# rk3399 PCIe rc device enumeration device resource identification analysis

Category column:  PCI/PCIe topic    Article tags:   driver development     PCIe RC     PCIe host

> **PCI/PCIe topic**  The column includes this content
>
> 7 subscribe    4 articles     subscribe
>

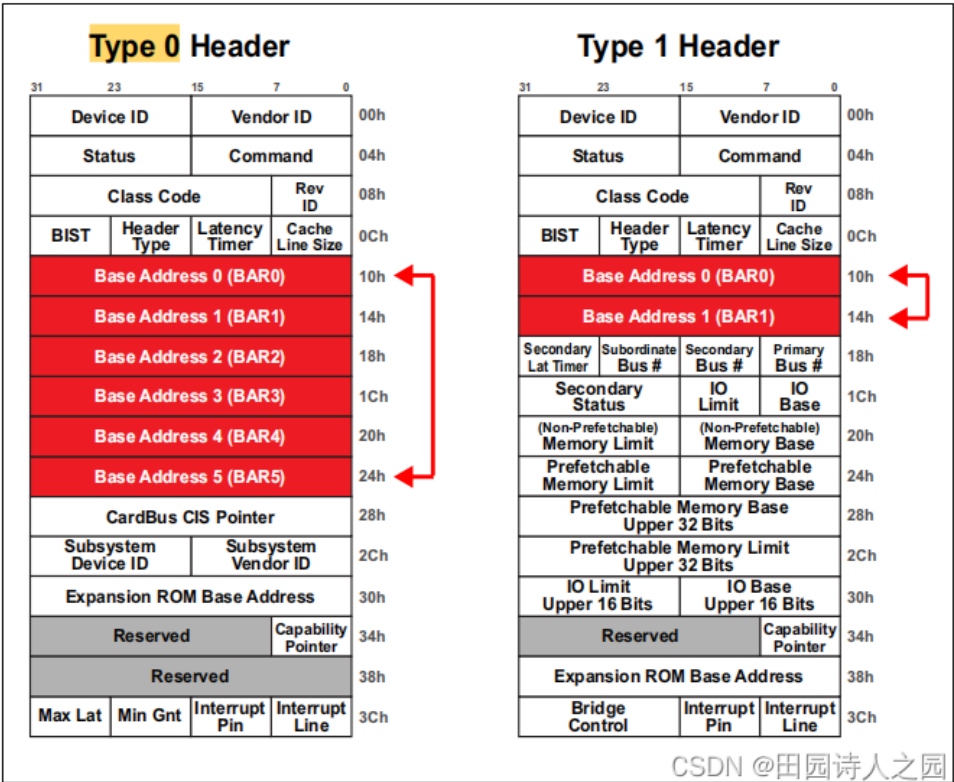**rk3399 PCIe rc device enumeration device resource identification analysis**

Based on the learning summary of Mr. Wei Dongshan's PCIe learning topic

## 1. Device configuration space

Figure 4-2 in "PCI Express Technology 3.0.pdf" shows the configuration space diagram of PCIe devices and switch (or bridge);



Figure 4-2: BARs in Configuration Space

The purpose of using PCI/PCIe is to simply access it: read and write PCI/PCIe devices like reading and writing memory.

Ask:

- Which addresses are used to read and write to the device?

- What is the range of these addresses?
- Is it accessed like memory, or is it accessed like IO?

Each PCI/PCIe device has a configuration space, which is a series of registers. For ordinary devices, its configuration space includes:

- Device ID
- Vendor ID
- Class Code: What kind of equipment? Storage device? display screen? wait
- 6 Base Address Registers:

## 1.1Device information

- Vendor ID: Vendor ID, the PCI SIG organization assigns a unique ID to each vendor

- Device ID: The manufacturer assigns a Device ID to a certain type of product

- Revision ID: The version number customized by the manufacturer, which can be considered as an extension of Device ID

- Header Type:

  - bit[7]: 1 - it is a multifunction device ("multi-function"), 0 - it is a single function device ("single-function")
  - bit[6:0]: 00h-normal equipment, 01h-switch or bridge equipment, this value also determines the meaning of the beginning of the offset address 10h in the configuration space
- Class Code: This is a read-only register , it contains 3 bytes, used to indicate the function of the device, it is divided into 3 parts

  - Highest byte: Indicates "base class", which is used to indicate that it belongs to the memory card, graphics card, etc.
  - Middle byte: Indicates "sub-class", subdivide the category
  - Lowest byte: used to represent the programming interface "Interface" at the register level
  - The example is as follows: When the Base Class is 01h, it means that it is a storage device, but sub-class and Interface subdivisions can also be used

| Base Class | Sub-Class | Interface | Meaning |
|---|---|---|---|
| 01h | 00h | 00h | SCSI bus controller |
| | 01h | xxh | IDE controller (see Note 1) |
| | 02h | 00h | Floppy disk controller |
| | 03h | 00h | IPI bus controller |
| | 04h | 00h | RAID controller |
| | 05h | 20h | ATA controller with ADMA interface– single stepping (see Note 2) |
| | | 30h | ATA controller with ADMA interface– continuous operation (see Note 2) |
| | 06h | 00h | Serial ATA controller–vendor specific interface |
| | | 01h | Serial ATA controller–AHCI 1.0 interface |
| | 07h | 00h | Serial Attached SCSI (SAS) controller |
| | 80h | 00h | Other mass storage controller |

## 1.2 Base Address

Ordinary PCI/PCIe devices have 6 base address registers, referred to as BAR:
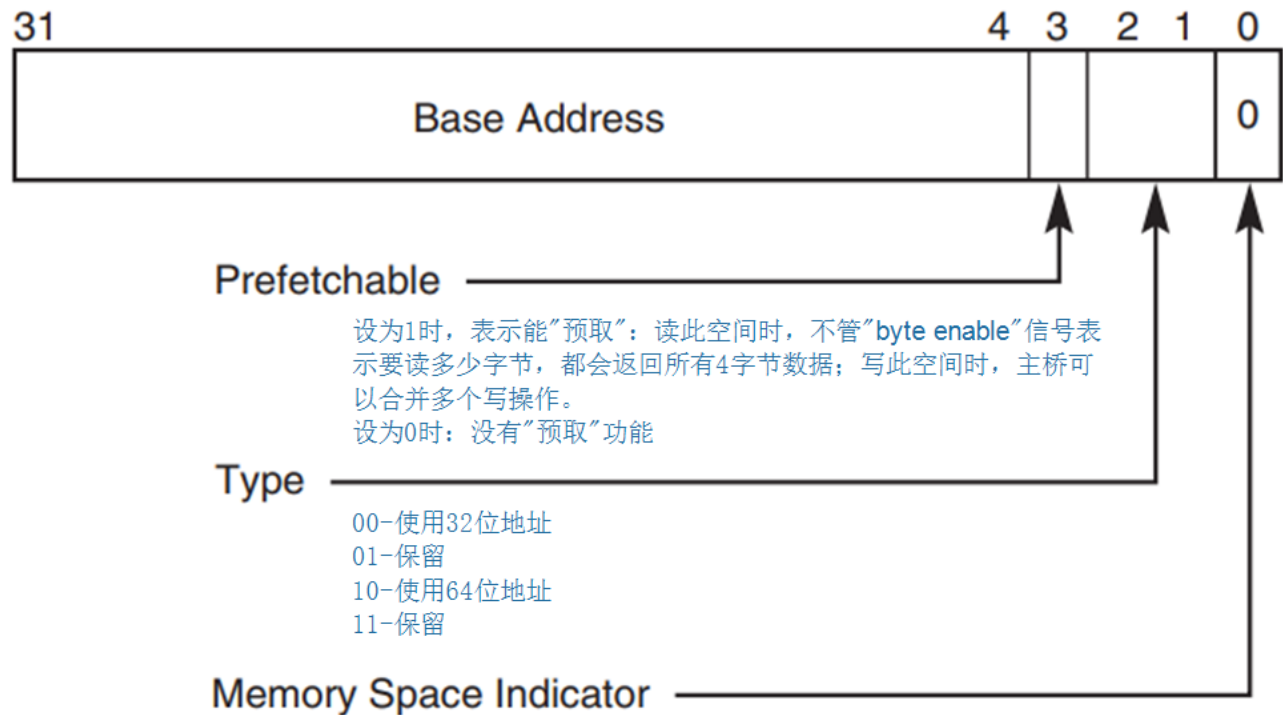
BARs are used to:

- What type of space is required for the statement: memory (32-bit address or 64-bit address), IO
- How much space is required to declare
- Save the PCI space base address assigned to it by the master

The address space can be divided into two categories: memory (Memory), IO:

- For memory, whatever value is written is what value is read out, which can be read in advance
- For IO, it reflects the current state of the hardware, and the value read at each moment is not necessarily the same

The format of BAR is as follows:

- for memory space



用于内存空间的**BAR**

- For IO space:



用于**IO**空间的**BAR**

How does BAR indicate how much space it wants to apply for? Take a 32-bit address as an example:

- Software writes 0xFFFFFFFF to BAR
- Software to read BAR
- The value read is assumed to be 0xFFF0,000?, and the lowest 4 bits are ignored, and the result is: 0xFFF0,0000
  - This means that the "Base Address" that can be written in BAR is only the highest 12 bits
  - It also means that the lowest 20 bits are a variable range, so the size of this space is 2^20=1M Byte The code example is as follows:
    pos indicates the bar register to be accessed

```
1        u32 l = 0, sz = 0, mask;
2        ...
3
```

```
4          mask = type ? PCI_ROM_ADDRESS_MASK : ~0;
5          ...
6          pci_read_config_dword(dev, pos, &l);
7          pci_write_config_dword(dev, pos, l | mask);
8          pci_read_config_dword(dev, pos, &sz);
           pci_write_config_dword(dev, pos, l);
```

If BAR indicates that it uses a 32-bit address, then BAR0~BAR5 can represent 6 address spaces respectively.

If BAR indicates that it uses a 64-bit address, then BAR0 and BAR1, BAR2 and BAR3, BAR4 and BAR5 represent three address spaces:

- The low-order BAR represents the lower 32 bits of the 64-bit address
- A high-order address represents the upper 32 bits of a 64-bit address.

## 2. The process of scanning the device

### 2.1 Core: Construct pci_dev

Scan the PCIe bus, and construct a corresponding pci_dev for each PCIe bridge and PCIe device:

- Fill in the members of pci_dev, such as VID, PID, Class, etc.
- Allocate address space, write to PCIe device

The pci_dev structure is as follows:

```
255: /*
256:  * The pci_dev structure is used to describe PCI devices.
257:  */
258: struct pci_dev {
259:     struct list_head bus_list;  /* node in per-bus list */
260:     struct pci_bus  *bus;        /* bus this device is on */
261:     struct pci_bus  *subordinate;   /* bus this device bridges to */
262:
263:     void        *sysdata;   /* hook for sys-specific extension */
264:     struct proc_dir_entry *procent; /* device entry in /proc/bus/pci */
265:     struct pci_slot *slot;      /* Physical slot this device is in */
266:
267:     unsigned int    devfn;       /* encoded device & function index */
268:     unsigned short  vendor;
269:     unsigned short  device;
270:     unsigned short  subsystem_vendor;       1.设备信息
271:     unsigned short  subsystem_device;
272:     unsigned int    class;       /* 3 bytes: (base,sub,prog-if) */
273:     u8      revision;   /* PCI revision, low byte of class word */
274:     u8      hdr_type;   /* PCI header type (`multi' flag masked out) */
275:
276:     /* 省略 */
277:
278:     /*
279:      * Instead of touching interrupt line and base address registers
280:      * directly, use the values stored here. They might be different!
281:      */
282:     unsigned int    irq;                     2.用到的资源
283:     struct resource resource[DEVICE_COUNT_RESOURCE]; /* I/O and memory
284:
```
CSDN @田园诗人之园

Corresponding to the device information in the pci_dev structure: it can be obtained by reading the configuration space of the PCI device.

