# kernel hacker 修炼之道——李万鹏

**男儿立志出乡关， 学不成名死不还。 埋骨何须桑梓地， 人生无处不青山。 ——西乡隆盛诗**

kernel hacker修炼之道之PCI subsystem(四)

# kernel hacker修炼之道之PCI subsystem(四)

## 作者 李万鹏

首先声明代码的版本是linux3.2.7，平台PowerPC_64

第一步分析PCI core对PCI device，PCI bus的遍历过程。PCI subsystem入口函数pcibios_init，在arch/powerpc/kernel/pci_64.c文件中：

```
47 static int __init pcibios_init(void)
48 {
49 struct pci_controller *hose, *tmp;
50
51 printk(KERN_INFO "PCI: Probing PCI hardware\n");
52
53 /* For now, override phys_mem_access_prot. If we need it,g
54 * later, we may move that initialization to each ppc_md
55 */
56 ppc_md.phys_mem_access_prot = pci_phys_mem_access_prot;
57
58 if (pci_probe_only)
59 pci_add_flags(PCI_PROBE_ONLY);
60
61 /* On ppc64, we always enable PCI domains and we keep domain 0
62 * backward compatible in /proc for video cards
63 */
64 pci_add_flags(PCI_ENABLE_PROC_DOMAINS | PCI_COMPAT_DOMAIN_0);
65
66 /* Scan all of the recorded PCI controllers. */
67 list_for_each_entry_safe(hose, tmp, &hose_list, list_node) {
68 pcibios_scan_phb(hose);
69 pci_bus_add_devices(hose->bus);
70 }
71
72 /* Call common code to handle resource allocation */
73 pcibios_resource_survey();
74
75 printk(KERN_DEBUG "PCI: Probing PCI hardware done\n");
76
77 return 0;
78 }
79
```

80subsys_initcall(pcibios_init);

**可以看到这个函数是被subsys_initcall调用的，查看其定义：**

```
#define subsys_initcall(fn)   __define_initcall("4",fn,4)
177#define __define_initcall(level,fn,id) \
178        static initcall_t __initcall_##fn##id __used
179        __attribute__((__section__(".initcall" level
```

在链接脚本vmlinux.lds中有这么一段：

```
    __inicall_start = .;
    .initcall.init : AT(ADDR(.initcall.init) - 0xC0000000) {
    *(.initcall1.init)
    *(.initcall2.init)
    *(.initcall3.init)
    *(.initcall4.init)
    *(.initcall5.init)
    *(.initcall6.init)
    *(.initcall7.init)
    }
    __initcall_end = .;
```

所以，__define_initcall把pcibios_init这个函数指针放到了.initca
start_kernel()->rest_init()->kernel_init()->do_basic_set
pci_probe_only是PowerPC_64下的一个全局变量，如果这个全局变量为1

```
1697void __devinit pcibios_scan_phb(struct pci_controller *hose)
1698{
1699 struct pci_bus *bus;
1700 struct device_node *node = hose->dn;
1701 int mode;
1702
1703 pr_debug("PCI: Scanning PHB %s\n",
1704 node ? node->full_name : "<NO NAME>");
1705
1706 /* Create an empty bus for the toplevel */
1707 bus = pci_create_bus(hose->parent, hose->first_busno, hose->ops, hose);
1708 if (bus == NULL) {
1709 pr_err("Failed to create bus for PCI domain %04x\n",
1710 hose->global_number);
1711 return;
1712 }
1713 bus->secondary = hose->first_busno;
1714 hose->bus = bus;
1715
1716 /* Get some IO space for the new PHB */
1717 pcibios_setup_phb_io_space(hose);
```

```
1718
1719 /* Wire up PHB bus resources */
1720 pcibios_setup_phb_resources(hose);
1721
1722 /* Get probe mode and perform scan */
1723 mode = PCI_PROBE_NORMAL;
1724 if (node && ppc_md.pci_probe_mode)
1725 mode = ppc_md.pci_probe_mode(bus);
1726 pr_debug(" probe mode: %d\n", mode);
1727 if (mode == PCI_PROBE_DEVTREE) {
1728 bus->subordinate = hose->last_busno;
1729 of_scan_bus(node, bus);
1730 }
1731
1732 if (mode == PCI_PROBE_NORMAL)
1733 hose->last_busno = bus->subordinate = pci_scan_child_bus(bus);
1734
1735 /* Configure PCI Express settings */
1736 if (bus && !pci_has_flag(PCI_PROBE_ONLY)) {
1737 struct pci_bus *child;
1738 list_for_each_entry(child, &bus->children, node) {
1739 struct pci_dev *self = child->self;
1740 if (!self)
1741 continue;
1742 pcie_bus_configure_settings(child, self->pcie_mpss);
1743 }
1744 }
1745}
```

