

```
c010606f: 55      push   %ebp
c0106070: 89 e5   mov    %esp,%ebp
c0106072: 56     push   %esi
c0106073: 53     push   %ebx
c0106074: 83 ec 0c sub    $0xc,%esp
c0106077: 68 44 a2 12 c0 push  $0xc012a24
c010607c: e8 ed 1c 00 00 call   c0107d6e <printf>
c0106081: be 00 60 10 00 mov    $0x106000,%esi
c0106086: c1 ee 0c shr    $0xc,%esi
c0106089: 83 c4 08 add    $0x8,%esp
c010608c: 56     push   %esi
c010608d: 68 6c a2 12 c0 push  $0xc012a26c
c0106092: e8 d7 1c 00 00 call   c0107d6e <printf>
c0106097: 83 c4 10 add    $0x10,%esp
c010609a: bb 00 00 00 00 mov    $0x0,%ebx
c010609f: eb 14   jmp    c01060b5 <_kernel_post_init+0x46>
c01060a1: 89 d8   mov    %ebx,%eax
c01060a3: c1 e0 0c shl    $0xc,%eax
c01060a6: 83 ec 0c sub    $0xc,%esp
c01060a9: 50     push   %eax
c01060aa: e8 9f 0a 00 00 call   c0106b4e <vmm_unmap_page>
c01060af: 83 c0   add    $0x1,%ebx
c01060b2: 83 c4 10 add    $0x10,%esp
c01060b5: 39 f3   cmp    %eax,%ebx
c01060b7: 72 e8   jb     c01060d1 <_kernel_post_init+0x32>
c01060b9: e8 54 1c 00 00 call   c0107d6e <printf>
c01060be: 85 c0   test   %eax,%eax
c01060c0: 74 07   jbe   c01060c1 <_kernel_post_init+0x5a>
c01060c2: 8d 65 fa mov    %ebp,%ebx
c01060c5: 5b     pop    %ebx
c01060c6: 5e     pop    %esi
c01060c7: 5d     pop    %ebp
c01060c8: c3     ret
c01060c9: 83 ec 04 sub    $0x4,%esp
c01060cc: 6a 40   push  $0x40
c01060ce: 68 af a0 12 c0 push  $0xc012a0af
```

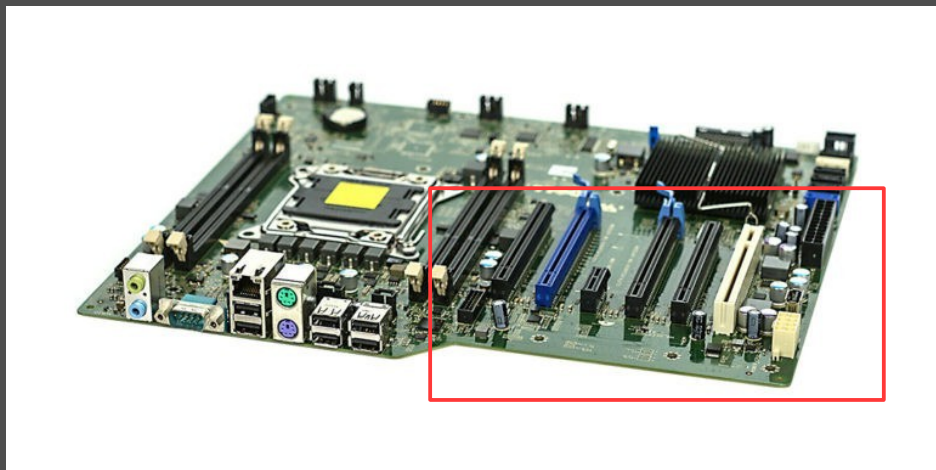
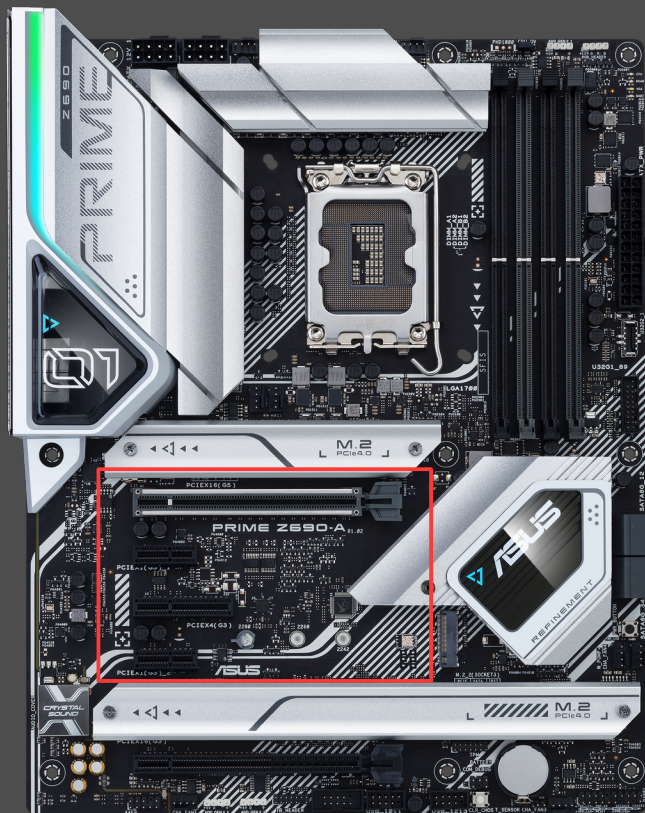
# LunaixOS

