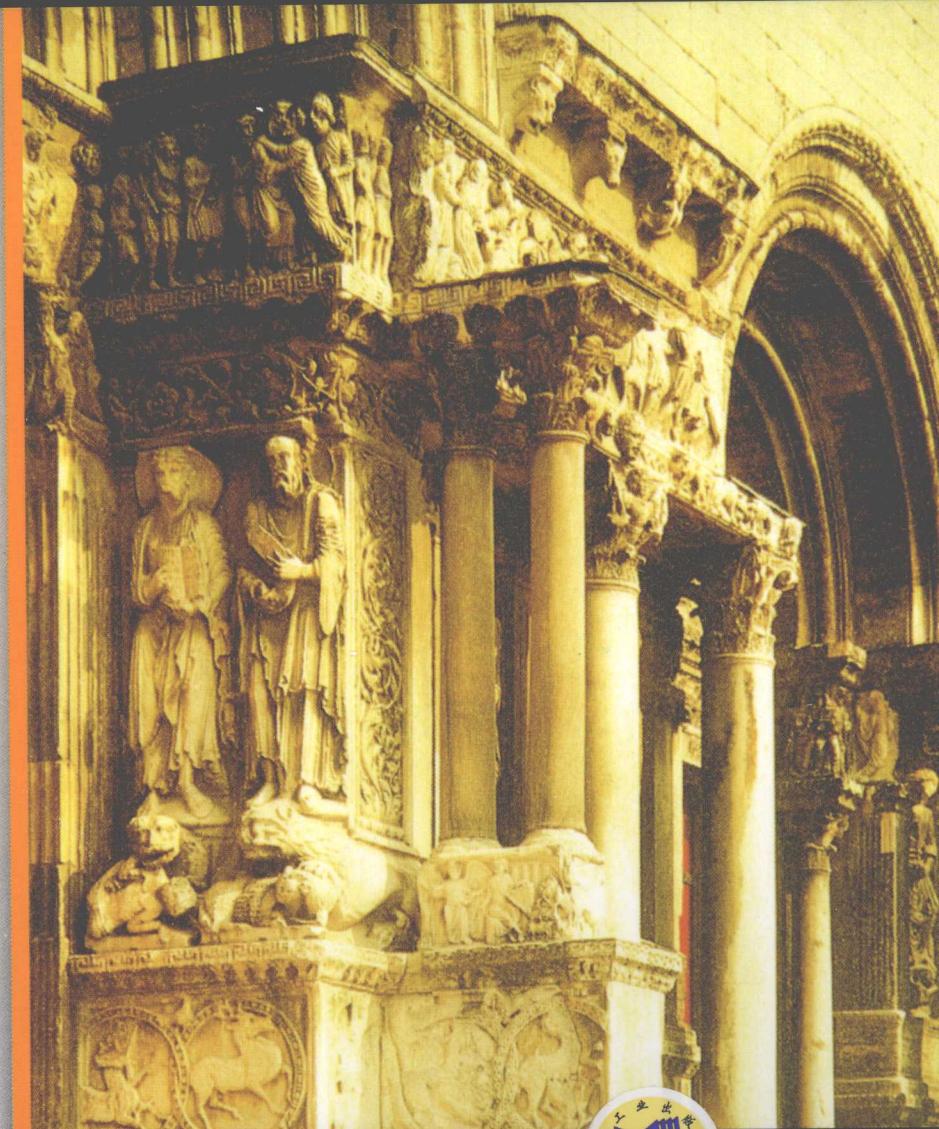


# PCI Express

## 体系结构导读

■ 王齐 编著



# PCI Express 体系结构导读

王齐 编著



机 械 工 业 出 版 社

本书讲述了与 PCI 及 PCI Express 总线相关的最为基础的内容，并介绍了一些必要的、与 PCI 总线相关的处理器体系结构知识，这也是本书的重点所在。深入理解处理器体系结构是理解 PCI 与 PCI Express 总线的重要基础。

读者通过对本书的学习，可超越 PCI 与 PCI Express 总线自身的内容，理解在一个通用处理器系统中局部总线的设计思路与实现方法，从而理解其他处理器系统使用的局部总线。本书适用于希望多了解一些硬件的软件工程师，以及希望多了解一些软件的硬件工程师，也可供电子工程和计算机类的研究生自学参考。

### 图书在版编目(CIP)数据

PCI Express 体系结构导读 / 王齐编著. —北京 : 机械工业出版社, 2010. 3  
ISBN 978-7-111-29822-9

I. ①P… II. ①王… III. ②总线—结构 IV. ③TP336

中国版本图书馆 CIP 数据核字(2010)第 028735 号

机械工业出版社（北京市百万庄大街 22 号 邮政编码 100037）

责任编辑：车 忱

责任印制：洪汉军

三河市宏达印刷有限公司印刷

2010 年 3 月第 1 版 · 第 1 次印刷

184mm × 260mm · 28.5 印张 · 704 千字

0001 - 3500 册

标准书号：ISBN 978-7-111-29822-9

定价：55.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

电话服务 网络服务

社服务中心：(010) 88361066

销售一部：(010) 68326294

销售二部：(010) 88379649

读者服务部：(010) 68993821

门户网：<http://www.cmpbook.com>

教材网：<http://www.cmpedu.com>

封面无防伪标均为盗版

# 序

PCI Express 总线是新一代的 I/O 局部总线标准，是取代 PCI 总线的革命性总线架构。PCI 总线曾经是 PC 体系结构发展史上的一个里程碑，但是随着技术的不断发展，新涌现出的一些外部设备对传输速度和带宽有更高的要求，如千兆和万兆以太网、4Gb/8Gb 的 Fiber Channel 和高速显示设备等。同时有些外部设备对总线的服务质量还有更严格的要求。PCI 总线在设计之初并没有考虑这些因素，因此并不能完全满足这些外部设备的需要。

PCI Express 总线正是在这种背景下应运而生的。在 2001 年的春季英特尔开发者论坛上，英特尔公布了取代 PCI 总线的第三代 I/O 技术，当时被称为“3GIO”。经 PCI-SIG 审核，于 2002 年 7 月正式公布了第一版规范，并更名为 PCI Express。从 2004 年开始，PCI Express 总线逐渐全面取代 PCI 和 AGP 总线，成为新的局部总线工业标准。