Corresponding to the resources in the pci_dev structure, this course will not analyze irq first. For the resource structure, each member corresponds to a BAR.

The resource structure is as follows. It should be noted that the start, end, etc. recorded in it are viewed from the perspective of the CPU. That is to say, if the memory address and IO address are recorded, it is the CPU address, not the PCI address. And these addresses are physical addresses, and to use them in software, ioremap must be executed first.

```
18:  struct resource {
19:      resource_size_t start;
20:      resource_size_t end;
21:      const char *name;
22:      unsigned long flags;
23:      struct resource *parent, *sibling, *child;
24:  };
25:
```

CSDN @田园诗人之园

## 2.2 Code Analysis

We need to find these 4 core codes:

- assign pci_dev
- Read the configuration space of the PCIe device and fill in the device information in pci_dev
- According to the BAR of the PCIe device, know what type of address it wants to apply for and how big it is
- Assign address, write to BAR

The key code is divided into two parts:

- Read the information and know how much space the PCIe device wants to apply for

```
1   rockchip_pcie_probe
2       bus = pci_scan_root_bus(&pdev->dev, 0, &rockchip_pcie_ops, rockchip, &res);
3                   pci_scan_root_bus_msi
4               pci_scan_child_bus
5                   pci_scan_slot
6                       dev = pci_scan_single_device(bus, devfn);
7                                       dev = pci_scan_device(bus, devfn);
8                                           struct pci_dev *dev;
9                                           dev = pci_alloc_dev(bus);
10                                          pci_setup_device
11                          pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
12                                          struct resource *res = &dev->resource[pos];
13                                          __pci_read_base
14                                                      pci_read_config_dword(dev
15                                                      pci_write_config_dword(de
16                                                      pci_read_config_dword(dev
17                                                      pci_write_config_dword(de
18                  pci_device_add(dev, bus);
```

- allocate space

```
1   rockchip_pcie_probe
2           pci_bus_size_bridges(bus);
3           pci_bus_assign_resources(bus);
4                   __pci_bus_assign_resources
5               pbus_assign_resources_sorted
6                   /* pci_dev->resource[]里记录有想申请的资源的大小,
7                    * 把这些资源按对齐的要求排序
8                    * 比如资源A要求1K地址对齐, 资源B要求32地址对齐
9                    * 那么资源A排在资源B前面, 优先分配资源A
10                   */
11                   list_for_each_entry(dev, &bus->devices, bus_list)
12                       __dev_sort_resources(dev, &head);
13                               // 分配资源
14                               __assign_resources_sorted
15                   assign_requested_resources_sorted(head, &local_fail_head);
```

### 2.2.1 Allocating the pci_dev structure

```
rockchip_pcie_probe
    bus = pci_scan_root_bus(&pdev->dev, 0, &rockchip_pcie_ops, rockchip, &res);
        pci_scan_root_bus_msi
            pci_scan_child_bus
                pci_scan_slot
                    dev = pci_scan_single_device(bus, devfn);
                        dev = pci_scan_device(bus, devfn);
                            struct pci_dev *dev;
                            dev = pci_alloc_dev(bus);
                            pci_setup_device
                                pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
                    pci_device_add(dev, bus);                    CSDN @田园诗人之园
```

**2.2.2 Read device information**

```
rockchip_pcie_probe
    bus = pci_scan_root_bus(&pdev->dev, 0, &rockchip_pcie_ops, rockchip, &res);
        pci_scan_root_bus_msi
            pci_scan_child_bus
                pci_scan_slot
                    dev = pci_scan_single_device(bus, devfn);
                        dev = pci_scan_device(bus, devfn);
                            struct pci_dev *dev;
                            dev = pci_alloc_dev(bus);
                            pci_setup_device
                                pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
                    pci_device_add(dev, bus);                    CSDN @田园诗人之园
```

