane
- Lane polarity
- Symbol lock per Lane
- Lane ordering within a Link
- Link width negotiation
- Lane-to-Lane de-skew within a multi-Lane Link.

training DS 会先 TS1 or TS2, SB10B 但不适用

Table 4-2: TS1 Ordered Set

Symbol Number	Encoded Values	Description
0	K28.5	COM for Symbol alignment
1	D0.0 - D31.7, K23.7	Link Number within component
2	D0.0 - D31.0, K23.7	Lane Number within Port
3	D0.0 - D31.7	N_FTS. This is the number of Fast Training Sequences required by the Receiver to obtain reliable bit and Symbol lock.

Table 4-3: TS2 Ordered Set

Symbol Number	Encoded Values	Description
0	K28.5	COM for Symbol alignment
1	D0.0 - D31.7, K23.7	Link Number within component
2	D0.0 - D31.0, K23.7	Lane Number within Port
3	D0.0 - D31.7	N_FTS. This is the number of Fast Training Sequences required by the Receiver to obtain reliable bit and Symbol lock.

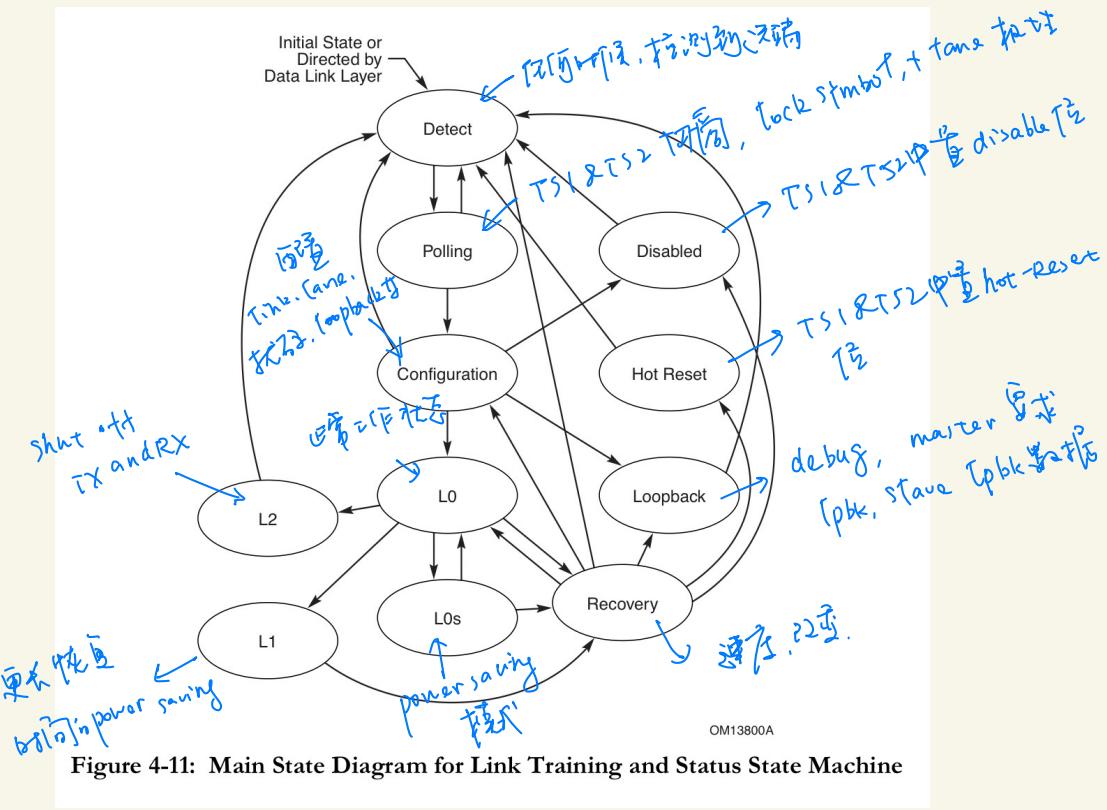
Symbol Number	Encoded Values	Description
4	D2.0, D2.2, D2.4, D2.6, D6.0, D6.2, D6.4, D6.6	<p>Data Rate Identifier</p> <p>Bit 0 – Reserved, set to 0b.</p> <p>Bit 1 – When set to 1b, indicates 2.5 GT/s data rate supported.</p> <p>Bit 2 – When set to 1b, indicates 5.0 GT/s data rate supported.</p> <p>Devices that advertise the 5 GT/s data rate must also advertise support for the 2.5 GT/s data rate (i.e., set Bit 1 to 1b).</p> <p>Bit 3:5 – Reserved, must be set to 0b.</p> <p>Bit 6 (Autonomous Change) –</p> <p>Downstream component: Autonomous Change&gt;Selectable De-emphasis: When set to 1b in Configuration state and LinkUp = 1b, indicates that the speed or Link width change initiated by the Downstream component is not caused by a Link reliability issue.</p> <p>In Recovery state, this bit indicates the de-emphasis preference of the Downstream component.</p> <p>In Polling.Active substate, this bit specifies the de-emphasis level the Upstream component must operate in if it enters Polling.Compliance and operates in 5.0 GT/s data rate.</p> <p>In Configuration.Linkwidth.Start substate with LinkUp = 0b and in the Loopback.Entry substate, this bit specifies the de-emphasis level the Upstream component must operate in if it enters Loopback state (from Configuration) and operates in 5.0 GT/s data rate. For de-emphasis, a value of 1b indicates -3.5 dB de-emphasis and a value of 0b indicates -6 dB de-emphasis.</p> <p>This bit is reserved in all other states for a Downstream component.</p> <p>Upstream component: In Polling.Active,</p> <p>Configuration.Linkwidth.Start, and Loopback.Entry substates, this bit specifies the de-emphasis level the Downstream component must operate in 5.0 GT/s data rate if it enters Polling.Compliance and Loopback states, respectively. A value of 1b indicates -3.5 dB de-emphasis and a value of 0b indicates -6 dB de-emphasis.</p> <p>This bit is reserved for all other states.</p> <p>Bit 7 (speed_change) – When set to 1b, indicates a request to change the speed of operation. This bit can be set to 1b only during Recovery.RcvrLock state.</p> <p>All Lanes under the control of a common LTSSM must transmit the same value in this Symbol. Transmitters must advertise all supported data rates in Polling.Active and Configuration.LinkWidth.Start substates, including data rates they do not intend to operate on.</p>