与 PCI 总线的共享并行架构不同，PCI Express 总线使用高速串行传送方式，能够支持更高的频率，连接的设备不再像 PCI 总线那样共享总线带宽。除此之外 PCI Express 总线还引入了一些新特性，如流量控制机制、服务质量管理、热插拔支持、数据完整性和新型错误处理机制等。而且 PCI Express 总线在系统软件级与 PCI 总线保持兼容，最大程度上降低了系统软件从原有的 PCI 总线体系结构移植到 PCI Express 总线体系结构的难度。

目前关于 PCI Express 总线规范的文献和书籍已有多种，但多集中在介绍规范本身。对于广大的开发者来说，能够从处理器系统的角度了解 PCI Express 总线功能，无疑更有实用价值。无论是系统外部设备的开发、驱动程序的编写，还是其他系统软件的开发，处理器系统始终处于核心位置。

本书正是从处理器系统的视角来讲述 PCI Express 总线的体系结构，较为细致地介绍了如何使用 FPGA 实现一个基于 PCIe 总线的外部设备，以及基于 Linux 系统的 PCL/PCI Express 总线驱动程序和设备驱动程序。本书对于 PCI Express 总线相关的软硬件开发人员具有很高的参考价值。

PCI Express 总线规范仍在不断发展。总的来说，PCI Express 总线规范提出的最新技术概念往往在英特尔的 x86 处理器系统和外部设备中最先出现。本书的作者王齐先生，目前工作于英特尔开源技术研究中心，对处理器体系结构和 Linux 系统核心技术均有深入研究，相信本书的读者能够从他的经验分享中获益。

杨继国  
英特尔开源技术中心

# 前　　言

PCI 与 PCI Express (PCIe) 总线在处理器系统中得到了大规模应用。PCISIG 也制定了一系列 PCI 与 PCI Express 总线相关的规范，这些规范所涉及的内容庞杂广泛。对于已经理解了 PCI 与 PCI Express 总线的工程师，这些规范便于他们进一步获得必要的细节知识。对于刚刚接触 PCI 与 PCI Express 总线的工程师，这些规范性的文档并不适合阅读。在阅读这些规范时，工程师还需要具备一些与体系结构相关的基础知识，这恰是规范并不涉及的内容。对于多数工程师，规范文档适于查阅，而不便于学习。

本书将以处理器体系结构为主线介绍 PCI Express 总线的组成，以便读者进一步理解 PCI Express 总线协议。本书并不是关于 PCI 和 PCI Express 总线的百科全书，因为读者完全可以通过阅读 PCI 和 PCI Express 总线规范获得细节信息。本书侧重的是 PCI 和 PCI Express 总线中与处理器体系结构相关的内容。

本书不会对 PCI 总线的相关规范进行简单重复，部分内容并不在 PCI 总线规范定义的范围内，例如 HOST 主桥和 RC。PCI 总线规范并没有规定处理器厂商如何实现 HOST 主桥和 RC，不同的处理器厂商实现的 HOST 主桥和 RC 有较大差异，而这些内容正是本书所讨论的重点。此外本书还讲述了一些在 PCI 总线规范中提及，但是容易被忽略的一些重要概念。

本书共由三篇组成。第 I 篇（第 1~3 章）介绍 PCI 总线的基础知识。第 II 篇（第 4~13 章）介绍 PCI Express 总线的相关概念。第 II 篇的内容以第 I 篇为基础。从系统软件的角度来看，PCI Express 总线向前兼容 PCI 总线，理解 PCI Express 总线必须建立在深刻理解 PCI 总线的基础之上。读者需要按照顺序阅读这两篇。

第 1 章主要说明 PCI 总线涉及的一些基本知识。有些知识稍显过时，但是在 PCI 总线中出现的一些数据传送方式，如 Posted、Non-Posted 和 Split 数据传送方式，依然非常重要，也是读者需要掌握的。

第 2 章重点介绍 PCI 桥。PCI 桥是 PCI 及 PCI Express 体系结构的精华所在，本章还使用了一定篇幅介绍了非透明桥。非透明桥不是 PCI 总线定义的标准桥片，但是在处理器系统之间的互联中得到了广泛的应用。

第 3 章详细阐述 PCI 总线的数据传送方式，与 Cache 相关的内容和预读机制是本章的重点。目前 PCI 与 PCI Express 对预读机制的支持并不理想。但是在可以预见的将来，PCI Express 总线将充分使用智能预读机制进一步提高总线的利用率。

第 4 章是 PCI Express 篇的综述。第 5 章以 Intel 的笔记本平台 Montevina 为例说明 RC 的各个组成模块。实际上 RC 这个概念，只有在 x86 处理器平台中才真正存在。其他处理器系统中，并不存在严格意义上的 RC。

第 6、7 章分别介绍 PCI Express 总线的事务层、数据链路层和物理层。物理层是 PCI Express 总线的真正核心，也是中国工程师最没有机会接触的内容。这也是我们这一代工程师的遗憾与无奈。第 8 章简要说明了 PCI Express 总线的链路训练与电源管理。

第 9 章主要讨论的是通用流量控制的管理方法与策略。PCI Express 总线的流量控制机

制仍需完善，其中不等长的报文长度也是限制 PCI Express 总线流量控制进一步提高的重要因素。

第 10 章重点介绍 MSI 和 MSI-X 中断机制。MSI 中断机制在 PCI 总线中率先提出，但是在 PCI Express 总线中也得到大规模普及。目前 x86 架构多使用 MSI-X 中断机制，而在许多嵌入式处理器中仍然使用 MSI 中断机制。

第 11 章的篇幅很短，重点介绍 PCI 和 PCI Express 总线中的序。有志于学习处理器体系结构的工程师务必掌握这部分内容。在处理器体系结构中有关 Cache 和数据传送序的内容非常复杂，掌握这些内容也是系统工程师进阶所必须的。

第 12 章讲述了笔者的一个实际设计——Capric 卡，简单介绍了 Linux 设备驱动程序的实现过程，并对 PCI Express 总线的延时与带宽进行了简要分析。

第 13 章介绍 PCI 总线与虚拟化相关的一些内容。虚拟化技术已崭露头角，与虚拟化相关的一系列内容将对处理器体系结构产生深远的影响。目前虚拟化技术已经在 x86 和 PowerPC 处理器中得到了广泛的应用。