## 从零开始

# 自制操作系统



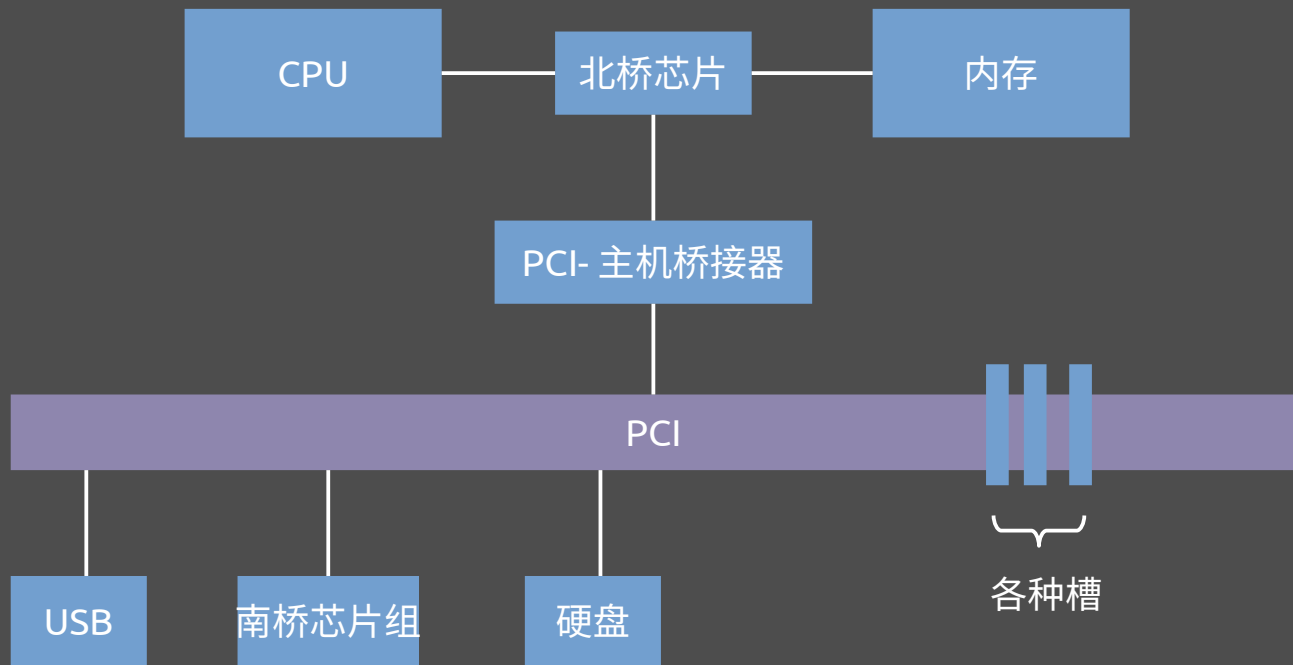
接入 **PCI EXPRESS**



但凡攒过机的，多多少少都知道这个东西。

# PCI 可不只是那几个槽!

PCI：外围设备互联总线





我们主要通过 PCI 去.....

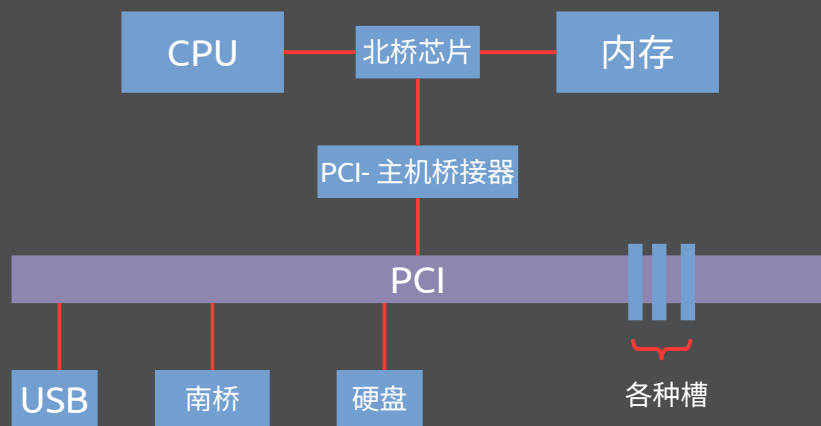
1. 发现设备
2. 初始化和配置设备
3. 操控设备

不需要关心.....