Symbol Number	Encoded Values	Description
4	D2.0, D2.2, D2.4, D2.6, D6.0, D6.2, D6.4, D6.6	<p>Data Rate Identifier</p> <p>Bit 0 – Reserved, set to 0b</p> <p>Bit 1 – When set to 1b, indicates 2.5 GT/s data rate supported.</p> <p>Bit 2 – When set to 1b, indicates 5.0 GT/s data rate supported.</p> <p>Devices that advertise the 5 GT/s data rate must also advertise support for the 2.5 GT/s data rate (i.e., set Bit 1 to 1b).</p> <p>Bit 3:5 – Reserved, must be set to 0b.</p> <p>Bit 6 (Autonomous Change/Link upConfigure Capability&gt;Selectable De-emphasis) – This bit is used for Link upConfigure capability notification, as well as Selectable De-emphasis by an Upstream component, and for Link upConfigure capability notification as well as autonomous bandwidth change notification by a Downstream component. In Configuration Complete substate (both for Upstream and Downstream components): When set to 1b, indicates the capability of the component to upConfigure the Link to a previously negotiated Link width during the current LinkUp = 1b state. The device advertising this optional capability must be capable of supporting at least a 1x Link in the Lane 0 assigned in the Configuration state.</p> <p>In Recovery state for a Downstream component: When set to 1b by the Downstream component, indicates that the speed or Link width change initiated by the Downstream component is not caused by a Link reliability issue.</p> <p>In Recovery state for an Upstream component: This bit must be set to 1b if the Upstream component wants the Link to operate in -3.5 dB in 5.0 GT/s speed and reset to 0b if the Upstream component wants the Link to operate in -6 dB in 5.0 GT/s speed if it is advertising 5.0 GT/s data rates (i.e., bit 2 of Symbol 4 is 1b).</p> <p>In Polling.Configuration substate both by Upstream and Downstream components: This indicates the de-emphasis level of the component in 5.0 GT/s data rate if the receiving device enters Loopback from Configuration. A value of 1b indicates -3.5 dB de-emphasis and a value of 0b indicates -6 dB de-emphasis.</p> <p>This bit is reserved in all the other states not covered above for an Upstream or Downstream component.</p> <p>Bit 7 (speed_change) – When set to 1b, indicates a request to change the speed of operation.</p> <p>This bit can be set to 1b only during Recovery.RcvrCtg state. All Lanes under the control of a common LTSSM must transmit the same value in this Symbol. Transmitters must advertise all supported data rates in Polling Configuration substate, including data rates they do not intend to operate on.</p>

Symbol Number	Encoded Values	Description
5	D0.0, D1.0, D2.0, D4.0, D8.0, D16.0, D20.0	<p>Training Control</p> <p>Bit 0 – Hot Reset</p> <p>Bit 0 = 0b, De-assert</p> <p>Bit 0 = 1b, Assert</p> <p><u>Bit 1 – Disable Link</u></p> <p>Bit 1 = 0b, De-assert</p> <p>Bit 1 = 1b, Assert</p> <p><u>Bit 2 – Loopback</u></p> <p>Bit 2 = 0b, De-assert</p> <p>Bit 2 = 1b, Assert</p> <p><u>Bit 3 – Disable Scrambling</u></p> <p>Bit 3 = 0b, De-assert</p> <p>Bit 3 = 1b, Assert</p> <p><u>Bit 4 – Compliance Receive</u></p> <p>Bit 4 = 0b, De-assert</p> <p>Bit 4 = 1b, Assert</p> <p>Components that support 5 GT/s data rate must implement this bit as specified. Components that support only 2.5 GT/s data rate may optionally implement this bit in a Receiver. If not implemented for components that support only 2.5 GT/s data rate, this bit will be reserved and must behave as if the component received a 0b in this bit position.</p> <p><u>Bit 5:7 – Reserved</u></p> <p>Set to 0b</p>
6 – 15	D10.2	TS1 Identifier

Symbol Number	Encoded Values	Description
5	D0.0, D1.0, D2.0, D4.0, D8.0	<p>Training Control</p> <p>Bit 0 – Hot Reset</p> <p>Bit 0 = 0b, De-assert</p> <p>Bit 0 = 1b, Assert</p> <p><u>Bit 1 – Disable Link</u></p> <p>Bit 1 = 0b, De-assert</p> <p>Bit 1 = 1b, Assert</p> <p><u>Bit 2 – Loopback</u></p> <p>Bit 2 = 0b, De-assert</p> <p>Bit 2 = 1b, Assert</p> <p><u>Bit 3 – Disable Scrambling</u></p> <p>Bit 3 = 0b, De-assert</p> <p>Bit 3 = 1b, Assert</p> <p><u>Bit 4:7 – Reserved</u></p> <p>Set to 0b</p>
6 – 15	D5.2	TS2 Identifier

1. 每 1538 个 symbol 中要有 1538 个 symbol order.
2. 1280 个 128 symbol 由 TS1 & TS2. 变化为  
在 Electrical IDLE.
3. 划分有线性翻转. 非 TS1 & TS2 中的 symbol D6-15.  
是 D10.2 与 D21.5 (D10.2 的反转) 或 D5.2 与 D26.5
4. FTSL (fast training) 因子 RX 方向接收



1. training 的时候，每条 lane 都以 20GT/s 的速率输出
  2. pCF-E link 包含一个或多个 lane
  3. TS2 中的连接和 TS1 中不同

Table 4-7: Link Status Mapped to the LTSSM

LTSSM State	Link Width	Link Speed	LinkUp	Link Training	Receiver Error	In-Band Presence <sup>26</sup>
Detect	Undefined	Undefined	0b	0b	No action	0b
Polling	Undefined	Set to 2.5 GT/s on entry from Detect. Link speed may change on entry to Polling Compliance.	0b	0b	No action	1b
Configuration	Set	No action	0b/1b <sup>27</sup>	1b	Set on 8b/10b Error	1b
Recovery	No action	Set to new speed when speed changes	1b	1b	No action	1b
L0	No action	No action	1b	0b	Set on 8b/10b Error or optionally on Framing Violation, Loss of Symbol Lock, Lane Desync Error, or Elasticity Buffer Overflow/Underflow	1b

<sup>26</sup> In-band refers to the fact that no sideband signals are used to calculate the presence of a powered up device on the other end of a Link.