第Ⅲ篇以 Linux 系统为实例说明 PCI 总线在处理器系统中的使用方法，也许有许多读者对这一篇有着浓厚的兴趣。Linux 无疑是一个非常优秀的操作系统。但是需要提醒系统工程师，Linux 系统仅是一个完全开源的操作系统。对于有志于学习处理器体系结构的工程师，学习 Linux 系统是必要的，但是仅靠学习 Linux 系统并不足够。

通常说来，理解处理器体系结构至少需要了解两三种处理器，并了解它们在不同处理器上的实现。尺有所短，寸有所长。不同的处理器和操作系统所应用的领域并不完全相同。也是因为这个原因，本书以 PowerPC 和 x86 处理器为基础对 PCI 和 PCI Express 总线进行说明。

本书在写作过程中得到了我的同事和在处理器及操作系统行业奋战多年的朋友们的帮助。在 Linux 系统中许多与处理器和 PCI 总线相关的模块，都有着他们的辛勤付出。刘建国和郭超宏先生审阅了本书的第 I 篇。马明辉先生审阅了本书的第 II 篇。张巍、余珂与刘劲松先生审阅了第 13 章。吴晓川、王勇、丁建峰、李力与吴强先生共同审阅了全书。

本书第 12 章中出现的 Capric 和 Cornus 卡由郭冠军和高健协助完成。看着他们通过对 PCI Express 总线理解的逐渐深入，最终设计出一个具有较高性能的 Cornus 卡，备感欣慰。此外杨强浩先生也参与了 Capric 和 Cornus 卡的原始设计与方案制定，在此对他及他的团队在这个过程中给予的帮助表示感谢，我们也一道通过这两块卡的制作进一步领略了 PCI Express 总线的技术之美。

一个优秀的协议，从制定到广大技术人员理解其精妙之处，再到协议应用到一个个优秀产品中，需要更多的人参与、投入、实践，这也是编写此书最大的动力源泉。本书的完成与我的妻子范淑琴的激励直接相关，Capricornus 也是她的星座。还需要感谢本书的编辑车忱与策划时静，正是他们的努力使得本书提前问世。

对本书尚留疑问的读者，可通过我的邮箱 sailing.w@gmail.com 与我联系。最后希望这本书对您有所帮助。

作 者

# 目 录

序  
前言

## 第 I 篇 PCI 体系结构概述

<b>第 1 章 PCI 总线的基本知识 .....</b>	<b>3</b>	<b>第 2 章 PCI 总线的桥与配置 .....</b>	<b>28</b>
1.1 PCI 总线的组成结构 .....	5	2.1 存储器域与 PCI 总线域 .....	28
1.1.1 HOST 主桥 .....	6	2.1.1 CPU 域、DRAM 域与存储器域 .....	29
1.1.2 PCI 总线 .....	7	2.1.2 PCI 总线域 .....	30
1.1.3 PCI 设备 .....	7	2.1.3 处理器域 .....	30
1.1.4 HOST 处理器 .....	8	2.2 HOST 主桥 .....	32
1.1.5 PCI 总线的负载 .....	8	2.2.1 PCI 设备配置空间的访问机制 .....	33
1.2 PCI 总线的信号定义 .....	9	2.2.2 存储器域地址空间到 PCI 总线 域地址空间的转换 .....	35
1.2.1 地址和数据信号 .....	9	2.2.3 PCI 总线域地址空间到存储器 域地址空间的转换 .....	37
1.2.2 接口控制信号 .....	10	2.2.4 x86 处理器的 HOST 主桥 .....	40
1.2.3 仲裁信号 .....	12	2.3 PCI 桥与 PCI 设备的配置空间 .....	42
1.2.4 中断请求等其他信号 .....	12	2.3.1 PCI 桥 .....	42
1.3 PCI 总线的存储器读写总线事务 .....	13	2.3.2 PCI Agent 设备的配置空间 .....	44
1.3.1 PCI 总线事务的时序 .....	14	2.3.3 PCI 桥的配置空间 .....	50
1.3.2 Posted 和 Non-Posted 传送方式 .....	15	2.4 PCI 总线的配置 .....	53
1.3.3 HOST 处理器访问 PCI 设备 .....	16	2.4.1 Type 01h 和 Type 00h 配置请求 .....	53
1.3.4 PCI 设备读写主存储器 .....	18	2.4.2 PCI 总线配置请求的转换原则 .....	55
1.3.5 Delayed 传送方式 .....	19	2.4.3 PCI 总线树 Bus 号的初始化 .....	57
1.4 PCI 总线的中断机制 .....	21	2.4.4 PCI 总线 Device 号的分配 .....	59
1.4.1 中断信号与中断控制器的连接 关系 .....	21	2.5 非透明 PCI 桥 .....	60
1.4.2 中断信号与 PCI 总线的连接 关系 .....	22	2.5.1 Intel 21555 中的配置寄存器 .....	62
1.4.3 中断请求的同步 .....	23	2.5.2 通过非透明桥片进行数据传递 .....	63
1.5 PCI-X 总线简介 .....	25	2.6 小结 .....	65
1.5.1 Split 总线事务 .....	25	第 3 章 PCI 总线的数据交换 .....	66
1.5.2 总线传送协议 .....	26	3.1 PCI 设备 BAR 空间的初始化 .....	66
1.5.3 基于数据块的突发传送 .....	26	3.1.1 存储器地址与 PCI 总线地址 的转换 .....	66
1.6 小结 .....	27		

3.1.2 PCI 设备 BAR 寄存器和 PCI 桥	81
Base、Limit 寄存器的初始化	67
3.2 PCI 设备的数据传递	69
3.2.1 PCI 设备的正向译码与负向	
译码	69
3.2.2 处理器到 PCI 设备的数据传送	71
3.2.3 PCI 设备的 DMA 操作	72
3.2.4 PCI 桥的 Combining、Merging	
和 Collapsing	73
3.3 与 Cache 相关的 PCI 总线事务	74
3.3.1 Cache 一致性的基本概念	75
3.3.2 PCI 设备对不可 Cache 的存储器	
空间进行 DMA 读写	80
3.3.3 PCI 设备对可 Cache 的存储器	
空间进行 DMA 读写	81
3.3.4 PCI 设备进行 DMA 写时发生	
Cache 命中	82
3.3.5 DMA 写时发生 Cache 命中	
的优化	85
3.4 预读机制	86
3.4.1 指令预读	86
3.4.2 数据预读	89
3.4.3 软件预读	91
3.4.4 硬件预读	93
3.4.5 PCI 总线的预读机制	94
3.5 小结	98