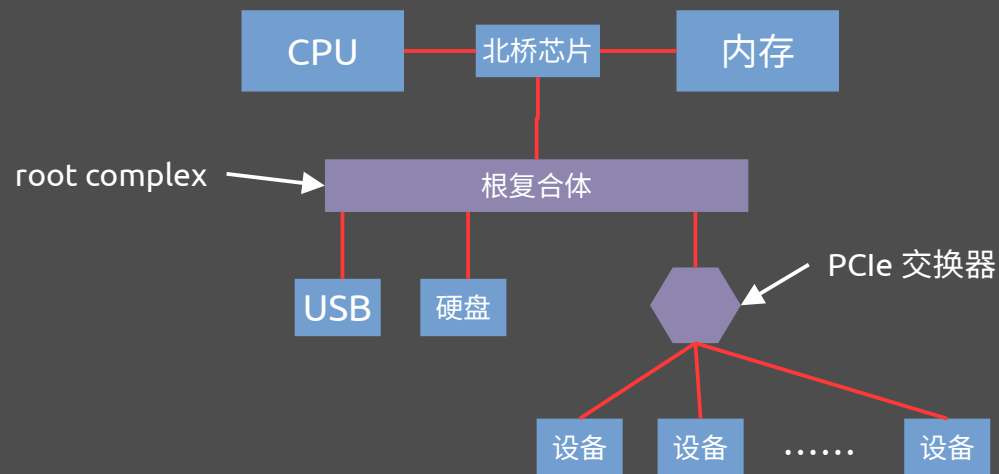
1. 如何协调设备之间的通讯
2. 如何具体的读写 PCI 总线

## LunaixOS 将实现 PCI 3.0 的规范标准

PCI 是实打实的总线结构



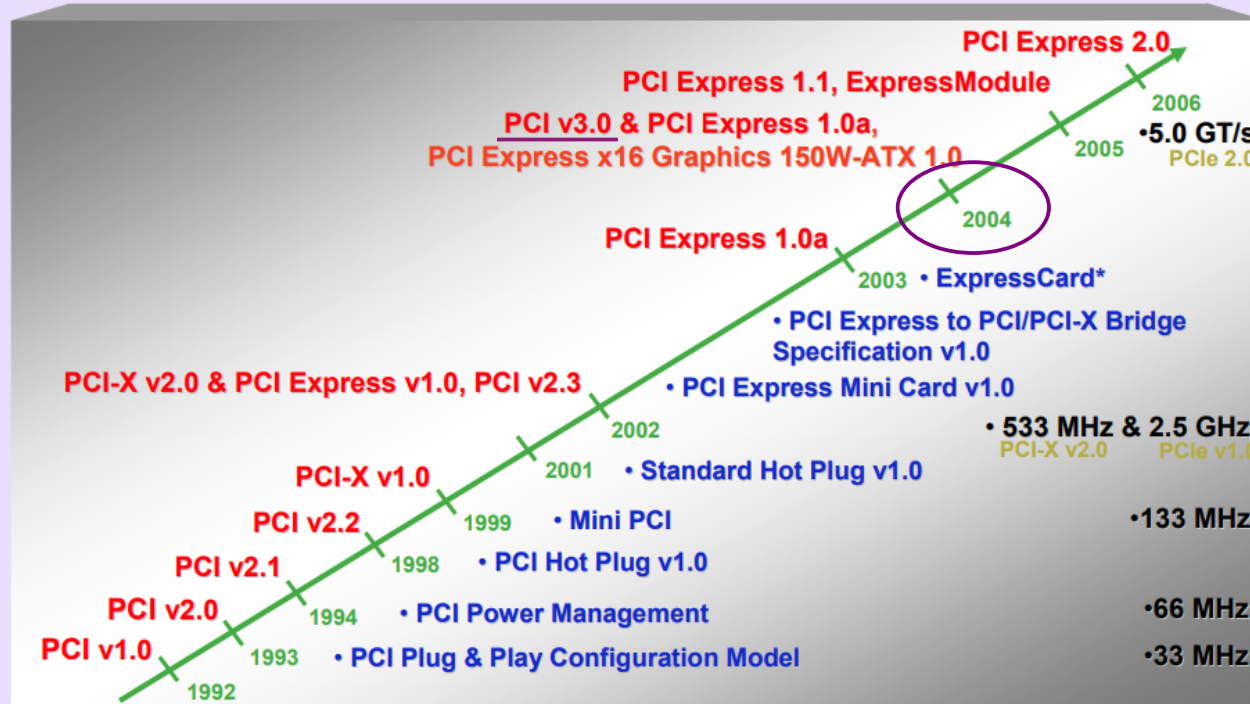
PCIe 基于报文交换机制，类似网络



从软件层面上来看，PCI 和 PCIe 是可以做到无缝兼容的



## PCI Family History



\* "ExpressCard" is a trademark, or registered trademark of PCMCIA in the United States and/or other countries.

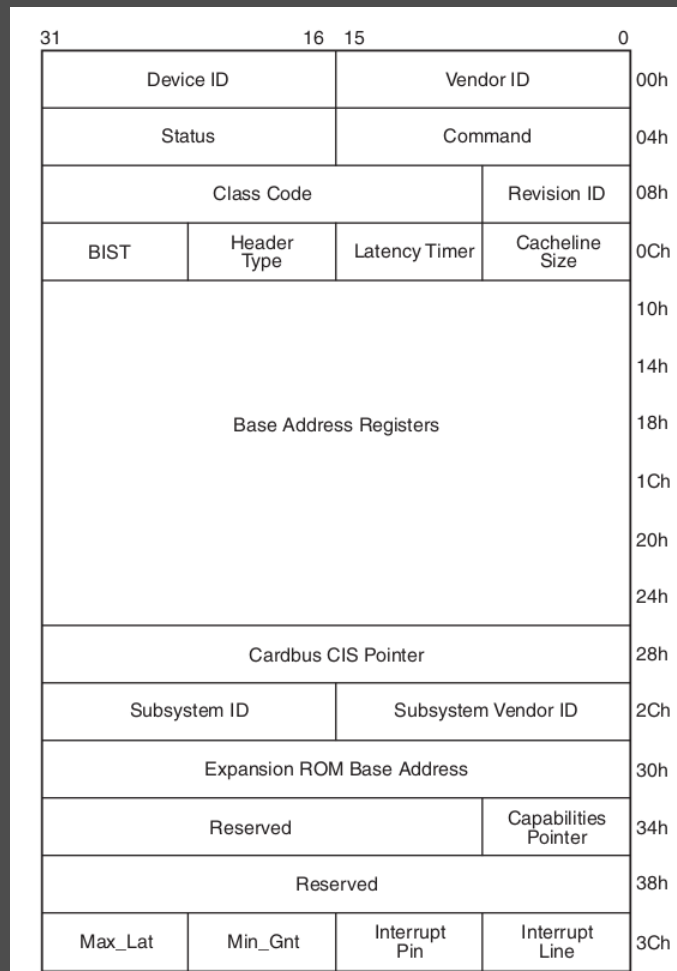
配置空间是.....

一组 32 位寄存器组成的区域。

大小 256 字节（64 个寄存器）。

在前 64 字节包含了一个头部。

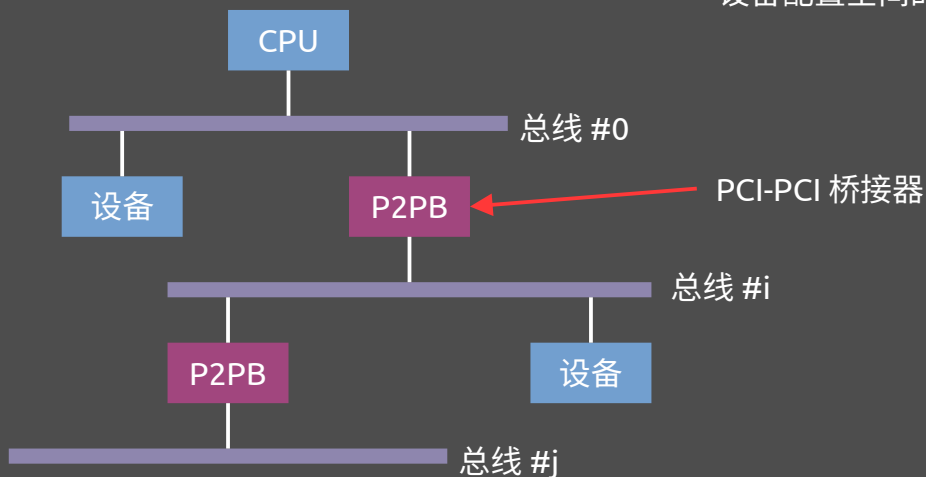
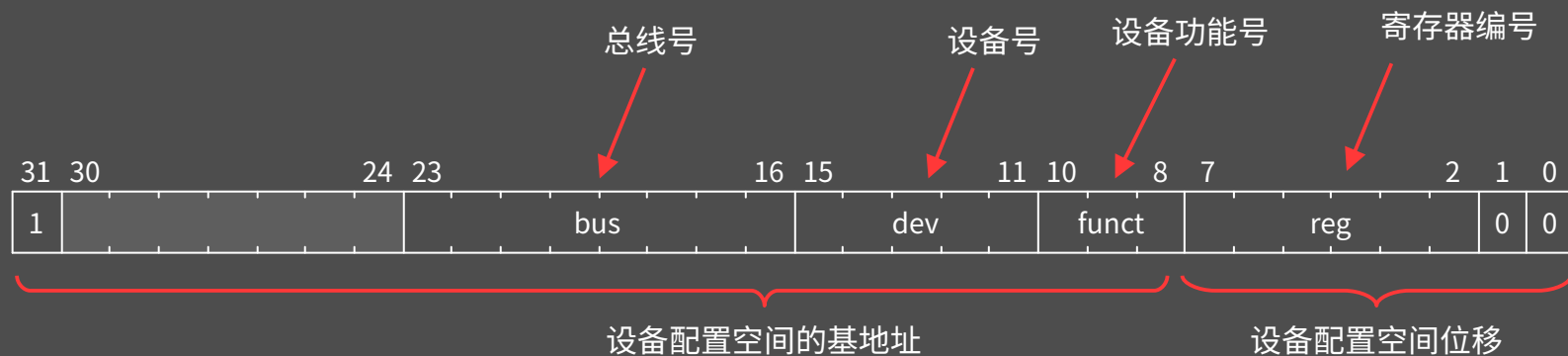
包括了设备的所有信息以及 PCI 相关设置