<sup>27</sup> LinkUp will always be 0 if coming into Configuration via Detect -> Polling -> Configuration and LinkUp will always be 1 if coming into Configuration from any other state.

LTSSM State	Link Width	Link Speed	LinkUp	Link Training	Receiver Error	In-Band Presence <sup>26</sup>
L0s	No action	No action	1b	0b	No action	1b
L1	No action	No action	1b	0b	No action	1b
L2	No action	No action	1b	0b	No action	1b
Disabled	Undefined		0b	0b	Optional: Set on 8b/10b Error	1b
Loopback	No action	Link speed may change on entry to Loopback from Configuration.	0b	0b	No action	1b
Hot Reset	No action	No action	0b	0b	Optional: Set on 8b/10b Error	1b

The state machine rules for configuring and operating a PCI Express Link are defined in the following sections.

## clock compensation

TX:

1. 随着 Lane [0] 时发送数据.
2. skip order sets : com + 3 个连续 & skip.
3. [1180-1528] 加入 skip

RX:

1. 可以通过 skip 信号. com + 1-5 skip.  
当 Lane 0 为 1528 时 Lane 4 为 1180 时忽略 skip.
2. 1180-1528 为 [1] 间隙
3. tolerance 为 ± 3 max\_payload\_size.

full swing (mandatory)

de-emphasis

low swing

power sensitive - shorter channel

TX 需支持 voltage margining, per Com2E 2012-02-28  
须支持 per Lane 6 operational

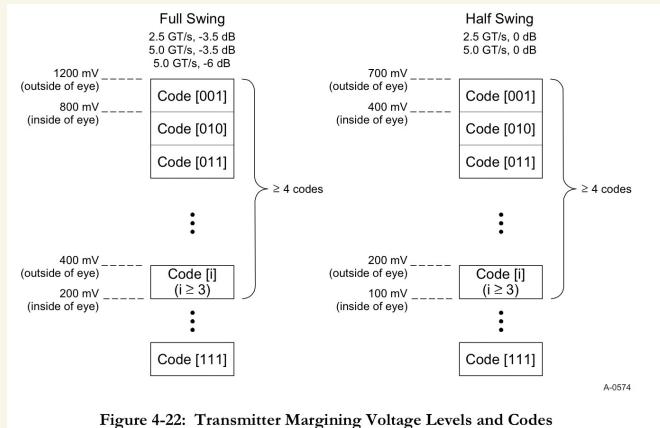


Figure 4-22: Transmitter Margining Voltage Levels and Codes

pin: 引脚. pad: 通常无法直接测量. 是由江炜提供 pm: 版图

将 30 kHz 以下信号注入到 DL. 通过 30Hz 时钟采样

TX:

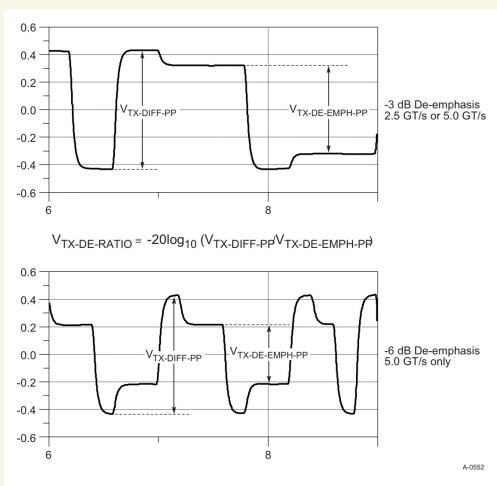


Figure 4-30: Full Swing Tx Parameters Showing De-emphasis

TX de-emphasis

10K  
并联,

$$\frac{V_{TX-DIFF-PP}}{-20 \log_{10} \sqrt{V_{TX-DE-EMPH-PP}}}$$

RX

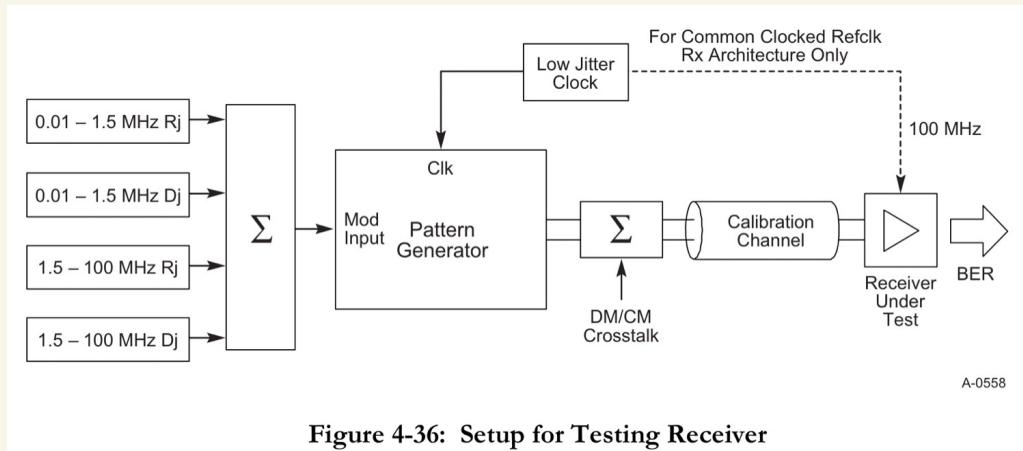


Figure 4-36: Setup for Testing Receiver

RX側通過發送EVE symbols來退出electrical idle.