## 第 II 篇 PCI Express 体系结构概述

<b>第 4 章 PCIe 总线概述</b>	101
4.1 PCIe 总线的基础知识	101
4.1.1 端到端的数据传递	101
4.1.2 PCIe 总线使用的信号	103
4.1.3 PCIe 总线的层次结构	107
4.1.4 PCIe 链路的扩展	108
4.1.5 PCIe 设备的初始化	110
4.2 PCIe 体系结构的组成部件	112
4.2.1 基于 PCIe 架构的处理器系统	112
4.2.2 RC 的组成结构	117
4.2.3 Switch	118
4.2.4 VC 和端口仲裁	120
4.2.5 PCIe-to-PCL/PCI-X 桥片	122
4.3 PCIe 设备的扩展配置空间	123
4.3.1 Power Management Capability	
结构	124
4.3.2 PCI Express Capability 结构	127
4.3.3 PCI Express Extended Capabilities	
结构	133
4.4 小结	139
<b>第 5 章 Montevina 的 MCH 和 ICH</b>	140
5.1 PCI 总线 0 的 Device 0 设备	141
5.1.1 EPBAR 寄存器	144
5.1.2 MCHBAR 寄存器	144
5.1.3 其他寄存器	144
5.2 Montevina 平台的存储器空间的	
组成结构	145
5.2.1 Legacy 地址空间	147
5.2.2 DRAM 域	147
5.2.3 存储器域	148
5.3 存储器域的 PCI 总线地址	
空间	150
5.3.1 PCI 设备使用的地址空间	150
5.3.2 PCIe 总线的配置空间	151
5.4 小结	154
<b>第 6 章 PCIe 总线的事务层</b>	155
6.1 TLP 的格式	155
6.1.1 通用 TLP 头的 Fmt 字段和 Type	
字段	156
6.1.2 TC 字段	158
6.1.3 Attr 字段	159
6.1.4 通用 TLP 头中的其他字段	160
6.2 TLP 的路由	161
6.2.1 基于地址的路由	161

6.2.2	基于 ID 的路由	164	8.1.3	Receiver Detect 识别逻辑	217
6.2.3	隐式路由	166	8.2	LTSSM 状态机	218
<b>6.3</b>	<b>存储器、I/O 和配置读写</b>		8.2.1	Detect 状态	220
	请求 TLP	167	8.2.2	Polling 状态	221
6.3.1	存储器读写请求 TLP	168	8.2.3	Configuration 状态	223
6.3.2	完成报文	172	8.2.4	Recovery 状态	228
6.3.3	配置读写请求 TLP	174	8.2.5	LTSSM 的其他状态	231
6.3.4	消息请求报文	175	8.3	PCIe 总线的 ASPM	232
6.3.5	PCIe 总线的原子操作	177	8.3.1	与电源管理相关的链路状态	232
6.3.6	TLP Processing Hint	178	8.3.2	L0 状态	233
<b>6.4</b>	<b>TLP 中与数据负载相关的参数</b>		8.3.3	L0s 状态	234
	参数	181	8.3.4	L1 状态	235
6.4.1	Max_Payload_Size 参数	181	8.3.5	L2 状态	236
6.4.2	Max_Read_Request_Size 参数	182	8.4	PCI PM 机制	237
6.4.3	RCB 参数	183	8.4.1	PCIe 设备的 D-State	237
<b>6.5</b>	<b>小结</b>		8.4.2	D-State 的状态迁移	238
	184	8.5	小结		240
<b>第 7 章</b>	<b>PCIe 总线的数据链路层与物理层</b>		<b>第 9 章</b>	<b>流量控制</b>	241
	物理层	185	9.1	流量控制的基本原理	242
<b>7.1</b>	<b>数据链路层的组成结构</b>		9.1.1	Rate-Based 流量控制	243
	数据链路层的状态	186	9.1.2	Credit-Based 流量控制	244
7.1.1	DL_Down 和		9.2	Credit-Based 机制使用的算法	246
	DL_Up 状态	189	9.2.1	N123 算法和 N123 + 算法	249
7.1.3	DLLP 的格式	189	9.2.2	N23 算法	250
<b>7.2</b>	<b>ACK/NAK 协议</b>		9.2.3	流量控制机制的缓冲管理	252
	发送端如何使用 ACK/NAK		9.3	PCIe 总线的流量控制	254
7.2.1	协议	192	9.3.1	PCIe 总线流量控制的缓存	
7.2.2	接收端如何使用 ACK/NAK			管理	255
	协议	195	9.3.2	Current 节点的 Credit	257
7.2.3	数据链路层发送报文的顺序	199	9.3.3	VC 的初始化	259
<b>7.3</b>	<b>物理层简介</b>		9.3.4	PCIe 设备如何使用 FCP	261
	PCIe 链路的差分信号	200	9.4	小结	262
7.3.1	物理层的组成结构	204	<b>第 10 章</b>	<b>MSI 和 MSI-X 中断机制</b>	263
7.3.3	8/10b 编码与解码	206	10.1	MSI/MSI-X Capability 结构	263
<b>7.4</b>	<b>小结</b>		10.1.1	MSI Capability 结构	264
	210	10.1.2	MSI-X Capability 结构	266	
<b>第 8 章</b>	<b>PCIe 总线的链路训练与电源管理</b>		10.2	PowerPC 处理器如何处理	
	链路训练使用的字符序列	211		MSI 中断请求	268
<b>8.1</b>	<b>PCIe 链路训练简介</b>		10.2.1	MSI 中断机制使用的寄存器	270
	Electrical Idle 状态	216			