In the `pci_scan_device` function, it will try to read the VID and PID first, and then continue to call if it succeeds `pci_setup_device`:

```
1632: static struct pci_dev *pci_scan_device(struct pci_bus *bus, int devfn)
1633: {
1634:     struct pci_dev *dev;
1635:     u32 l;
1636:
1637:     if (!pci_bus_read_dev_vendor_id(bus, devfn, &l, 60*1000))
1638:         return NULL;
1639:                                                          1.先读取VID、PID
1640:     dev = pci_alloc_dev(bus);
1641:     if (!dev)
1642:         return NULL;
1643:
1644:     dev->devfn = devfn;
1645:     dev->vendor = l & 0xffff;                2.成功的话，
1646:     dev->device = (l >> 16) & 0xffff;           分配pci_dev并记录VID、PID
1647:
1648:     pci_set_of_node(dev);
1649:
1650:     if (pci_setup_device(dev)) {      3.进一步设置pci_dev
1651:         pci_bus_put(dev->bus);
1652:         kfree(dev);
1653:         return NULL;
1654:     }
1655:
1656:     return dev;
1657: } « end pci_scan_device »
1658:
```

Internally `pci_setup_device`, other information continues to be read:

```
1180: int pci_setup_device(struct pci_dev *dev)
1181: {
1182:     u32 class;
1183:     u16 cmd;
1184:     u8 hdr_type;
1185:     int pos = 0;
1186:     struct pci_bus_region region;
1187:     struct resource *res;
1188:                                         1.读取Head Type，分辨是桥还是普通设备
1189:     if (pci_read_config_byte(dev, PCI_HEADER_TYPE, &hdr_type))
1190:         return -EIO;
1191:
1192:     dev->sysdata = dev->bus->sysdata;
1193:     dev->dev.parent = dev->bus->bridge;
1194:     dev->dev.bus = &pci_bus_type;
1195:     dev->hdr_type = hdr_type & 0x7f;
1196:     dev->multifunction = !!(hdr_type & 0x80);
1197:     dev->error_state = pci_channel_io_normal;
1198:     set_pcie_port_type(dev);
1199:
1200:     pci_dev_assign_slot(dev);
1201:     /* Assume 32-bit PCI; let 64-bit PCI cards (which are far rarer)
1202:        set this higher, assuming the system even supports it.  */
1203:     dev->dma_mask = 0xffffffff;
1204:
1205:     dev_set_name(&dev->dev, "%04x:%02x:%02x.%d", pci_domain_nr(dev->bus),
1206:             dev->bus->number, PCI_SLOT(dev->devfn),
1207:             PCI_FUNC(dev->devfn));          2.读取Class Code
1208:
1209:     pci_read_config_dword(dev, PCI_CLASS_REVISION, &class);
1210:     dev->revision = class & 0xff;
1211:     dev->class = class >> 8;               /* upper 3 bytes */
1212:
```
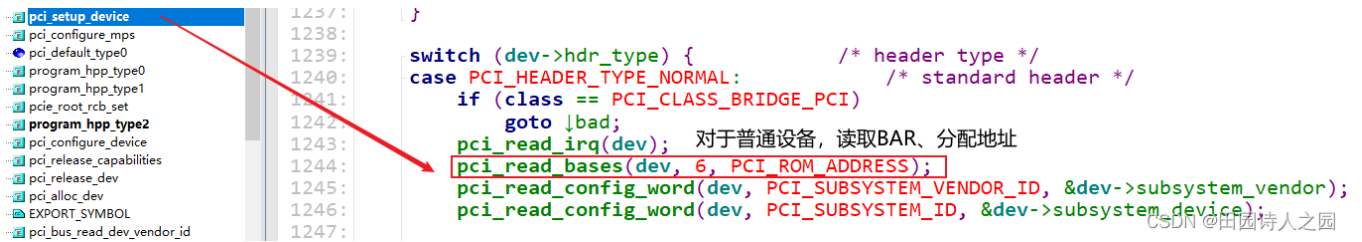
**2.2.3 Read BAR**

```
pci_setup_device          1237:    }
pci_configure_mps         1238:
pci_default_type0         1239:    switch (dev->hdr_type) {              /* header type */
program_hpp_type0         1240:    case PCI_HEADER_TYPE_NORMAL:         /* standard header */
program_hpp_type1         1241:        if (class == PCI_CLASS_BRIDGE_PCI)
pcie_root_rcb_set         1242:            goto ↓bad;
program_hpp_type2         1243:        pci_read_irq(dev);   对于普通设备，读取BAR、分配地址
pci_configure_device      1244:        pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
pci_release_capabilities  1245:        pci_read_config_word(dev, PCI_SUBSYSTEM_VENDOR_ID, &dev->subsystem_vendor);
pci_release_dev           1246:        pci_read_config_word(dev, PCI_SUBSYSTEM_ID, &dev->subsystem_device);
pci_alloc_dev             1247:
EXPORT_SYMBOL
pci_bus_read_dev_vendor_id
```

```
rockchip_pcie_probe
    bus = pci_scan_root_bus(&pdev->dev, 0, &rockchip_pcie_ops, rockchip, &res);
        pci_scan_root_bus_msi
            pci_scan_child_bus
                pci_scan_slot
                    dev = pci_scan_single_device(bus, devfn);
                        dev = pci_scan_device(bus, devfn);
                            struct pci_dev *dev;
                            dev = pci_alloc_dev(bus);
                            pci_setup_device
                                pci_read_bases(dev, 6, PCI_ROM_ADDRESS);
                        pci_device_add(dev, bus);
```