PCI-E Link里所有的Lane都是AL耦合的。

BEACON 是在PCI-E从L2回傳正確的  
到Bridge & switch to downstream  
並且由Bridge生於 upstream

# PLL 架構

1. common refclk

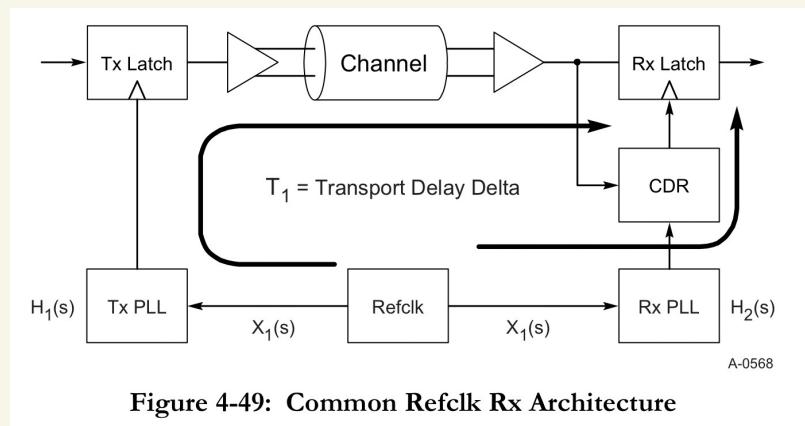


Figure 4-49: Common Refclk Rx Architecture

2. data clocked

又稱：兩埠[並行] CDR 'Tx' 出來的時鐘

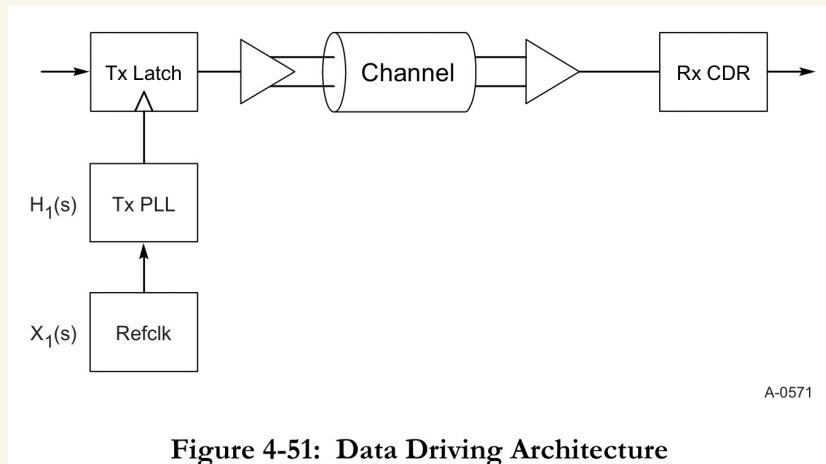


Figure 4-51: Data Driving Architecture

3. separate refclk , Tx & Rx 同不同，refclk

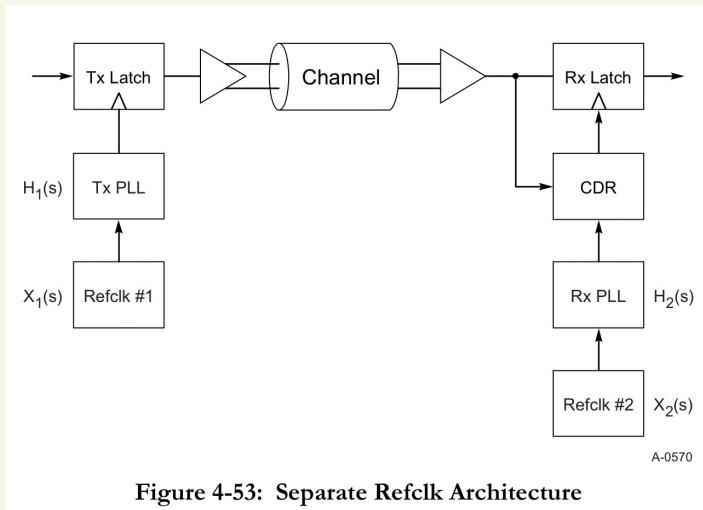
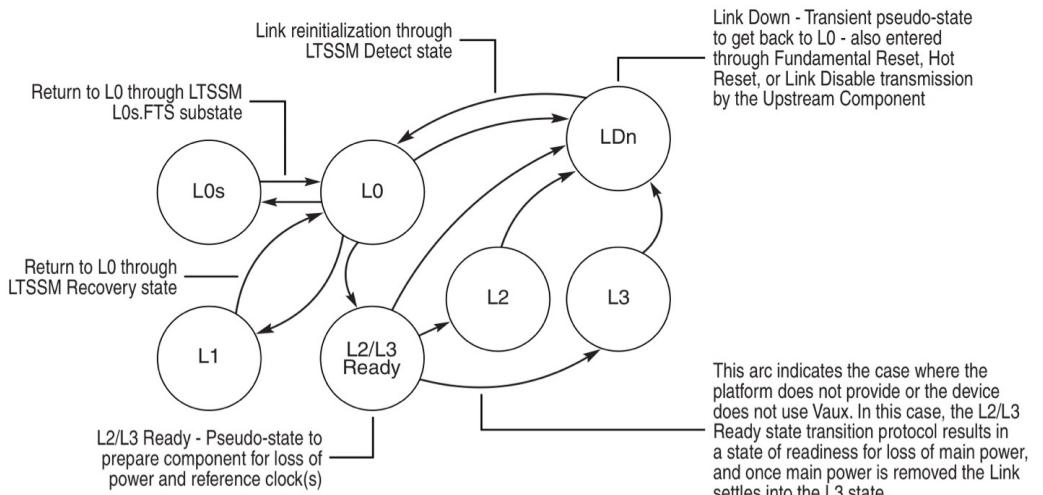


Figure 4-53: Separate Refclk Architecture

PCI Express-PM provides the following services:

- ❑ A mechanism to identify power management capabilities of a given Function
- ❑ The ability to transition a Function into a certain power management state
- ❑ Notification of the current power management state of a Function
- ❑ The option to wakeup the system on a specific event



OM13819B

Figure 5-1: Link Power Management State Flow Diagram

1. 软件中可手动设置 PCI-E power state, 通过 configure write .
2. 同时 PMT (LVDP) 中有 支持 power

# 中斷

## 1. INTx 模式

→ 向後兼容但沒溫

## 2. MSI & MSI-X

PCI-E 沒有必須支持，實現此功能是通過 memory write transaction. 同時其為 edge-triggered.