首先调用函数pci_create_bus创建toplevel bus，也就是每个pci
domain，phb直接引出的那个bus。这里把这条toplevel bus所在的phb的下游
bus设置为hose->first_busno，其实就是这条bus的number。为新PHB分配一
些空间，设置PHB。如果探测模式是PCI_PROBE_NORMAL，那么由PCI core进行
遍历。调用pci_scan_child_bus函数扫描这个toplevel bus上的所有设备，
以及之下的整个pci sub tree。注意这里的toplevel bus不一定就是0号
bus，因为可能会有几个pci domain。

下面来看toplevel bus是如何创建的：

```
1525struct pci_bus * pci_create_bus(struct device *parent,
1526 int bus, struct pci_ops *ops, void *sysdata)
1527{
1528 int error;
1529 struct pci_bus *b, *b2;
1530 struct device *dev;
1531
1532 b = pci_alloc_bus();
1533 if (!b)
1534 return NULL;
1535
1536 dev = kzalloc(sizeof(*dev), GFP_KERNEL);
1537 if (!dev){
1538 kfree(b);
1539 return NULL;
1540 }
1541
1542 b->sysdata = sysdata;
1543 b->ops = ops;
1544
1545 b2 = pci_find_bus(pci_domain_nr(b), bus);
1546 if (b2) {
1547 /* If we already got to this bus through a different bridge, ignore it */
1548 dev_dbg(&b2->dev, "bus already known\n");
1549 goto err_out;
```

```
1550 }
1551
1552 down_write(&pci_bus_sem);
1553 list_add_tail(&b->node, &pci_root_buses);
1554 up_write(&pci_bus_sem);
1555
1556 dev->parent = parent;
1557 dev->release = pci_release_bus_bridge_dev;
1558 dev_set_name(dev, "pci%04x:%02x", pci_domain_nr(b), bus);
1559 error = device_register(dev);
1560 if (error)
1561 goto dev_reg_err;
1562 b->bridge = get_device(dev);
1563 device_enable_async_suspend(b->bridge);
1564 pci_set_bus_of_node(b);
1565
1566 if (!parent)
1567 set_dev_node(b->bridge, pcibus_to_node(b));
1568
1569 b->dev.class = &pcibus_class;
1570 b->dev.parent = b->bridge;
1571 dev_set_name(&b->dev, "%04x:%02x", pci_domain_nr(b), bus);
1572 error = device_register(&b->dev);
1573 if (error)
1574 goto class_dev_reg_err;
1575
1576 /* Create legacy_io and legacy_mem files for this bus */
1577 pci_create_legacy_files(b);
1578
1579 b->number = b->secondary = bus;
1580 b->resource[0] = &ioport_resource;
1581 b->resource[1] = &iomem_resource;
1582
1583 return b;
1584
1585 class_dev_reg_err:
1586 device_unregister(dev);
1587 dev_reg_err:
1588 down_write(&pci_bus_sem);
1589 list_del(&b->node);
1590 up_write(&pci_bus_sem);
1591 err_out:
1592 kfree(dev);
1593 kfree(b);
1594 return NULL;
1595 }
```

这里的那个bus 是hose->first_busno传过来的，也就是说，已经由firmware分配好的。这里调用pci_find_bus查看一下，如果这个bus已经从不同的bridge分配过了，那么忽略它。否则将这个刚建好的pci_bus添加到pci_root_buses链表上，这个链表是所有toplevel bus的天堂。总线号和总线所在的桥的下游的总线号都是hose->first_buso。toplevel bus的两个指针分别指向了ioport_resousrce和iomem_resource。