10.2.2 系统软件如何初始化 PCIe 设备的 MSI Capability 结构	272	12.1.4 DMA 读	302
10.3 x86 处理器如何处理 MSI-X 中断请求	274	12.1.5 中断请求	303
10.3.1 Message Address 字段和 Message Data 字段的格式	274	12.2 Capric 卡的数据传递	303
10.3.2 FSB Interrupt Message 总线事务	277	12.2.1 DMA 写使用的 TLP	304
10.4 小结	278	12.2.2 DMA 读使用的 TLP	308
<b>第 11 章 PCI/PCIe 总线的序</b>	279	12.2.3 Capric 卡的中断请求	317
11.1 生产/消费者模型	279	12.3 基于 PCIe 总线的设备驱动	317
11.1.1 生产/消费者的工作原理	280	12.3.1 Capric 卡驱动程序的加载与卸载	318
11.1.2 生产/消费者模型在 PCI/PCIe 总线中的实现	281	12.3.2 Capric 卡的初始化与关闭	319
11.2 PCI 总线的死锁	283	12.3.3 Capric 卡的 DMA 读写操作	324
11.2.1 缓冲管理引发的死锁	283	12.3.4 Capric 卡的中断处理	327
11.2.2 数据传送序引发的死锁	283	12.3.5 存储器地址到 PCI 总线地址的转换	328
11.3 PCI 总线的序	284	12.3.6 存储器与 Cache 的同步	330
11.3.1 PCI 总线序的通用规则	284	12.4 Capric 卡的延时与带宽	334
11.3.2 Delayed 总线事务的传送规则	285	12.4.1 TLP 的传送开销	335
11.3.3 PCI 总线事务通过 PCI 桥的顺序	286	12.4.2 PCIe 设备的 DMA 读写延时	338
11.3.4 LOCK, Delayed 和 Posted 总线事务间的关系	289	12.4.3 Capric 卡的优化	342
11.4 PCIe 总线的序	290	12.5 小结	343
11.4.1 TLP 传送的序	290	<b>第 13 章 PCIe 总线与虚拟化技术</b>	344
11.4.2 ID-Base Ordering	294	13.1 IOMMU	344
11.4.3 MSI 报文的序	295	13.1.1 IOMMU 的工作原理	345
11.5 小结	296	13.1.2 IA 处理器的 VT-d	347
<b>第 12 章 PCIe 总线的应用</b>	297	13.1.3 AMD 处理器的 IOMMU	349
12.1 Capric 卡的工作原理	297	13.2 ATS (Address Translation Services)	352
12.1.1 BAR 空间	298	13.2.1 TLP 的 AT 字段	353
12.1.2 Capric 卡的初始化	301	13.2.2 地址转换请求	354
12.1.3 DMA 写	302	13.2.3 Invalidate ATC	356

### 第 III 篇 Linux 与 PCI 总线

<b>第 14 章 Linux PCI 的初始化过程</b>	365	初始化	365
14.1 Linux x86 对 PCI 总线的		14.1.1 pcibus_class_init 与 pci_driver_init	

函数 .....	368	BAR 寄存器 .....	407
14.1.2 pci_arch_init 函数 .....	369	14.4 Linux PowerPC 如何初始化 PCI	
14.1.3 pci_slot_init 和 pci_subsys_init 函数 .....	372	总线树 .....	412
14.1.4 与 PCI 总线初始化相关的其他 函数 .....	373	14.5 小结 .....	416
14.2 x86 处理器的 ACPI .....	374	<b>第 15 章 Linux PCI 的中断处理 .....</b>	417
14.2.1 ACPI 驱动程序与 AML 解释器 .....	377	15.1 PCI 总线的中断路由 .....	417
14.2.2 ACPI 表 .....	380	15.1.1 PCI 设备如何获取 irq 号 .....	419
14.2.3 ACPI 表的使用实例 .....	382	15.1.2 PCI 中断路由表 .....	426
14.3 基于 ACPI 机制的 Linux PCI 的 初始化 .....	388	15.1.3 PCI 插槽使用的 irq 号 .....	428
14.3.1 基本的准备工作 .....	388	15.2 使用 MSI/MSIX 中断机制 申请中断向量 .....	432
14.3.2 Linux PCI 初始化 PCI 总线号 ...	393	15.2.1 Linux 如何使能 MSI 中断 机制 .....	432
14.3.3 Linux PCI 检查 PCI 设备使用的 BAR 空间 .....	404	15.2.2 Linux 如何使能 MSI-X 中断 机制 .....	437
14.3.4 Linux PCI 分配 PCI 设备使用的		15.3 小结 .....	440
		参考文献 .....	441

# 第 I 篇 PCI 体系结构概述

PCI (Peripheral Component Interconnect) 总线的诞生与 PC (Personal Computer) 的蓬勃发展密切相关。在处理器体系结构中，PCI 总线属于局部总线 (Local Bus)。局部总线作为系统总线的延伸，其主要功能是连接外部设备。

处理器主频的不断提升，要求速度更快、带宽更高的局部总线。起初 PC 使用 8 位的 XT 总线作为局部总线，很快升级到 16 位的 ISA (Industry Standard Architecture) 总线，并逐步发展到 32 位的 EISA (Extended Industry Standard Architecture)、VESA (Video Electronics Standards Association) 和 MCA (Micro Channel Architecture) 总线。

PCI 总线规范在 20 世纪 90 年代提出。这条总线推出之后，很快得到了各大主流半导体厂商的认同，并迅速统一了当时并存的各类局部总线，EISA、VESA 等其他 32 位总线很快就被 PCI 总线淘汰了。从那时起，PCI 总线一直在处理器体系结构中占有重要地位。

在此后相当长的一段时间里，处理器系统的大多数外部设备都是直接或者间接地与 PCI 总线相连。即使目前 PCI Express 总线逐步取代 PCI 总线成为 PC 局部总线的主流，也不能掩盖 PCI 总线的光芒。从软件层面上看，PCI Express 总线与 PCI 总线基本兼容；从硬件层面上看，PCI Express 总线在很大程度上继承了 PCI 总线的设计思路。因此 PCI 总线依然是软硬件工程师在进行处理器系统的开发与设计时必须掌握的一条局部总线。

PCI 总线规范由 Intel 的 IAL (Intel Architecture Lab) 提出，其 V1.0 规范在 1992 年 6 月 22 日正式发布。IAL 是 Intel 的一个重要实验室，USB (Universal Serial Bus)、AGP (Accelerated Graphics Port)、PCI Express 总线规范和 PC 的南北桥结构都是由这个实验室提出的。

IAL 起初的研究领域包括硬件和软件，但是 IAL 在软件领域的研究遭到了 Microsoft 的抵触，IAL 提出的许多软件规范并不被 Microsoft 认可，于是 IAL 更专注硬件领域，并在 PC 体系架构上取得了一个又一个突破。IAL 是现代 PC 体系架构的重要奠基者。2001 年，IAL 由于其创始人的离去而临时解散。2005 年，Intel 重建了这个实验室。