If the Root Port is enabled for edge-triggered interrupt signaling using MSI or MSI-X, an interrupt message must be sent every time the logical AND of the following conditions transitions from FALSE to TRUE:

- The associated vector is unmasked (not applicable if MSI does not support PVM).
- The PME Interrupt Enable bit in the Root Control register is set to 1b.
- The PME Status bit in the Root Status register is set.

Note that PME and Hot-Plug Event interrupts (when both are implemented) always share the same MSI or MSI-X vector, as indicated by the Interrupt Message Number field in the PCI Express Capabilities register.

# 錯誤報告和日志

Baseline 和 AER

1. 錯誤分為 fatal, non-fatal, Uncorrectable, Correctable

2. Correctable errors

由 software 處理

3. Uncorrectable errors

{ fatal → 需要重置  
    no-fatal → 需要恢復

# 错误信号实现

1. completion status  
通过 unsupported request complete status
2. error message - 一种 TLP message  
request

Table 6-1: Error Messages

Error Message	Description
ERR_COR	This Message is issued when the Function or Device detects a correctable error on the PCI Express interface. Refer to Section 6.2.2.1 for the definition of a correctable error.
ERR_NONFATAL	This Message is issued when the Function or Device detects a Non-fatal, uncorrectable error on the PCI Express interface. Refer to Section 6.2.2.2 for the definition of a Non-fatal, uncorrectable error.
ERR_FATAL	This Message is issued when the Function or Device detects a Fatal, uncorrectable error on the PCI Express interface. Refer to Section 6.2.2.2.1 for the definition of a Fatal, uncorrectable error.

root complext 通过 message header  
根据 Request ID 判别出错误。  
error pollution

For errors detected in the Transaction layer, it is permitted and recommended that no more than one error be reported for a single received TLP, and that the following precedence (from highest to lowest) be used:

- Receiver Overflow
- Flow Control Protocol Error
- ECRC Check Failed
- Malformed TLP
- Unsupported Request (UR), Completer Abort (CA), or Unexpected Completion<sup>44</sup>
- Poisoned TLP Received

## Error Trapping

若 function 有误，会将 device 之 fit 取出  
function 由 registers 中读取。

## Interrupt Generation

1 - level-triggered interrupt → Intx message

2 - edge-triggered interrupt →  
Intx & Intx message.

(由) vector 地址 → ATR interrupt message.

**Table 6-2: Physical Layer Error List**

Error Name	Error Type	Detecting Agent Action <sup>54</sup>	References
Receiver Error	Correctable	<i>Receiver:</i> Send ERR_COR to Root Complex.	Section 4.2.1.3 Section 4.2.2 Section 4.2.4.6 Section 4.2.6

**Table 6-3: Data Link Layer Error List**

Error Name	Error Type (Default Severity)	Detecting Agent Action <sup>54</sup>	References
Bad TLP	Correctable	<i>Receiver:</i> Send ERR_COR to Root Complex.	Section 3.5.3.1
Bad DLLP		<i>Receiver:</i> Send ERR_COR to Root Complex.	Section 3.5.2.1
Replay Timeout		<i>Transmitter:</i> Send ERR_COR to Root Complex.	Section 3.5.2.1
REPLAY NUM Rollover		<i>Transmitter:</i> Send ERR_COR to Root Complex.	Section 3.5.2.1
Data Link Layer Protocol Error	Uncorrectable (Fatal)	If checking, send ERR_FATAL to Root Complex.	Section 3.5.2.1

Error Name	Error Type (Default Severity)	Detecting Agent Action <sup>54</sup>	References
Surprise Down		If checking, send ERR_FATAL to Root Complex.	Section 3.5.2.1

Table 6-4: Transaction Layer Error List

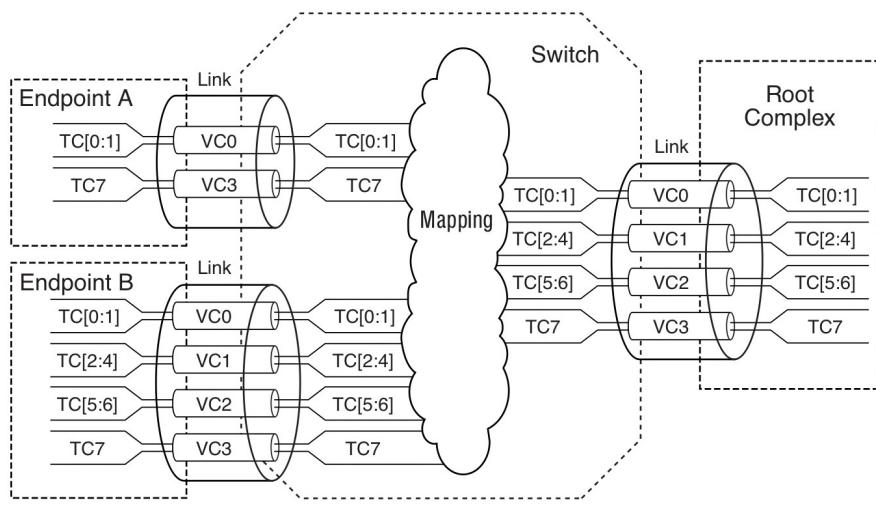
Error Name	Error Type (Default Severity)	Detecting Agent Action <sup>54</sup>	References
Poisoned TLP Received	Uncorrectable (Non-Fatal)	<p><i>Receiver:</i></p> <p>Send ERR_NONFATAL to Root Complex or ERR_COR for the Advisory Non-Fatal Error cases described in Sections 6.2.3.2.4.2 and 6.2.3.2.4.3.</p> <p>Log the header of the Poisoned TLP.<sup>55</sup></p>	Section 2.7.2.2
ECRC Check Failed		<p><i>Receiver (if ECRC checking is supported):</i></p> <p>Send ERR_NONFATAL to Root Complex or ERR_COR for the Advisory Non-Fatal Error case described in Section 6.2.3.2.4.2.</p> <p>Log the header of the TLP that encountered the ECRC error.</p>	Section 2.7.1
Unsupported Request (UR)		<p><i>Request Receiver:</i></p> <p>Send ERR_NONFATAL to Root Complex or ERR_COR for the Advisory Non-Fatal Error case described in Section 6.2.3.2.4.1.</p> <p>Log the header of the TLP that caused the error.</p>	Section 2.2.8.6, Section 2.3.1, Section 2.3.2, Section 2.7.2.2, Section 2.9.1, Section 5.3.1, Section 6.2.3.1, Section 6.2.6, Section 6.2.8.1, Section 6.5.7, Section 7.3.1, Section 7.3.3, Section 7.5.1.1, Section 7.5.1.2