PCI core会使用Depth first的方式遍历整个子树：

```
1482 unsigned int __devinit pci_scan_child_bus(struct pci_bus *bus)
1483 {
1484 unsigned int devfn, pass, max = bus->secondary;
1485 struct pci_dev *dev;
1486
1487 dev_dbg(&bus->dev, "scanning bus\n");
1488
```

```
1489 /* Go find them, Rover! */
1490 for (devfn = 0; devfn < 0x100; devfn += 8)
1491 pci_scan_slot(bus, devfn);
1492
1493 /* Reserve buses for SR-IOV capability. */
1494 max += pci_iov_bus_range(bus);
1495
1496 /*
1497 * After performing arch-dependent fixup of the bus, look behind
1498 * all PCI-to-PCI bridges on this bus.
1499 */
1500 if (!bus->is_added) {
1501 dev_dbg(&bus->dev, "fixups for bus\n");
1502 pcibios_fixup_bus(bus);
1503 if (pci_is_root_bus(bus))
1504 bus->is_added = 1;
1505 }
1506
1507 for (pass=0; pass < 2; pass++)
1508 list_for_each_entry(dev, &bus->devices, bus_list) {
1509 if (dev->hdr_type == PCI_HEADER_TYPE_BRIDGE ||
1510 dev->hdr_type == PCI_HEADER_TYPE_CARDBUS)
1511 max = pci_scan_bridge(bus, dev, max, pass);
1512 }
1513
1514 /*
1515 * We've scanned the bus and so we know all about what's on
1516 * the other side of any bridges that may be on this bus plus
1517 * any devices.
1518 *
1519 * Return how far we've got finding sub-buses.
1520 */
1521 dev_dbg(&bus->dev, "bus scan returning with max=%02x\n", max);
1522 return max;
1523}
```

由于一条总线上可以有32个slot，256个logical device，所以这里先扫描这个bus上的所有logical device，如果读取他们的configuration space，填充pci_dev的各个域。然后遍历这个bus上的所有桥，遍历两遍：

第一次： 一旦遇到BIOS未初始化的PCI桥，既向上返回一层，目的是统计BIOS分配到的最大总线号；

第二次： 一旦遇到BIOS已初始化的PCI桥，既向下深入一层，目的是继续分配总线号；

这里的具体细节会在下一篇blog中介绍，下面来看负责扫描slot的函数pci_scan_slot：

```
1295int pci_scan_slot(struct pci_bus *bus, int devfn)
1296{
1297 unsigned fn, nr = 0;
1298 struct pci_dev *dev;
1299 unsigned (*next_fn)(struct pci_dev *, unsigned) = no_next_fn;
1300
1301 if (only_one_child(bus) && (devfn > 0))
1302 return 0; /* Already scanned the entire slot */
1303
1304 dev = pci_scan_single_device(bus, devfn);
```

```
1305 if (!dev)
1306 return 0;
1307 if (!dev->is_added)
1308 nr++;
1309
1310 if (pci_ari_enabled(bus))
1311 next_fn = next_ari_fn;
1312 else if (dev->multifunction)
1313 next_fn = next_trad_fn;
1314
1315 for (fn = next_fn(dev, 0); fn > 0; fn = next_fn(dev, fn)) {
1316 dev = pci_scan_single_device(bus, devfn + fn);
1317 if (dev) {
1318 if (!dev->is_added)
1319 nr++;
1320 dev->multifunction = 1;
1321 }
1322 }
1323
1324 /* only one slot has pcie device */
1325 if (bus->self && nr)
1326 pcie_aspm_init_link_state(bus->self);
1327
1328 return nr;
1329 }
```

这里调用pci_scan_single_device扫描单个设备：

```
1225 struct pci_dev *__ref pci_scan_single_device(struct pci_bus *bus, int devfn)
1226 {
1227 struct pci_dev *dev;
1228
1229 dev = pci_get_slot(bus, devfn);
1230 if (dev) {
1231 pci_dev_put(dev);
1232 return dev;
1233 }
1234
1235 dev = pci_scan_device(bus, devfn);
1236 if (!dev)
1237 return NULL;
1238
1239 pci_device_add(dev, bus);
1240
1241 return dev;
1242 }
```

通过bus和device:function号获得pci_dev，然后扫描这个logical device，主要就是获得厂商ID和设备ID，然后通过device的bus_list bit field把pci_dev添加到bus的devices链表上。

```
1120 static struct pci_dev *pci_scan_device(struct pci_bus *bus, int devfn)
1121 {
1122 struct pci_dev *dev;
1123 u32 l;
1124 int delay = 1;
1125
1126 if (pci_bus_read_config_dword(bus, devfn, PCI_VENDOR_ID, &l))
1127 return NULL;
1128
1129 /* some broken boards return 0 or ~0 if a slot is empty: */
1130 if (l == 0xffffffff || l == 0x00000000 ||
1131 l == 0x0000ffff || l == 0xffff0000)
1132 return NULL;
1133
1134 /* Configuration request Retry Status */
1135 while (l == 0xffff0001) {
1136 msleep(delay);
1137 delay *= 2;
1138 if (pci_bus_read_config_dword(bus, devfn, PCI_VENDOR_ID, &l))
```

```
1139 return NULL;
1140 /* Card hasn't responded in 60 seconds? Must be stuck. */
1141 if (delay > 60 * 1000) {
1142 printk(KERN_WARNING "pci %04x:%02x:%02x.%d: not "
1143 "responding\n", pci_domain_nr(bus),
1144 bus->number, PCI_SLOT(devfn),
1145 PCI_FUNC(devfn));
1146 return NULL;
1147 }
1148 }
1149
1150 dev = alloc_pci_dev();
1151 if (!dev)
1152 return NULL;
1153
1154 dev->bus = bus;
1155 dev->devfn = devfn;
1156 dev->vendor = l & 0xffff;
1157 dev->device = (l >> 16) & 0xffff;
1158
1159 pci_set_of_node(dev);
1160
1161 if (pci_setup_device(dev)) {
1162 kfree(dev);
1163 return NULL;
1164 }
1165
1166 return dev;
1167 }
```