A-0191

Figure 6-1: Type 00h Configuration Space Header

每个设备在 PCI 上都有一个唯一的地址



PCI 是树状结构，最多 256 条总线



PCI - 主机适配器提供的 I/O 口

```
#define PCI_CONFIG_ADDR 0xcf8 ← 地址 I/O  
#define PCI_CONFIG_DATA 0xcfc ← 寄存器 I/O
```

```
inline pci_reg_t  
pci_read_cspace(uint32_t base, int offset)  
{  
    io_outl(PCI_CONFIG_ADDR, base | (offset & ~0x3));  
    return io_inl(PCI_CONFIG_DATA);  
}  
  
inline void  
pci_write_cspace(uint32_t base, int offset, pci_reg_t data)  
{  
    io_outl(PCI_CONFIG_ADDR, base | (offset & ~0x3));  
    io_outl(PCI_CONFIG_DATA, data);  
}
```

向适配器写入地址



穷举出所有可能的设备地址，进行探测。

通过检测生产商 ID 来确定设备是否存在。

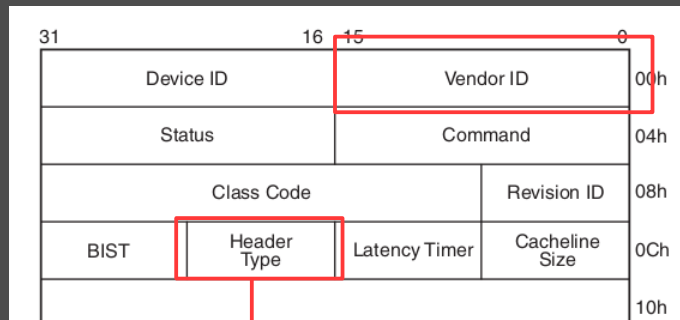
0xFFFF 则为不存在。

256 条总线

每条总线最多 32 个设备

每个设备最多 8 个功能

常数时间：  $T(n) = 256 * 32 * 8 = 65536$



检查第七位是否置位来判断该设备是不是多功能设备。

大多数设备都只有一个功能。

平均时间：  $T(n) = 2048$



```
void
pci_probe()
{
    // 暴力扫描所有PCI设备
    // XXX: 尽管最多会有256条PCI总线, 但就目前而言, 只考虑bus #0就足够了
    for (int bus = 0; bus < 1; bus++) {
        for (int dev = 0; dev < 32; dev++) {
            pci_probe_device(bus, dev, 0);
        }
    }
}
```

```
void
pci_probe_device(int bus, int dev, int funct)
{
    uint32_t base = PCI_ADDRESS(bus, dev, funct);
    pci_reg_t reg1 = pci_read_cspace(base, 0);

    // Vendor=0xffff则表示设备不存在
    if (PCI_DEV_VENDOR(reg1) == PCI_VENDOR_INVLD) {
        return;
    }

    pci_reg_t hdr_type = pci_read_cspace(base, 0xc);
    hdr_type = (hdr_type >> 16) & 0xff;

    // 防止堆栈溢出
    // QEMU的ICH9/Q35实现似乎有点问题, 对于多功能设备的每一个功能的header type
    // 都将第七位置位。而virtualbox就没有这个毛病。
    if ((hdr_type & 0x80) && funct == 0) {
        hdr_type = hdr_type & ~0x80;
        // 探测多用途设备 (multi-function device)
        for (int i = 1; i < 7; i++) {
            pci_probe_device(bus, dev, i);
        }
    }

    if (hdr_type != PCI_TDEV) {
        // XXX: 目前忽略所有桥接设备, 比如PCI-PCI桥接器, 或者是CardBus桥接器
        return;
    }

    pci_reg_t intr = pci_read_cspace(base, 0x3c);
    pci_reg_t class = pci_read_cspace(base, 0x8);

    struct pci_device* device = lxmalloc(sizeof(struct pci_device));
    *device = (struct pci_device){ .cspace_base = base,
                                  .class_info = class,
                                  .device_info = reg1,
                                  .intr_info = intr };

    pci_probe_msi_info(device);

    llist_append(&pci_devices, &device->dev_chain);
}
```