Error Name	Error Type (Default Severity)	Detecting Agent Action <sup>54</sup>	References
Completion Timeout		<i>Requester:</i> Send ERR_NONFATAL to Root Complex or ERR_COR for the Advisory Non-Fatal Error case described in Section 6.2.3.2.4.4.	Section 2.8
Completer Abort		<i>Completer:</i> Send ERR_NONFATAL to Root Complex or ERR_COR for the Advisory Non-Fatal Error case described in Section 6.2.3.2.4.1. Log the header of the Request that encountered the error.	Section 2.3.1
Unexpected Completion		<i>Receiver:</i> Send ERR_COR to Root Complex. This is an Advisory Non-Fatal Error case described in Section 6.2.3.2.4.5. Log the header of the Completion that encountered the error.	Section 2.3.2
ACS Violation		<i>Receiver (if checking):</i> Send ERR_NONFATAL to Root Complex. Log the header of the Request TLP that encountered the error.	
Receiver Overflow		<i>Receiver (if checking):</i> Send ERR_FATAL to Root Complex.	Section 2.6.1.2
Flow Control Protocol Error		<i>Receiver (if checking):</i> Send ERR_FATAL to Root Complex.	Section 2.6.1

Error Name	Error Type (Default Severity)	Detecting Agent Action <sup>54</sup>	References
Malformed TLP		<i>Receiver:</i> Send ERR_FATAL to Root Complex. Log the header of the TLP that encountered the error.	Section 2.2.2, Section 2.2.3, Section 2.2.5, Section 2.2.7, Section 2.2.8.1, Section 2.2.8.2, Section 2.2.8.3, Section 2.2.8.4, Section 2.2.8.5, Section 2.2.9, Section 2.3, Section 2.3.1, Section 2.3.1.1, Section 2.3.2, Section 2.5, Section 2.5.3, Section 2.6.1, Section 2.6.1.2, Section 6.3.2

For all errors listed above, the appropriate status bit(s) must be set upon detection of the error. For Unsupported Request (UR), additional detection and reporting enable bits apply (see Section 6.2.5).

# VC support — vc & port & function-based



OM13828

Figure 6-4: TC Filtering Example

TLP 在传输过程中不具改变 TC 值

转发

1. TLP 中 address & router 信息决定 egress port.
2. TC & VC mapping 由 TLP 中的 TC 值决定

同一出口不同源到 traffic 竞争如何

1. flow control 是否满足  
ordering rule (不同 TLP 类型：前序保证)

2.

VC Resource	VC ID	Relative Priority	VC Arbitration	Usage Example
Extended VC Count = 7				
8th VC	VC 7	High		
7th VC	VC 6			
6th VC	VC 5			
5th VC	VC 4			
4th VC	VC 3			
3rd VC	VC 2			
2nd VC	VC 1			
1st VC	VC 0	Low		
		Priority Order	Strict Priority	For isochronous traffic
			Governed by VC Arbitration Capability field (e.g. WRR)	For QoS usage
			Low Priority Extended VC Count = 3	Default VC (PCI Express/PCI)

OM14287

Figure 6-9: VC ID and Priority Order – An Example

VC [向] 競争:

- 1. Strict priority - 平均: 优先级排序
- 2. Round Robin - 轮循.
- 3. Weighted RR - 带 weight: 权重.

↑ function ↑  
高优先级 VC 高带宽  
MFVL 和划

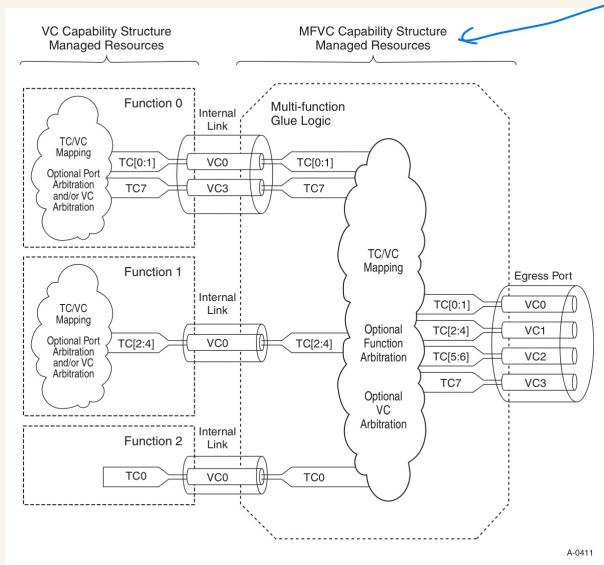


Figure 6-10: Multi-Function Arbitration Model

A-0411

## lock transaction

1. 为了解决冲突的purge问题.
2. 只有运行中的lock access.
3. lock后，其他向PL ↔ legacy PCI桥接的traffic被阻塞.
4. TUP类型是MRDUK & GCDUK.
5. 停止向地址RL & TL.

Conventional reset. — cold, warm, cold.

1. port 离线 & 外接串行复位.
  2. unreset tx phy层 bring up link  
Tx link training. 同时复位 VCO
  3. - 一些设置需要额外的 delay
- 除此之外，还有以功能复位 functional reset ← 由软件引起

# Hot plug

Table 6-5: Elements of Hot-Plug

Element	Purpose
Indicators	Show the power and attention state of the slot
Manually-operated Retention Latch (MRL)	Holds adapter in place
MRL Sensor	Allows the Port and system software to detect the MRL being opened
Electromechanical Interlock	Prevents removal of adapter from slot
Attention Button	Allows user to request hot-plug operations
Software User Interface	Allows user to request hot-plug operations
Slot Numbering	Provides visual identification of slots
Power Controller	Software-controlled electronic component or components that control power to a slot or adapter and monitor that power for fault conditions

Hot plug 相关的寄存器存在下述 downstream port 宝宝  
Device capabilities, slot capabilities, slot control,  
(slot status) ↗