PCI 总线 V1.0 规范仅针对在一个 PCB (Printed Circuit Board) 环境内的器件之间的互连，而 1993 年 4 月 30 日发布的 V2.0 规范增加了对 PCI 插槽的支持。1995 年 6 月 1 日，PCI V2.1 总线规范发布，这个规范具有里程碑意义。正是这个规范使得 PCI 总线大规模普及，至此 PCI 总线完成了对(E)ISA 和 MCA 总线的替换。

至 1996 年，VESA 总线也逐渐离开了人们的视线，当然 PCI 总线并不能完全提供显卡所需要的带宽，真正替代 VESA 总线的是 AGP 总线。随后 PCISIG (PCI Special Interest Group) 陆续发布了 PCI 总线 V2.2、V2.3 规范，并最终将 PCI 总线规范定格在 V3.0。

除了 PCI 总线规范外，PCISIG 还定义了一些与 PCI 总线相关的规范，如 PCMCIA（Personal Computer Memory Card International Association）规范和 MiniPCI 规范。其中 PCMCIA 规范主要针对 Laptop 应用，后来 PCMCIA 升级为 PC Card（Cardbus）规范，而 PC Card 又升级为 ExpressCard 规范。

PC Card 规范基于 32 位、33MHz 的 PCI 总线；而 ExpressCard 规范基于 PCI Express 和 USB 2.0。这两个规范都在 Laptop 领域中获得了成功。除了 PCMCIA 规范外，Mini PCI 总线也非常流行，与标准 PCI 插槽相比，Mini PCI 插槽占用面积较小，适用于一些对尺寸有要求的应用。

除了以上规范之外，PCISIG 还推出了一系列和 PCI 总线直接相关的规范。如 PCI-to-PCI 桥规范、PCI 电源管理规范、PCI 热插拔规范和 CompactPCI 总线规范。其中 PCI-to-PCI 桥规范最为重要，理解 PCI-to-PCI 桥是理解 PCI 体系结构的基础；而 CompactPCI 总线规范多用于具有背板结构的大型系统，并支持热插拔。

PCISIG 在 PCI 总线规范的基础上，进一步提出 PCI-X 规范。与 PCI 总线相比，PCI-X 总线规范可以支持 133MHz、266MHz 和 533MHz 的总线频率，并在传送规则上做了一些改动。虽然 PCI-X 总线还没有得到大规模普及就被 PCI Express 总线替代，但是在 PCI-X 总线中提出的许多设计思想仍然被 PCI Express 总线继承。

PCI 总线规范是 Intel 对 PC 领域做出的一个巨大贡献。Intel 也在 PCI 总线规范中留下了深深的印记，PCI 总线规范的许多内容都与基于 IA（Intel Architecture）架构的 x86 处理器密切相关。但是这并不妨碍其他处理器系统使用 PCI 总线，事实上 PCI 总线在非 x86 处理器系统上也取得了巨大的成功。目前绝大多数处理器系统都使用 PCI/PCI Express 总线连接外部设备，特别是一些通用外设。

随着时间的推移，PCI 和 PCI-X 总线逐步遇到瓶颈。PCI 和 PCI-X 总线使用单端并行信号进行数据传递，由于单端信号容易被外部系统干扰，其总线频率很难进一步提高。目前，为了获得更高的总线频率以提高总线带宽，高速串行总线逐步替代了并行总线。PCI Express 总线也逐渐替代 PCI 总线成为主流。但是从系统软件的角度上看，PCI Express 总线仍然基于 PCI 总线。理解 PCI Express 总线的一个基础是深入理解 PCI 总线，同时 PCI Express 总线也继承了 PCI 总线的许多概念。本篇将详细介绍与处理器体系结构相关的一些必备的 PCI 总线知识。

为简化起见，本篇主要介绍 PCI 总线的 32 位地址模式。在实际应用中，使用 64 位地址模式的 PCI 设备非常少。而且在 PCI Express 总线逐渐取代 PCI 总线的大趋势之下，将来也很难会有更多的使用 64 位地址的 PCI 设备。如果读者需要掌握 PCI 总线的 64 位地址模式，请自行阅读 PCI 总线的相关规范。实际上，如果读者真正掌握了 PCI 总线的 32 位地址模式之后，理解 64 位地址模式并不困难。

为节省篇幅，下文将 PCI Express 总线简称为 PCIe 总线，PCI-to-PCI 桥简称为 PCI 桥，PCI Express-to-PCI 桥简称为 PCIe 桥，Host-to-PCI 主桥简称为 HOST 主桥。值得注意的是许多书籍将 HOST 主桥称为 PCI 主桥或者 PCI 总线控制器。

# 第1章 PCI总线的基本知识

PCI总线作为处理器系统的局部总线，其主要目的是为了连接外部设备，而不是作为处理器的系统总线连接Cache和主存储器。但是PCI总线、系统总线和处理器体系结构之间依然存在着紧密的联系。

PCI总线作为系统总线的延伸，其设计考虑了许多与处理器相关的内容，如处理器的Cache共享一致性和数据完整性，以及如何与处理器进行数据交换等一系列内容。其中Cache共享一致性和数据完整性是现代处理器局部总线的设计的重点和难点，也是本书将重点讲述的主题之一。

孤立地研究PCI总线并不可取，因为PCI总线仅是处理器系统的一个部分。深入理解PCI总线需要了解一些与处理器体系结构相关的知识。这些知识是本书侧重描述的，同时也是PCI总线规范忽略的内容。脱离实际的处理器系统，不容易也不可能深入理解PCI总线规范。

对于今天的读者来说，PCI总线提出的许多概念略显过时，也有许多不足之处。但是在当年，PCI总线与之前存在的其他并行局部总线如ISA、EISA和MCA总线相比，具有许多突出的优点，是一个全新的设计。

## (1) PCI总线空间与处理器空间隔离

PCI设备具有独立的地址空间，即PCI总线地址空间，该空间与存储器地址空间通过HOST主桥隔离。处理器需要通过HOST主桥才能访问PCI设备，而PCI设备需要通过HOST主桥才能访问主存储器。在HOST主桥中含有许多缓冲，这些缓冲使得处理器总线与PCI总线工作在各自的时钟频率中，互不干扰。HOST主桥的存在也使得PCI设备和处理器可以方便地共享主存储器资源。