## PCI 相关的初始化: Command

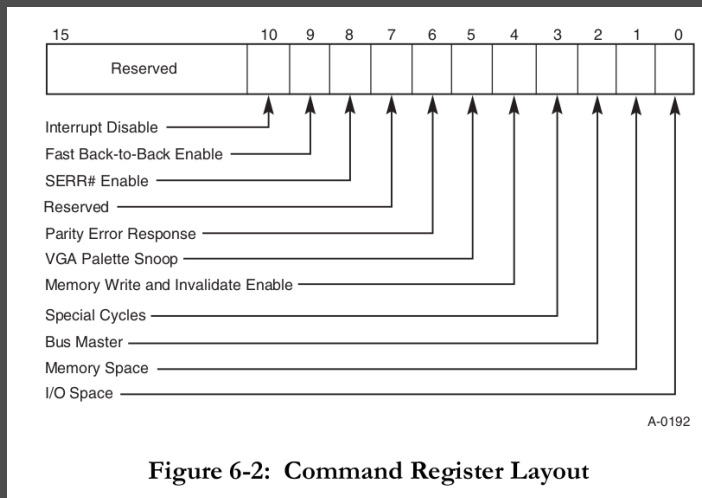


Figure 6-2: Command Register Layout

设备相关, 不同的设备, 对每个位的具体解释不一样。

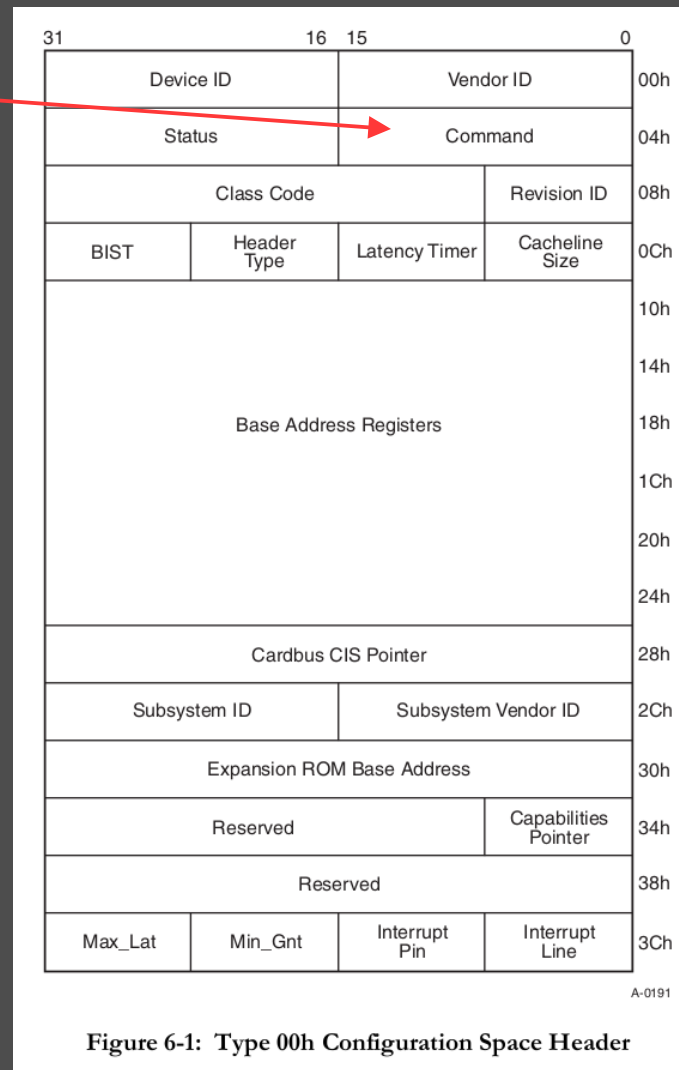
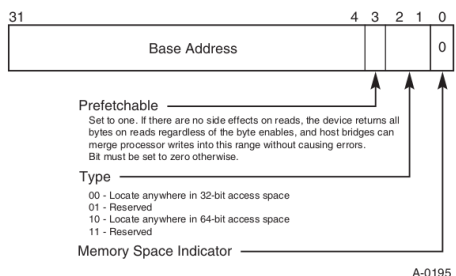


Figure 6-1: Type 00h Configuration Space Header

也是设备相关的问题。

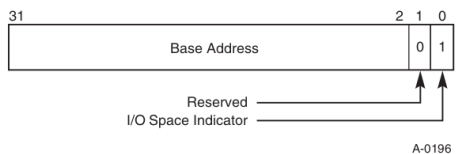
不同的设备会提供不同的寄存器，OS 通过往这些寄存器里写值，向设备发送命令。

设备相关的寄存器的基地址会在 **BAR** 中给出：



对于内存映射的寄存器  
(MMIO)