Speed management.  
板卡插入和拔出时自动触发并触发  
retrain

# ACS & PCI-E 60 ACL

- ACS Source Validation (V)
- ACS Translation Blocking (B)
- ACS P2P Request Redirect (R)
- ACS P2P Completion Redirect (C)
- ACS Upstream Forwarding (U)

- ACS P2P Egress Control (E)
- ACS Direct Translated P2P (T)

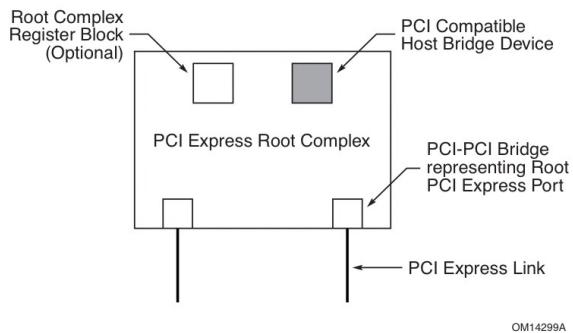


Figure 7-1: PCI Express Root Complex Device Mapping

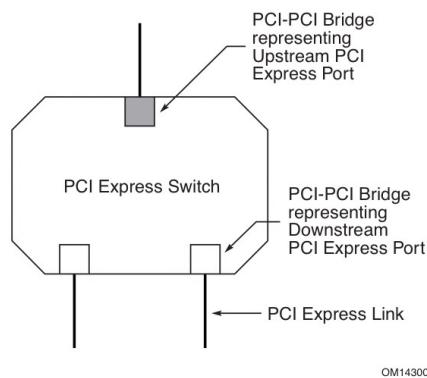


Figure 7-2: PCI Express Switch Device Mapping <sup>67</sup>

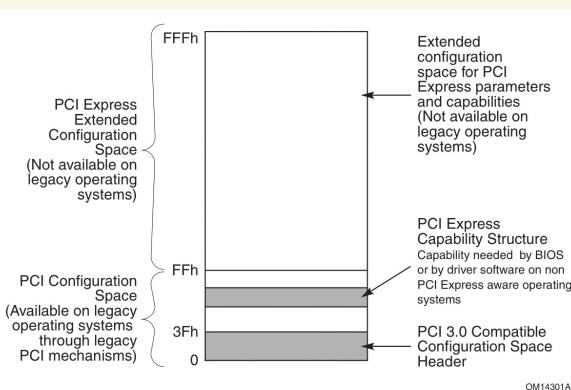


Figure 7-3: PCI Express Configuration Space Layout

RC 和 switch 間部分區  
(父型儲存) PCI.

1. PCI-Ex Link 由 RC 與 PCI bridge to second bus

2. 只有 RC, switch 才有 PCI-bridge 機制, EP 沒有

1. 前 256 字節為舊版 PCI
2. 0x100 - 0xFFFF 為 PCI-Ex 扩展 config space.
3. BDF 有 2<sup>16</sup> 個設備. 每個設備有 2<sup>12</sup> 大小的空間. 2<sup>16</sup> 個指標有 2<sup>28</sup> BYTES. 約 256 MB.

Table 7-1: Enhanced Configuration Address Mapping

Memory Address <sup>70</sup>	PCI Express Configuration Space
A[(20 + n - 1):20]	Bus Number $1 \leq n \leq 8$
A[19:15]	Device Number
A[14:12]	Function Number
A[11:8]	Extended Register Number
A[7:2]	Register Number
A[1:0]	Along with size of the access, used to generate Byte Enables

memory 地址 config 空间

1. A[63:120+n] 指定  
memory base address

2. 其他域如图所示

3. 将于 posted-write (通常为 memory write). Software  
必须通过 read 同一个地址来保证写结果4. configure write 和 memory read 时注意排序。  
所以需要有一种机制来保证芯片 config unit  
配置 BASE BAR memory ADDR. 然后是 read BAR  
范围地址 ← 需排序导致读不到内容因为在 vendor-specific 中的寄存器是自由的。  
memory-map allocate 前需完成这些配置。

- PCI-E 设备在 upstream port。由 PCI-E → Device number, 由 PCI-E → function number.
  - 从 downstream & root port 以 F 为 bus ID - [2] 为 device-number 为 0 (P2 为 PCI)
  - downstream ports 会分配给 device。
  - configuration-request Address:  
bus + dev + function + extended & normal register number.

Byte Offset	31	
00h	Device ID	
04h	Status	
08h	Class Code	
0Ch	Revision ID	
10h	BIST	
14h	Header Type	
18h	Master Latency Timer	
1Ch	Cache Line Size	
20h	Header Type Specific	
24h		
28h		
2Ch		
30h		
34h	Capabilities Pointer	
38h		
3Ch	Interrupt Pin      Interrupt Line	

B. Register and Field Types		
Register Attribute	Register Attribute	Description
<b>Hardware Initialization</b> – Registers are initialized by hardware mechanisms such as a掉电复位 or system boot.	<b>Sticky - Write-1-to-clear status</b> – Registers indicate status when read. A set bit indicates a status event which is Cleared by writing a 1b. Writing a 0b to RWIC/S bits has no effect.	
<b>Reset</b> – Registers are initialized by a硬复位 or power cycle.	<b>RWIC/S</b>	Bits are neither initialized nor modified by hot reset or FLR.
<b>RO</b> – Registers are read-only and cannot be modified by software.		Where noted, devices that consume AUX power must preserve sticky register values when ALU power consumption (via either AUX power or PME Enable) is enabled. In these cases, registers are neither initialized nor modified by hot, warm, or cold reset (see Section 6.6).
<b>RW</b> – Registers are read/write and are permitted to be modified by software.	<b>Write-1-to-clear status</b> – Registers indicate status when read. An Set bit indicates a status event which is Cleared by writing a 1b. Writing a 0b to RWIC/S bits has no effect.	
<b>RWIC</b>		
<b>Sticky - Read-only</b> – Registers are read-only and cannot be modified by software.	<b>Sticky - Read-only</b> – Registers are read-only and cannot be modified by software.	
<b>RWS</b>	<b>Reserved and Preserved</b> – Reserved for future RW implementations. Registers are read-only and must return zero when read. Software must preserve the value read for writes to bits.	
<b>RW</b>	<b>RsvdZ</b>	
<b>RWIS</b>	<b>Reserved and Zero</b> – Reserved for future RWIC implementations. Registers are read-only and must return zero when read. Software must use 0b for writes to bits.	