处理器访问PCI设备时，必须通过HOST主桥进行地址转换；而PCI设备访问主存储器时，也需要通过HOST主桥进行地址转换。HOST主桥的一个重要作用就是将处理器访问的存储器地址转换为PCI总线地址。PCI设备使用的地址空间是属于PCI总线域的，这与存储器地址空间不同。

x86处理器对PCI总线域与存储器域的划分并不明晰，这也使得许多程序员并没有准确地区分PCI总线域地址空间与存储器域地址空间。而本书将反复强调存储器地址和PCI总线地址的区别，因为这是理解PCI体系结构的重要内容。

PCI规范并没有对HOST主桥的设计进行约束。每一个处理器厂商使用的HOST主桥，其设计都不尽相同。HOST主桥是联系PCI总线与处理器的核心部件，掌握HOST主桥的实现机制是深入理解PCI体系结构的前提。

本书将以Freescale的PowerPC处理器和Intel的x86处理器为例，说明各自HOST主桥的实现方式，值得注意的是本书涉及的PowerPC处理器仅针对Freescale的PowerPC处理器，而不包含IBM和AMCC的Power和PowerPC处理器。而且如果没有特别说明，本书中涉及的x86处理器特指Intel的处理器，而不是其他厂商的x86处理器。

## (2) 可扩展性

PCI 总线具有很强的扩展性。在 PCI 总线中，HOST 主桥可以直接推出一条 PCI 总线，这条总线也是该 HOST 主桥管理的第一条 PCI 总线，该总线还可以通过 PCI 桥扩展出一系列 PCI 总线，并以 HOST 主桥为根节点，形成 1 棵 PCI 总线树。这些 PCI 总线都可以连接 PCI 设备，但是在 1 棵 PCI 总线树上，最多只能挂接 256 个 PCI 设备（包括 PCI 桥）。

在同一条 PCI 总线上的设备间可以直接通信，而并不会影响其他 PCI 总线上设备间的数据通信。隶属于同一棵 PCI 总线树上的 PCI 设备，也可以直接通信，但是需要通过 PCI 桥进行数据转发。

PCI 桥是 PCI 总线的一个重要组成部件，该部件的存在使得 PCI 总线极具扩展性。PCI 桥也是有别于其他局部总线的一个重要部件。在“以 HOST 主桥为根节点”的 PCI 总线树中，每一个 PCI 桥下也可以连接一个 PCI 总线子树，PCI 桥下的 PCI 总线仍然可以使用 PCI 桥继续进行总线扩展。

PCI 桥可以管理这个 PCI 总线子树，PCI 桥的配置空间含有一系列管理 PCI 总线子树的配置寄存器。在 PCI 桥的两端，分别连接了两条总线，分别是上游总线（Primary Bus）和下游总线（Secondary Bus）。其中与处理器距离较近的总线被称为上游总线，另一条被称为下游总线。这两条总线间的通信需要通过 PCI 桥进行。PCI 桥中的许多概念被 PCIe 总线采纳，理解 PCI 桥也是理解 PCIe 体系结构的基础。

## (3) 动态配置机制

PCI 设备使用的地址可以根据需要由系统软件动态分配。PCI 总线使用这种方式合理地解决了设备间的地址冲突，从而实现了“即插即用”功能。因此 PCI 总线不需要使用 ISA 或者 EISA 接口卡为解决地址冲突而使用的硬件跳线。

每一个 PCI 设备都有独立的配置空间，在配置空间中含有该设备在 PCI 总线中使用的基址，系统软件可以动态配置这个基址，从而保证每一个 PCI 设备使用的物理地址并不相同。PCI 桥的配置空间中含有其下 PCI 子树所能使用的地址范围。

## (4) 总线带宽

PCI 总线与之前的局部总线相比，极大提高了数据传送带宽，32 位/33 MHz 的 PCI 总线可以提供 132 MB/s 的峰值带宽，而 64 位/66 MHz 的 PCI 总线可以提供的峰值带宽为 532 MB/s。虽然 PCI 总线所能提供的峰值带宽远不能和 PCIe 总线相比，但是与之前的局部总线 ISA、EISA 和 MCA 总线相比，仍然具有极大的优势。

ISA 总线的最高主频为 8 MHz，位宽为 16，其峰值带宽为 16 MB/s；EISA 总线的最高主频为 8.33 MHz，位宽为 32，其峰值带宽为 33 MB/s；而 MCA 总线的最高主频为 10 MHz，最高位宽为 32，其峰值带宽为 40 MB/s。PCI 总线提供的峰值带宽远高于这些总线。

## (5) 共享总线机制

PCI 设备通过仲裁获得 PCI 总线的使用权后，才能进行数据传送，在 PCI 总线上进行数据传送，并不需要处理器进行干预。

PCI 总线仲裁器不在 PCI 总线规范定义的范围内，也不一定是 HOST 主桥和 PCI 桥的一部分。虽然绝大多数 HOST 主桥和 PCI 桥都包含 PCI 总线仲裁器，但是在某些处理器系统的设计中也可以使用独立的 PCI 总线仲裁器。如在 PowerPC 处理器的 HOST 主桥中含有 PCI 总线仲裁器，但是用户可以关闭这个总线仲裁器，而使用独立的 PCI 总线仲裁器。

PCI 设备使用共享总线方式进行数据传递，在同一条总线上，所有 PCI 设备共享同一总线带宽，这将极大地影响 PCI 总线的利用率。这种机制显然不如 PCIe 总线采用的交换结构，但是在 PCI 总线盛行的年代，半导体的工艺、设计能力和制作成本决定了采用共享总线方式是当时的最优选择。

#### (6) 中断机制

PCI 总线上的设备可以通过四根中断请求信号 INTA ~ D# 向处理器提交中断请求。与 ISA 总线上的设备不同，PCI 总线上的设备可以共享这些中断请求信号，不同的 PCI 设备可以将这些中断请求信号“线与”后，与中断控制器的中断请求引脚连接。PCI 设备的配置空间记录了该设备使用这四根中断请求信号的信息。

PCI 总线还进一步提出了 MSI (Message Signal Interrupt) 机制，该机制使用存储器写总线事务传递中断请求，并可以使用 x86 处理器 FSB (Front Side Bus) 总线提供的 Interrupt Message 总线事务，从而提高了 PCI 设备的中断请求效率。

