

# 码农故事

搜索

搜索

[首页](#)[← 上一篇](#)[下一篇 →](#)

## QEMU设备模拟

发表于 2014 年 10 月 15 日

### 设备模拟目的

我们好像不会干一件事而毫无目的，就算不停刷微信朋友圈也是为了打发你无聊的时间。

其实最装B的回答是：设备模拟的目的就是模拟设备。这话是屁话，不过也能说明些什么，确实是模拟设备，用软件的方式提供硬件设备具备的功能。

对于和PC机交互的硬件设备，主要要干两件事，一是提供IRQ中断，二是响应IO输入输出。IO包括PIO/MMIO/DMA等（DMA算不算IO？）

以i8254.c实现的pit为例，主要提供了IRQ注入和PIO响应，见初始化函数pit\_initfn：

```
1 static const MemoryRegionOps pit_ioport_ops = {
2     .read = pit_ioport_read,
3     .write = pit_ioport_write,
4     .impl = {
5         .min_access_size = 1,
6         .max_access_size = 1,
7     },
8     .endianness = DEVICE_LITTLE_ENDIAN,
9 };
10
11 static int pit_initfn(PITCommonState *pit)
12 {
13     PITChannelState *s;
14
15     s = &pit->channels[0];
16     /* the timer 0 is connected to an IRQ */
17     //这里有个irq_timer, 用于qemu_set_irq提供中断注入
18     s->irq_timer = qemu_new_timer_ns(vm_clock, pit_irq_timer, s);
19     qdev_init_gpio_out(&pit->dev.qdev, &s->irq, 1);
20
21     memory_region_init_io(&pit->iports, &pit_ioport_ops, pit, "pit", 4);
22     qdev_init_gpio_in(&pit->dev.qdev, pit_irq_control, 1);
23     return 0;
24 }
```

这里的pit\_ioport\_ops，主要注册GUEST操作系统读写PIO时候的回调函数。

### 模块注册

QEMU要模拟模块那么多，以程序员的喜好，至少得来一套管理这些模拟设备模块的接口，以示设计良好。

QEMU将被模拟的模块分为了四类：

```
1 typedef enum {
2     MODULE_INIT_BLOCK,
3     MODULE_INIT_MACHINE,
4     MODULE_INIT_QAPI,
5     MODULE_INIT_QOM,
6     MODULE_INIT_MAX
7 } module_init_type;
```

- BLOCK

比如磁盘IO就属于BLOCK类型，e.g: block\_init(bdrv\_qcow2\_init); block\_init(iscsi\_block\_init);

- MACHINE

PC虚拟machine\_init(pc\_machine\_init); XEN半虚拟化machine\_init(xenpv\_machine\_init); MIPS虚拟

```
machine_init(mips_machine_init);
```

#### ◦ QAPI

QEMU GUEST AGENT模块里面会执行QAPI注册的回调

#### ◦ QOM

QOM树大枝多，儿孙满堂，应该是这里面最复杂的一个，我们重点介绍。

e.g:

```
1 ObjectClass -> PCIDeviceClass //显卡type_init(cirrus_vga_register_types), 网卡type_init(rtl8139_register_ty
2 IDEDeviceClass //IDE硬盘或CD-ROM type_init(ide_register_types)
3 ISADeviceClass //鼠标键盘type_init(i8042_register_types), RTC时钟type_init(pit_register)
4 SysBusDeviceClass //MMIO IDE(IDE设备直接连接CPU bus而不是连接IDE controller)type_init(mmio_ide_
5 -> CRISCPUClass
```

注册QOM设备的时候，使用QEMU提供的宏，type\_init宏进行注册：

```
1 #define type_init(function) module_init(function, MODULE_INIT_QOM)
2 #define module_init(function, type)
3     static void __attribute__((constructor)) do_qemu_init_## function(void) {
4         register_module_init(function, type);
5     }
```

这和写Linux驱动类似，一般写在一个模块实现文件的最底部，以pit为例，写的是type\_init(pit\_register\_types)展开后为：

```
1 static void __attribute__((constructor)) do_qemu_init_pit_register_types(void)
2 {
3     register_module_init(pit_register_types, MODULE_INIT_QOM);
4 }
```

那么，这个do\_qemu\_init\_pit\_register\_types何时调用？

在gcc里面，给函数加上\_\_attribute\_\_((constructor))，表示此函数需要在main开始前自动调用，测试调用顺序是：全局对象构造函数 -> \_\_attribute\_\_((constructor)) -> main -> 全局对象析构函数 -> \_\_attribute\_\_((destructor))。