Figure 7-4: Common Configuration Space Header

## 0x04 中关锁：控制 (control)

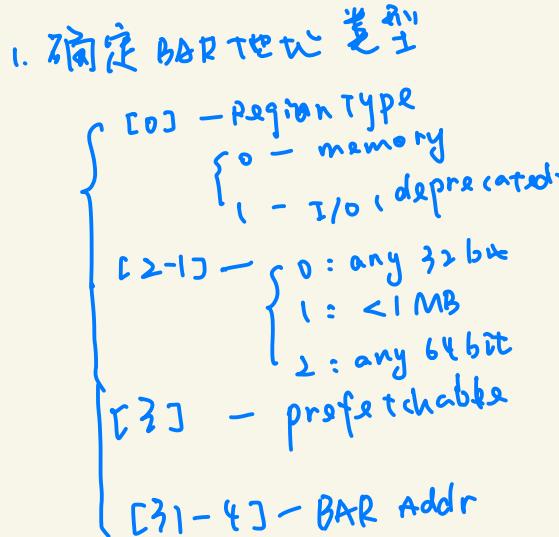
1. Bus master, 是否可以发起 I/O & mem 请求
2. parity error, 磁盘校验符错误报告
3. SERR enable, 系统错误报告
4. int disable, 是否关闭 Intx 中断  
(通过 message code [3f])

## 0x06 关键字锁：状态 (status)

1. Int status - 是否有 pending 中断
2. parity error - 是否有 parity error
3. INTB[5:4] 和 42# error.

Device ID		Vendor ID		Byte Offset 0 00h 04h 08h 0Ch 10h 14h 18h 1Ch 20h 24h 28h 2Ch 30h 34h 38h 3Ch	
Status		Command			
Class Code			Revision ID		
BIST	Header Type	Master Latency Timer	Cache Line Size		
<b>3个BAR地址，前两个BAR128Bit</b> <b>64bit</b>					
<b>6个BAR地址</b> <b>~32bit</b>					
Base Address Registers					
Cardbus CIS Pointer					
Subsystem ID	Subsystem Vendor ID				
Expansion ROM Base Address					
Reserved		Capabilities Pointer			
Reserved					
Max_Lat	Min_Gnt	Interrupt Pin	Interrupt Line		

Figure 7-5: Type 0 Configuration Space Header



2. 确定 BAR 地址大小  
向 BAR 寄存器写入值返回  
n口. 得到大小

Device ID		Vendor ID		Byte Offset 00h 04h 08h 0Ch 10h 14h 18h 20h 24h 28h 2Ch 30h 34h 38h 3Ch				
Status		Command						
Class Code			Revision ID					
BIST	Header Type	Primary Latency Timer	Cache Line Size					
<b>指向地址0</b>								
Base Address Register 0								
Base Address Register 1								
Secondary Latency Timer	Subordinate Bus Number	Secondary Bus Number	Primary Bus Number					
Secondary Status		I/O Limit	I/O Base					
Memory Limit		Memory Base						
Prefetchable Memory Limit		Prefetchable Memory Base						
Prefetchable Base Upper 32 Bits								
Prefetchable Limit Upper 32 Bits								
I/O Limit Upper 16 Bits		I/O Base Upper 16 Bits						
Reserved			Capability Pointer					
Expansion ROM Base Address								
Bridge Control	Interrupt Pin	Interrupt Line						

Figure 7-6: Type 1 Configuration Space Header

# Capabilities

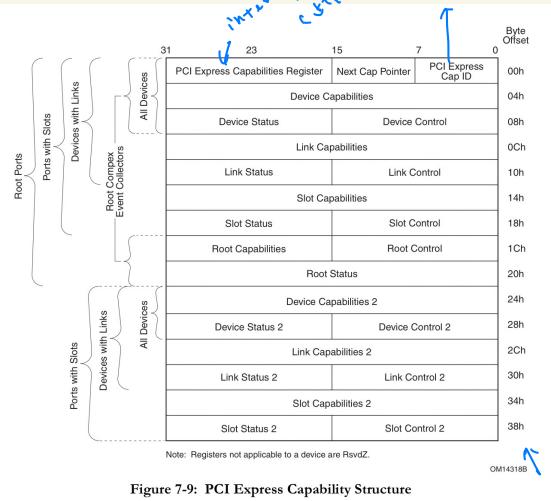


Figure 7-9: PCI Express Capability Structure

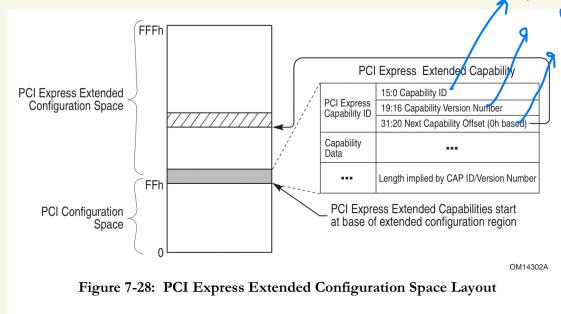


Figure 7-28: PCI Express Extended Configuration Space Layout

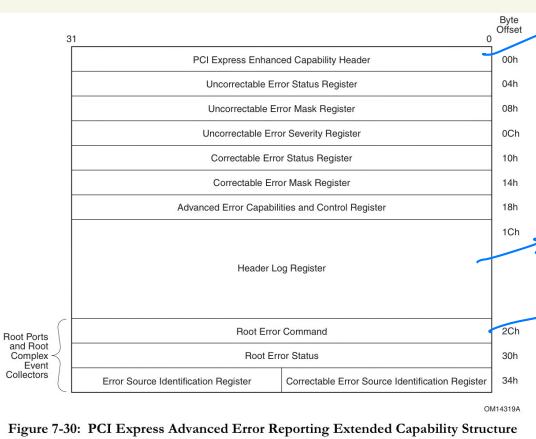


Figure 7-30: PCI Express Advanced Error Reporting Extended Capability Structure