只要不是全1或全0就算有效的厂商ID，设备ID，如果读出的是 0xffff0001，则重复读，如果超过60s，不再尝试。否则调用 pci_setup_device函数，这个函数主要读取configuration space，然后设置 device的相应域。

```
892 int pci_setup_device(struct pci_dev *dev)
893 {
894 u32 class;
895 u8 hdr_type;
896 struct pci_slot *slot;
897 int pos = 0;
898
899 if (pci_read_config_byte(dev, PCI_HEADER_TYPE, &hdr_type))
900 return -EIO;
901
902 dev->sysdata = dev->bus->sysdata;
903 dev->dev.parent = dev->bus->bridge;
904 dev->dev.bus = &pci_bus_type;
905 dev->hdr_type = hdr_type & 0x7f;
906 dev->multifunction = !!(hdr_type & 0x80);
907 dev->error_state = pci_channel_io_normal;
908 set_pcie_port_type(dev);
909
910 list_for_each_entry(slot, &dev->bus->slots, list)
911 if (PCI_SLOT(dev->devfn) == slot->number)
912 dev->slot = slot;
913
914 /* Assume 32-bit PCI; let 64-bit PCI cards (which are far rarer)
915 set this higher, assuming the system even supports it. */
916 dev->dma_mask = 0xffffffff;
917
918 dev_set_name(&dev->dev, "%04x:%02x:%02x.%d", pci_domain_nr(dev->bus),
919 dev->bus->number, PCI_SLOT(dev->devfn),
920 PCI_FUNC(dev->devfn));
921
922 pci_read_config_dword(dev, PCI_CLASS_REVISION, &class);
923 dev->revision = class & 0xff;
```

```
924 class >>= 8; /* upper 3 bytes */
925 dev->class = class;
926 class >>= 8;
927
928 dev_printk(KERN_DEBUG, &dev->dev, "[%04x:%04x] type %d class %#08x\n",
929 dev->vendor, dev->device, dev->hdr_type, class);
930
931 /* need to have dev->class ready */
932 dev->cfg_size = pci_cfg_space_size(dev);
933
934 /* "Unknown power state" */
935 dev->current_state = PCI_UNKNOWN;
936
937 /* Early fixups, before probing the BARs */
938 pci_fixup_device(pci_fixup_early, dev);
939 /* device class may be changed after fixup */
940 class = dev->class >> 8;
941
942 switch (dev->hdr_type) { /* header type */
943 case PCI_HEADER_TYPE_NORMAL: /* standard header */
944 if (class == PCI_CLASS_BRIDGE_PCI)
945 goto bad;
946 pci_read_irq(dev);
947 pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
948 pci_read_config_word(dev, PCI_SUBSYSTEM_VENDOR_ID, &dev->subsystem_vendor);
949 pci_read_config_word(dev, PCI_SUBSYSTEM_ID, &dev->subsystem_device);
950
951 /*
952 * Do the ugly legacy mode stuff here rather than broken chip
953 * quirk code. Legacy mode ATA controllers have fixed
954 * addresses. These are not always echoed in BAR0-3, and
955 * BAR0-3 in a few cases contain junk!
956 */
957 if (class == PCI_CLASS_STORAGE_IDE) {
958 u8 progif;
959 pci_read_config_byte(dev, PCI_CLASS_PROG, &progif);
960 if ((progif & 1) == 0) {
961 dev->resource[0].start = 0x1F0;
962 dev->resource[0].end = 0x1F7;
963 dev->resource[0].flags = LEGACY_IO_RESOURCE;
964 dev->resource[1].start = 0x3F6;
965 dev->resource[1].end = 0x3F6;
966 dev->resource[1].flags = LEGACY_IO_RESOURCE;
967 }
968 if ((progif & 4) == 0) {
969 dev->resource[2].start = 0x170;
970 dev->resource[2].end = 0x177;
971 dev->resource[2].flags = LEGACY_IO_RESOURCE;
972 dev->resource[3].start = 0x376;
973 dev->resource[3].end = 0x376;
974 dev->resource[3].flags = LEGACY_IO_RESOURCE;
975 }
976 }
977 break;
978
979 case PCI_HEADER_TYPE_BRIDGE: /* bridge header */
980 if (class != PCI_CLASS_BRIDGE_PCI)
981 goto bad;
982 /* The PCI-to-PCI bridge spec requires that subtractive
983 decoding (i.e. transparent) bridge must have programming
984 interface code of 0x01. */
985 pci_read_irq(dev);
986 dev->transparent = ((dev->class & 0xff) == 1);
987 pci_read_bases(dev, 2, PCI_ROM_ADDRESS1);
988 set_pcie_hotplug_bridge(dev);
989 pos = pci_find_capability(dev, PCI_CAP_ID_SSVID);
```

```
990 if (pos) {
991 pci_read_config_word(dev, pos + PCI_SSVID_VENDOR_ID, &dev->subsystem_vendor);
992 pci_read_config_word(dev, pos + PCI_SSVID_DEVICE_ID, &dev->subsystem_device);
993 }
994 break;
995
996 case PCI_HEADER_TYPE_CARDBUS: /* CardBus bridge header */
997 if (class != PCI_CLASS_BRIDGE_CARDBUS)
998 goto bad;
999 pci_read_irq(dev);
1000 pci_read_bases(dev, 1, 0);
1001 pci_read_config_word(dev, PCI_CB_SUBSYSTEM_VENDOR_ID, &dev->subsystem_vendor);
1002 pci_read_config_word(dev, PCI_CB_SUBSYSTEM_ID, &dev->subsystem_device);
1003 break;
1004
1005 default: /* unknown header */
1006 dev_err(&dev->dev, "unknown header type %02x, "
1007 "ignoring device\n", dev->hdr_type);
1008 return -EIO;
1009
1010 bad:
1011 dev_err(&dev->dev, "ignoring class %02x (doesn't match header "
1012 "type %02x)\n", class, dev->hdr_type);
1013 dev->class = PCI_CLASS_NOT_DEFINED;
1014 }
1015
1016 /* We found a fine healthy device, go go go... */
1017 return 0;
1018}
```
这里注意Header Type Register是8位的，低7位指明header type，最高1位指明这是多功能设备还是单功能设备，当最高位为1时是多功能设备，当最高位为0时是单功能设备。所以这里dev->hdr_type=hdr_type & 0x7f获得低7位，dev->multifunction=!!(hdr_type & 08)获得最高位。目前header type头主要有三种，可以从include/linux/pci_regs.h：

```
69#define PCI_HEADER_TYPE_NORMAL 0
70#define PCI_HEADER_TYPE_BRIDGE 1
71#define PCI_HEADER_TYPE_CARDBUS 2
```