Figure 6-5: Base Address Register for Memory



对于 I/O 映射的寄存器

Figure 6-6: Base Address Register for I/O

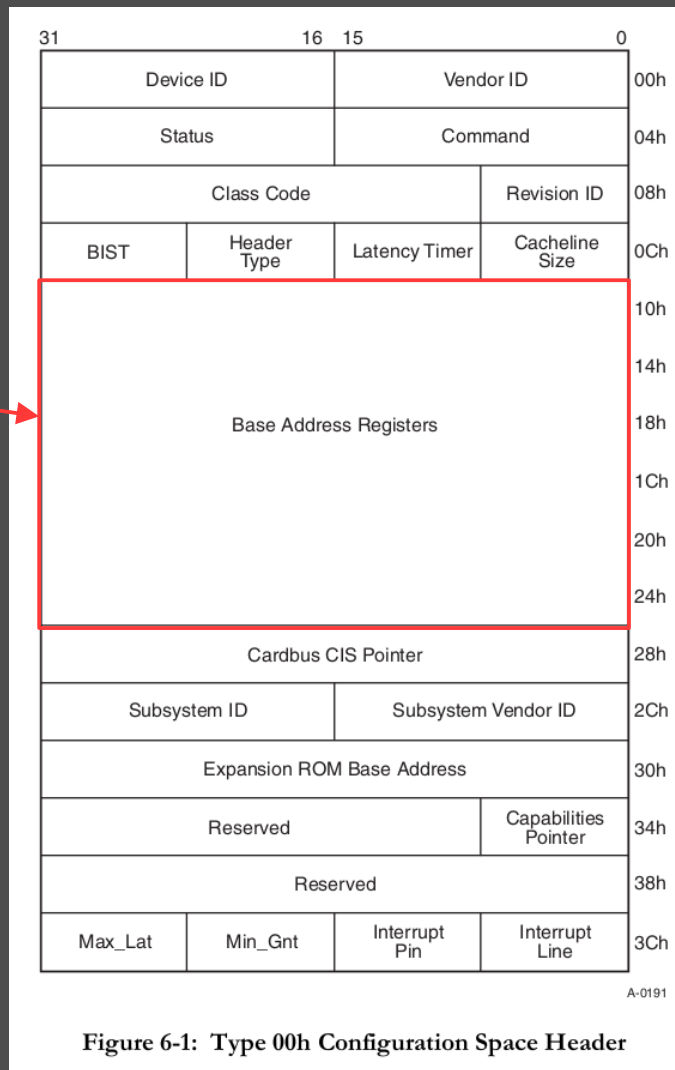


Figure 6-1: Type 00h Configuration Space Header

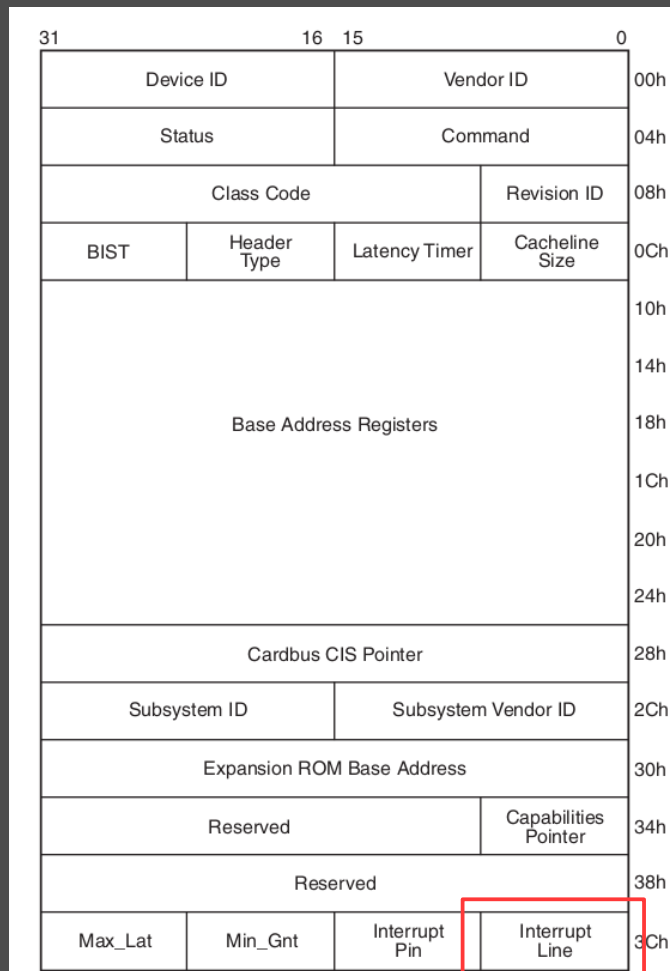
一个设备需要通过中断去和 CPU 通讯。

Interrupt Line：使用的 IRQ

那么我们可以通过 ACPI 查询到 IRQ → IOAPIC 针脚的映射关系，从而通过 IOAPIC 来管理这些中断。

但这种方法有两个很严重的问题.....