虽然从现代总线技术的角度来看，PCI 总线仍有许多不足之处，但也不能否认 PCI 总线已经获得了巨大的成功。不仅 x86 处理器将 PCI 总线作为标准的局部总线连接各类外部设备，PowerPC、MIPS 和 ARM<sup>○</sup>处理器也将 PCI 总线作为标准局部总线。除此之外，基于 PCI 总线的外部设备，如以太网控制器、声卡、硬盘控制器等，也已经成为主流。

### 1.1 PCI 总线的组成结构

如上文所述，PCI 总线作为处理器系统的局部总线，是处理器系统的一个组成部件，讲述 PCI 总线的组成结构不能离开处理器系统这个大环境。在一个处理器系统中，与 PCI 总线相关的模块如图 1-1 所示。

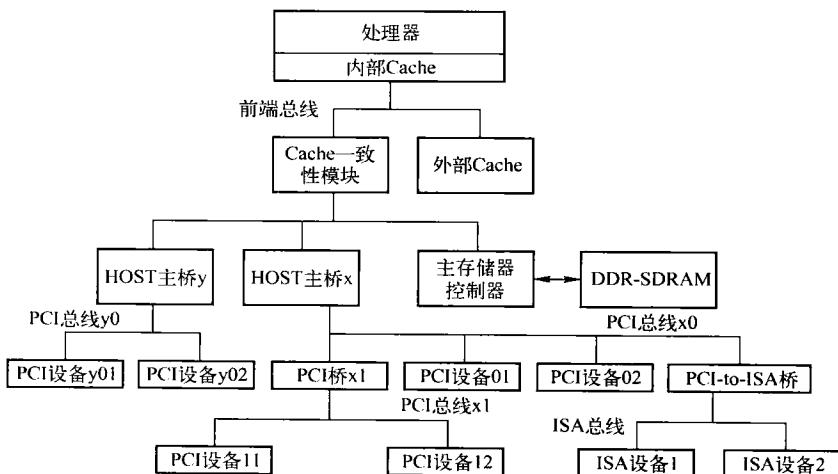


图 1-1 基于 PCI 总线的处理器系统

○ 在 ARM 处理器中，使用 SoC 平台总线，即 AMBA 总线，连接片内设备。但是某些 ARM 生产厂商，依然使用 AMBA-to-PCI 桥推出 PCI 总线，以连接 PCI 设备。

图中与 PCI 总线相关的模块包括：HOST 主桥、PCI 总线、PCI 桥和 PCI 设备。PCI 总线由 HOST 主桥和 PCI 桥推出，HOST 主桥与主存储器控制器在同一级总线上，因此 PCI 设备可以方便地通过 HOST 主桥访问主存储器，即进行 DMA 操作。

值得注意的是，PCI 设备的 DMA 操作需要与处理器系统的 Cache 进行一致性操作，当 PCI 设备通过 HOST 主桥访问主存储器时，Cache 一致性模块将进行地址监听，并根据监听的结果改变 Cache 的状态。

在一些简单的处理器系统中，可能不含有 PCI 桥，此时所有 PCI 设备都是连接在 HOST 主桥推出的 PCI 总线上。此外在一些处理器系统中可能含有多个 HOST 主桥，如图 1-1 所示的处理器系统中含有 HOST 主桥 x 和 HOST 主桥 Y。

### 1.1.1 HOST 主桥

HOST 主桥是一个很特别的桥片，其主要功能是隔离处理器系统的存储器域与处理器系统的 PCI 总线域，管理 PCI 总线域，并完成处理器与 PCI 设备间的数据交换。处理器与 PCI 设备间的数据交换主要由“处理器访问 PCI 设备的地址空间”和“PCI 设备使用 DMA 机制访问主存储器”这两部分组成。

为简便起见，下面将处理器系统的存储器域简称为存储器域，而将处理器系统的 PCI 总线域称为 PCI 总线域，存储器域和 PCI 总线域的详细介绍见第 2.1 节。值得注意的是，在一个处理器系统中，有几个 HOST 主桥，就有几个 PCI 总线域。

HOST 主桥在处理器系统中的位置并不相同，如 PowerPC 处理器将 HOST 主桥与处理器集成在一个芯片中。而有些处理器不进行这种集成，如 x86 处理器使用南北桥结构，处理器内核在一个芯片中，而 HOST 主桥在北桥中。但是从处理器体系结构的角度看，这些集成方式并不重要。

PCI 设备通过 HOST 主桥访问主存储器时，需要与处理器的 Cache 进行一致性操作，因此在设计 HOST 主桥时需要重点考虑 Cache 一致性操作。在 HOST 主桥中，还含有许多数据缓冲，以支持 PCI 总线的预读机制。

HOST 主桥是联系处理器与 PCI 设备的桥梁。在一个处理器系统中，每一个 HOST 主桥都管理了一棵 PCI 总线树，在同一棵 PCI 总线树上的所有 PCI 设备属于同一个 PCI 总线域。如图 1-1 所示，HOST 主桥 x 之下的 PCI 设备属于 PCI 总线 x 域，而 HOST 主桥 y 之下的 PCI 设备属于 PCI 总线 y 域。在这棵总线树上的所有 PCI 设备的配置空间都由 HOST 主桥通过配置读写总线周期访问。

如果 HOST 主桥支持 PCI V3.0 规范的 Peer-to-Peer 数据传送方式，那么分属不同 PCI 总线域的 PCI 设备可以直接进行数据交换。如图 1-1 所示，如果 HOST 主桥 y 支持 Peer-to-Peer 数据传送方式，PCI 设备 y01 可以直接访问 PCI 设备 01 或者 PCI 设备 11，而不需要通过处理器的参与。但是这种跨越总线域的数据传送方式在 PC 架构中并不常用，在 PC 架构中，重点考虑的是 PCI 设备与主存储器之间的数据交换，而不是 PCI 设备之间的数据交换。此外在 PC 架构中，具有两个 HOST 主桥的处理器系统也并不多见。

在 PowerPC 处理器中，HOST 主桥可以通过设置 Inbound 寄存器，使得分属于不同 PCI 总线域的设备可以直接通信。许多 PowerPC 处理器都具有多个 HOST 主桥，有关 PowerPC 处理器使用的 HOST 主桥详见第 2.2 节。