Revision ID register是8位的寄存器，Class Code Register是24位的寄存器，这里一次读取了长字，其中最低的8位是revision id，高24位是class code，这24位又分成3部分，class：sub-class：programming interface。假如读出的数是0x060401，那么06表示class(bridge device)，04表示sub-class(PCI-to-PCI bridge)，01是programming interface。其中head type 0 对应normal device，head type 1对应pci-to-pci bridge。注意如果这个不匹配回答因警告信息：

```
1011 dev_err(&dev->dev, "ignoring class %02x (doesn't match header "
1012 "type %02x)\n", class, dev->hdr_type);
```
跟桥有关的类和子类如下，这些都是在文件include/linux/pci_ids.h中定义的：

```
50#define PCI_BASE_CLASS_BRIDGE 0x06
51#define PCI_CLASS_BRIDGE_HOST 0x0600
52#define PCI_CLASS_BRIDGE_ISA 0x0601
53#define PCI_CLASS_BRIDGE_EISA 0x0602
54#define PCI_CLASS_BRIDGE_MC 0x0603
55#define PCI_CLASS_BRIDGE_PCI 0x0604
56#define PCI_CLASS_BRIDGE_PCMCIA 0x0605
57#define PCI_CLASS_BRIDGE_NUBUS 0x0606
58#define PCI_CLASS_BRIDGE_CARDBUS 0x0607
59#define PCI_CLASS_BRIDGE_RACEWAY 0x0608
60#define PCI_CLASS_BRIDGE_OTHER 0x0680
```
注意上边的type1与class：pci-to-pci bridge

匹配，与其他的不匹配，不要看错了。