1. IRQ 共用的情况
2. 竞态条件

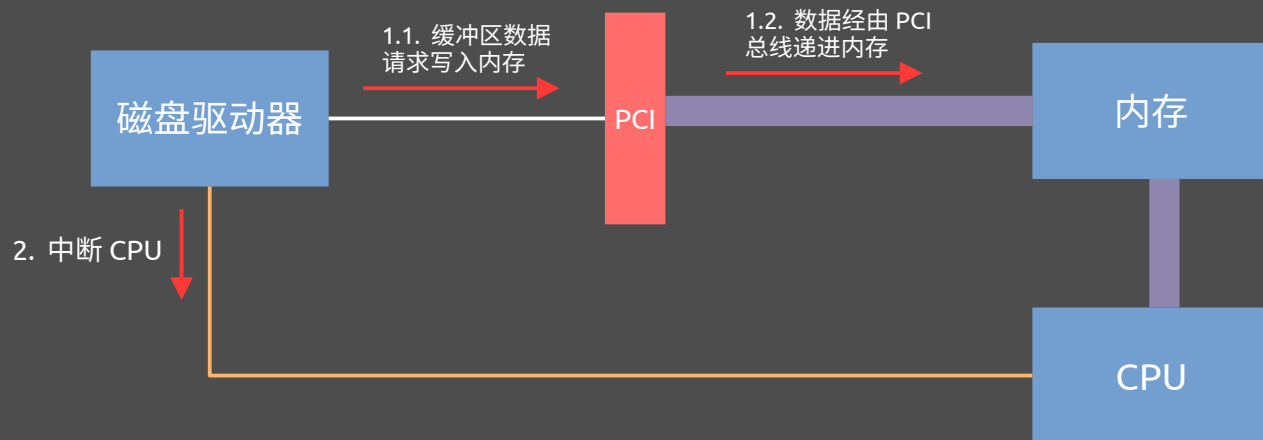


A-0191

Figure 6-1: Type 00h Configuration Space Header



# PCI 中断机制：竞态条件

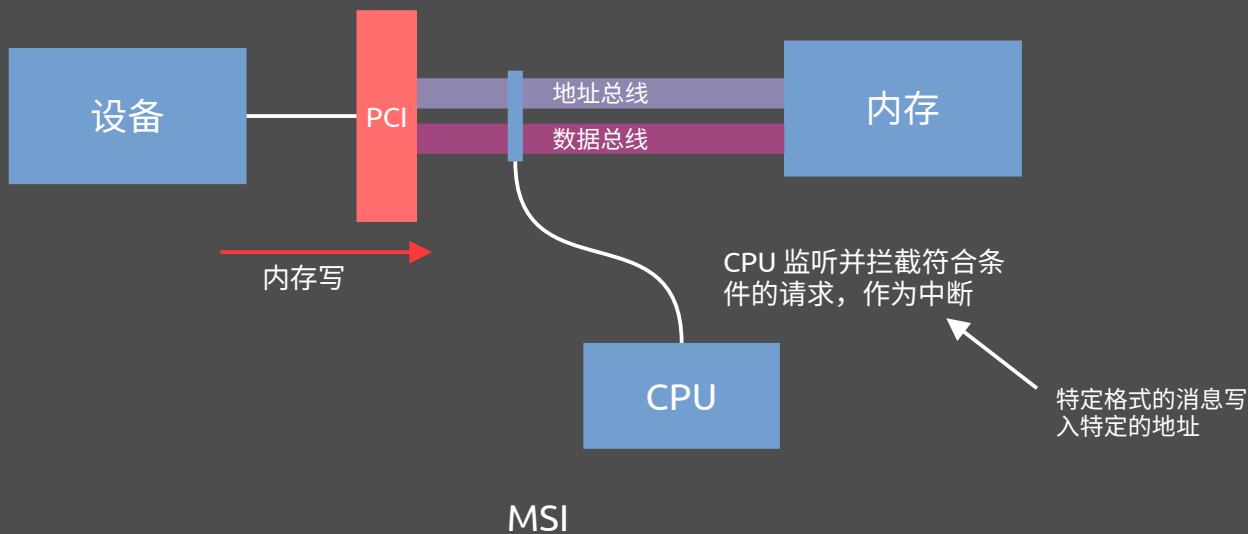


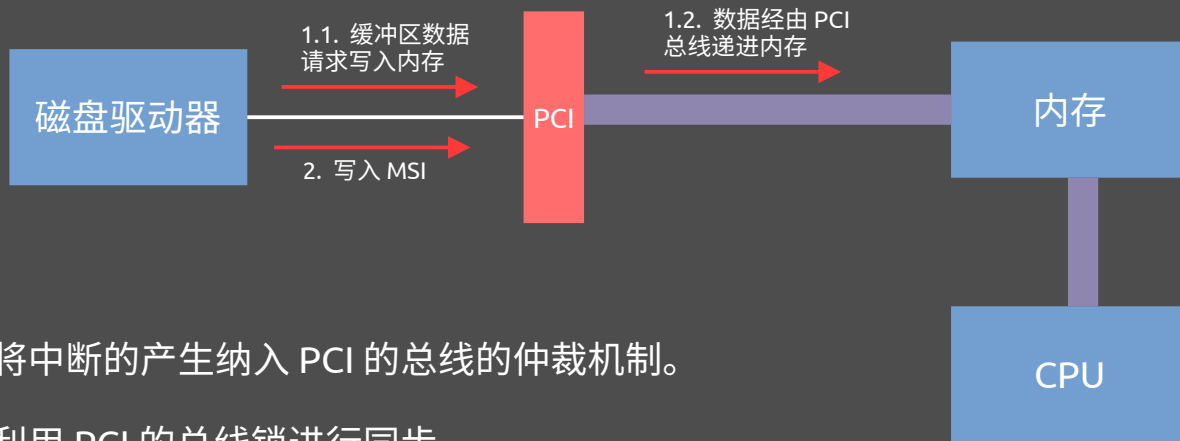
正常执行顺序： 1.1 → 1.2 → 2

如果 PCI 繁忙，决定稍后再写入数据：

1.1 → 2 → 1.1

# 另一种办法：MSI（信息中断）





将中断的产生纳入 PCI 的总线的仲裁机制。

利用 PCI 的总线锁进行同步。

→ 确保了 MSI 只会在数据送达内存后被送出

# 判断设备是否支持 MSI

在 PCI 下面：MSI 为可选功能。在 PCIe 里面则是强制要求。

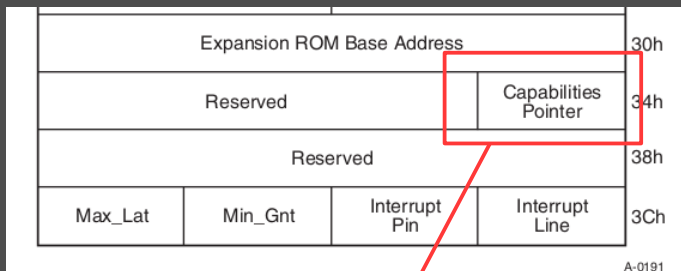
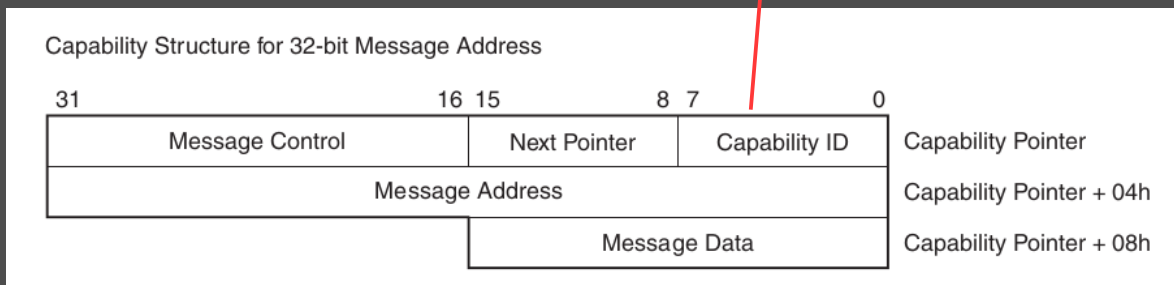


Figure 6-1: Type 00h Configuration Space Header

遍历“能力链表”去查找 MSI。

0x5 : MSI 功能





status 寄存器第五位  
告诉我们该设备是否  
存在能力链表

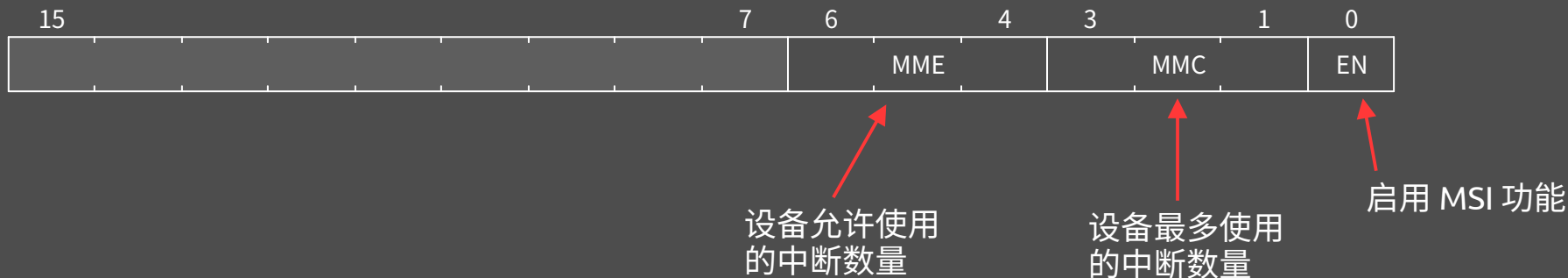
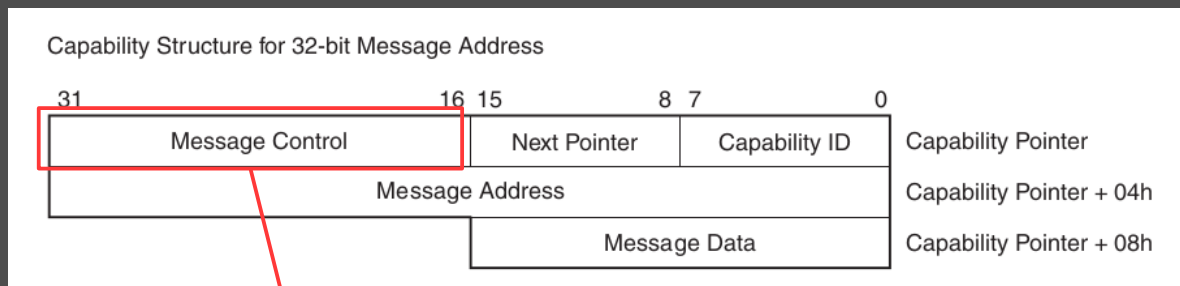
```
pci_reg_t status =  
    pci_read_cspace(device->cspace_base, PCI_REG_STATUS_CMD) >> 16;
```

```
if (!(status & 0x10)) {  
    device->msi_loc = 0;  
    return;  
}
```

```
pci_reg_t cap_ptr = pci_read_cspace(device->cspace_base, 0x34) & 0xff;  
uint32_t cap_hdr;
```

```
while (cap_ptr) {  
    cap_hdr = pci_read_cspace(device->cspace_base, cap_ptr);  
    if ((cap_hdr & 0xff) == 0x5) {  
        // MSI  
        device->msi_loc = cap_ptr;  
        return;  
    }  
    cap_ptr = (cap_hdr >> 8) & 0xff;  
}
```