pci_read_bases Function code analysis:

```
319: static void pci_read_bases(struct pci_dev *dev, unsigned int howmany, int rom)
320: {
321:     unsigned int pos, reg;
322:
323:     if (dev->non_compliant_bars)
324:         return;
325:
326:     for (pos = 0; pos < howmany; pos++) {            根据BAR设置pci_dev里的resource
327:         struct resource *res = &dev->resource[pos];
328:         reg = PCI_BASE_ADDRESS_0 + (pos << 2);
329:         pos += __pci_read_base(dev, pci_bar_unknown, res, reg);
330:     }
331:
```

pci_read_bases The function will be called again __pci_read_base , __pci_read_base just read the BAR and calculate the size of the space you want to apply for:

- Read BAR, retain the original value
- Write 0xFFFFFFFF to BAR
- After reading it out, parse out the required address space size and record it in pci_dev->resource[ ]
    - pci_dev->resource[ ].start = 0;
    - pci_dev->resource[ ].end = size - 1;

Posting what I said earlier will help you understand the code:

How does BAR indicate how much space it wants to apply for? Take a 32-bit address as an example:

- Software writes 0xFFFFFFFF to BAR
- Software to read BAR
- The value read is assumed to be 0xFFF0,000?, and the lowest 4 bits are ignored, and the result is: 0xFFF0,0000
    - This means that the "Base Address" that can be written in BAR is only the highest 12 bits
    - It also means that the lowest 20 bits can be changed, so the size of this space is 2^20=1M Byte

Below is __pci_read_bases the code analysis of the function.