## 设置一个 MSI：启用 MSI 功能



```
// manipulate the MSI_CTRL to allow device using MSI to request service.  
reg1 = (((reg1 >> 16) & ~0x70) | 0x1) << 16 | (reg1 & 0xffff);  
pci_write_cspace(device->cspace_base, device->msi_loc, reg1);
```



## 设置一个 MSI：让 CPU 识别 MSI

0xFEE00000

```
// Dest: APIC#0, Physical Destination, No redirection
uint32_t msi_addr = (__APIC_BASE_PADDR | 0x8);

// Edge trigger, Fixed delivery
uint32_t msi_data = vector;

pci_write_cspace(
    device->cspace_base, PCI_MSI_ADDR(device->msi_loc), msi_addr);
pci_write_cspace(
    device->cspace_base, PCI_MSI_DATA(device->msi_loc), msi_data & 0xffff);
```



```
void
pci_setup_msi(struct pci_device* device, int vector)
{
    // Dest: APIC#0, Physical Destination, No redirection
    uint32_t msi_addr = (__APIC_BASE_PADDR | 0x8);

    // Edge trigger, Fixed delivery
    uint32_t msi_data = vector;

    pci_write_cspace(
        device->cspace_base, PCI_MSI_ADDR(device->msi_loc), msi_addr);
    pci_write_cspace(
        device->cspace_base, PCI_MSI_DATA(device->msi_loc), msi_data & 0xffff);

    pci_reg_t reg1 = pci_read_cspace(device->cspace_base, device->msi_loc);

    // manipulate the MSI_CTRL to allow device using MSI to request service.
    reg1 = (((reg1 >> 16) & ~0x70) | 0x1) << 16 | (reg1 & 0xffff);
    pci_write_cspace(device->cspace_base, device->msi_loc, reg1);
}
```

# 一个额外的话题：PCIe 的支持

如同在一开始时所说，上述 PCI 的方法在 PCIe 中全部适用，不需要进行任何的更改。

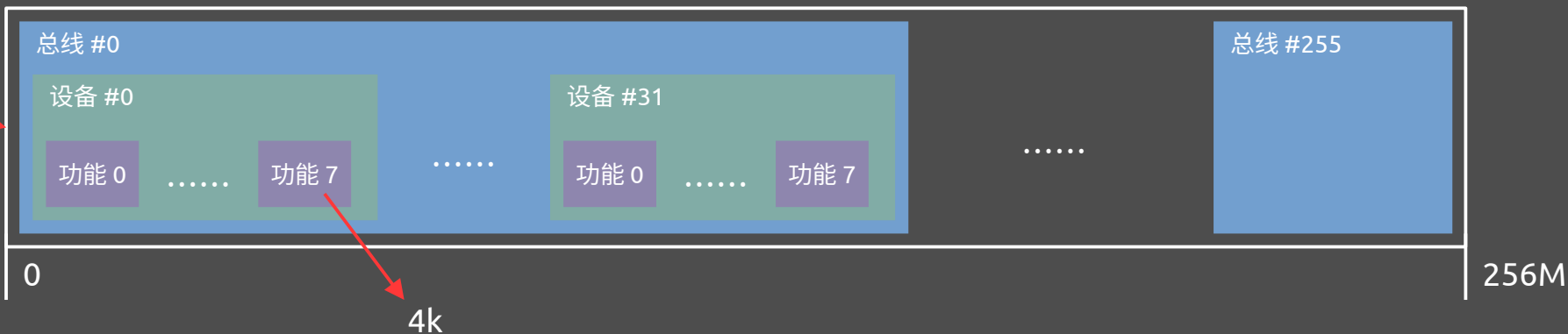
不同的一点：PCIe 的配置空间扩容到 4096 字节。需要使用 PCIe 专用的访问方式去访问。

映射到内存的配置空间

基地址通过查询 ACPI 的 MCFG 表可知

Table 4-2: MCFG Table to Support Enhanced Configuration Space Access

Field	Byte Length	Byte Offset	Description
Header			
Signature	4	0	"MCFG". Signature for the Memory mapped configuration space base address Description Table. (refer to Note 1)
Length	4	4	Length, in bytes, of the entire MCFG Description table including the memory mapped configuration space base address allocation structures.
Revision	1	8	1
Checksum	1	9	Entire table must sum to zero
OEMID	6	10	OEM ID
OEM Table ID	8	16	For the MCFG Description Table, the table ID is the manufacture model ID
OEM Revision	4	24	OEM revision of MCFG table for supplied OEM Table ID
Creator ID	4	28	Vendor ID of utility that created the table
Creator Revision	4	32	Revision of utility that created the table
Reserved	8	36	Reserved
Configuration space base address allocation structure [n]	---	44	A list of the memory mapped configuration base address allocation structures. This list will contain one entry corresponding to each PCI Segment Group present in the platform. The structure of this entry is defined in Table 4-3.



下一步，SATA 驱动!