- Get the size (original data, further analysis is required): For example, in the following code, sz is assigned a value of 0xFFF0,000?, which requires further analysis

```
174: int __pci_read_base(struct pci_dev *dev, enum pci_bar_type type,
175:                 struct resource *res, unsigned int pos)
176: {
177:     u32 l, sz, mask;
178:     u64 l64, sz64, mask64;
179:     u16 orig_cmd;
180:     struct pci_bus_region region, inverted_region;
181:
182:     mask = type ? PCI_ROM_ADDRESS_MASK : ~0;        1. mask = 0xFFFFFFFF
183:
184:     /* No printks while decoding is disabled! */
185:     if (!dev->mmio_always_on) {
186:         pci_read_config_word(dev, PCI_COMMAND, &orig_cmd);
187:         if (orig_cmd & PCI_COMMAND_DECODE_ENABLE) {
188:             pci_write_config_word(dev, PCI_COMMAND,
189:                 orig_cmd & ~PCI_COMMAND_DECODE_ENABLE);
190:         }
191:     }
192:
193:     res->name = pci_name(dev);
194:
195:     pci_read_config_dword(dev, pos, &l);          2. 读出原值, l=原值
196:     pci_write_config_dword(dev, pos, l | mask);   3. 写入0xffffffff
197:     pci_read_config_dword(dev, pos, &sz);         4. 再读出，表示大小(后面解析)
198:     pci_write_config_dword(dev, pos, l);          5. 写入原值
199:
200:     /* 省略 */   6.解析出l64, sz64: l64=0, sz64为大小
201:
202:     region.start = l64;                           7.设置pci_dev->resource[ ]
203:     region.end = l64 + sz64;                        .start = 0
204:                                                     .size  = 大小 - 1
205:     pcibios_bus_to_resource(dev->bus, res, &region);
206:     pcibios_resource_to_bus(dev->bus, &inverted_region, res);
```

### 2.2.4 Allocating Address Space

The function call in this part of the code is very deep, we can grasp two problems:

- Where is the address space allocated from?
  - In the device tree, the corresponding relationship between CPU address and PCI address is specified, which are recorded in pci_bus as "resources"
  - When reading BAR, it records the size of the space it wants to apply for in pci_dev->resource[]
- The allocated base address, to be written into BAR

The code calling relationship is as follows:

- Sort the resources to be applied for according to the alignment requirements, and then call assign_requested_resources_sorted, the code is as follows:

```
1
2   /* 把要申请的资源，按照对齐要求排序
3    * 然后调用assign_requested_resources_sorted
4    */
5
6   rockchip_pcie_probe
7         pci_bus_size_bridges(bus);
8         pci_bus_assign_resources(bus);
9             __pci_bus_assign_resources
10            pbus_assign_resources_sorted
11                /* pci_dev->resource[]里记录有想申请的资源的大小,
12                 * 把这些资源按对齐的要求排序
13                 * 比如资源A要求1K地址对齐，资源B要求32地址对齐
14                 * 那么资源A排在资源B前面，优先分配资源A
15                 */
16                list_for_each_entry(dev, &bus->devices, bus_list)
17                    __dev_sort_resources(dev, &head);
```

```
18              // 分配资源
19              __assign_resources_sorted
          assign_requested_resources_sorted(head, &local_fail_head);
```

- assign_requested_resources_sorted function does two things

  - allocate address space

  - Write the PCI address corresponding to this space into the BAR of the PCIe device

  - code show as below:

```
1   assign_requested_resources_sorted(head, &local_fail_head);
2       pci_assign_resource
3           ret = _pci_assign_resource(dev, resno, size, align);
4                   // 分配地址空间
5               __pci_assign_resource
6                   pci_bus_alloc_resource
7                       pci_bus_alloc_from_region
8                           /* Ok, try it out.. */
9                           ret = allocate_resource(r, res, size, ...);
10                              err = find_resource(root, new, size,...);
11                                  __find_resource
12
13                                      // 从资源链表中分配地址空间
14                                      // 设置pci_dev->resource[]
15                                      new->start = alloc.start;
16                                      new->end = alloc.end;
17              // 把对应的PCI地址写入BAR
18              pci_update_resource(dev, resno);
19                  pci_std_update_resource
20                      /* 把CPU地址转换为PCI地址: PCI地址 = CPU地址 - offset
21                       * 写入BAR
22                       */
23                      pcibios_resource_to_bus(dev->bus, &region, res);
24                      new = region.start;
25                      reg = PCI_BASE_ADDRESS_0 + 4 * resno;
26                      pci_write_config_dword(dev, reg, new);
```

**garden of idyllic p...**  (focus on)