调用register\_module\_init就是将pit\_register\_types回调函数插入utilmodule.c里定义的init\_type\_list[MODULE\_INIT\_QOM]链表内。

```
1 void register_module_init(void (*fn)(void), module_init_type type)
2 {
3     ModuleEntry *e;
4     ModuleTypeList *l;
5     e = g_malloc0(sizeof(*e));
6     e->init = fn; //init指针被设置为fn
7     l = find_type(type);
8     QTAILQ_INSERT_TAIL(l, e, node);
9 }
```

通过下面main函数的部分代码可以看出，模块初始化顺序是QOM->MACHINE->BLOCK，至于QAPI，在这个流程里没看到。

```
1 void main()
2 {
3     module_call_init(MODULE_INIT_QOM); //初始化设备
4     qemu_add_opts //初始化默认选项
5     module_call_init(MODULE_INIT_MACHINE); //初始化机器类型
6     machine = find_default_machine(); //这里对machine赋值，下面还会通过参数更改machine
7     vtp_script_execute(g_qemu_start_hook_path, g_fairsched_string, TYPE_START); //开机启动脚本的调用
8     深度分析启动参数
9     bdrv_init_with_whitelist -> bdrv_init -> module_call_init(MODULE_INIT_BLOCK); //初始化BLOCK设备
10    machine->init(&args); //初始化machine
11
12    qemu_run_machine_init_done_notifiers(); //初始化成功回调通知
13    qemu_system_reset(VMRESET_SILENT); //system reset 启动运行
14    if (loadvm) {
15        load_vmstate(loadvm);
16    } else if (loadstate) {
17        load_state_from_blockdev(loadstate);
18    }
19
20    resume_all_vcpus();
21    main_loop(); //进入主循环
22 }
```

在main函数进来的时候，首先调用module\_call\_init(MODULE\_INIT\_QOM);

```
1 void module_call_init(module_init_type type)
2 {
3     ModuleTypeList *l;
4     ModuleEntry *e;
5     l = find_type(type);
6     QTAILQ_FOREACH(e, l, node) {
7         e->init(); //这里，就是调用刚才注册的回调，例如，对于kvm-pit来说，调用的是pit_register
8     }
9 }
```

此module\_call\_init将依次调用注册的回调，如PIT的pit\_register\_types：

```
1 static const TypeInfo pit_info = {
```

```

2  .name         = "isa-pit",           //做为type_table的key
3  .parent       = "pit-common",       //父类型, 这个比较重要, 如果本TypeInfo没有设置class_size, 会根据parent获取par
4  .instance_size = sizeof(PITCommonState), //分配实例的大小
5  .class_init   = pit_class_init,     //初始化函数
6  };
7
8  static void pit_register_types(void)
9  {
10     type_register_static(&pit_info);
11 }

```

pit\_register\_types又进一步调用type\_register\_static -> type\_register -> type\_register\_internal, 这个函数完成的功能其实只是在qomobject.c的type\_table里插入了一个HASH键值对, 以TypeInfo的name为KEY, malloc了一个TypeInfo结构的超集TypeImpl为VALUE, 在以name为KEY回溯parent时需要TypeImpl, 其实这个hash也可以做成一个tree。

## QOM的Object模型

以pit为例, 通过回溯parent你可以看到, 其定义TypeInfo最终形成一个继承关系:

"isa-pit" -> "pit-common" -> "isa-device" -> "device" -> "object"

qomobject.c

```

1  static TypeInfo object_info = {
2      .name = "object",
3      .instance_size = sizeof(Object),
4      .instance_init = object_instance_init,
5      .abstract = true,
6  };

```

hwqdev.c

```

1  static const TypeInfo device_type_info = {
2      .name = "device",
3      .parent = "object",
4      .instance_size = sizeof(DeviceState),
5      .instance_init = device_initfn,
6      .instance_finalize = device_finalize,
7      .class_base_init = device_class_base_init,
8      .class_init = device_class_init,
9      .abstract = true,
10     .class_size = sizeof(DeviceClass),
11 };

```

hwisa-bus.c

```

1  static const TypeInfo isa_device_type_info = {
2      .name = "isa-device",
3      .parent = "device",
4      .instance_size = sizeof(ISADevice),
5      .abstract = true,
6      .class_size = sizeof(ISADeviceClass),
7      .class_init = isa_device_class_init,
8  };

```

hwi8254\_common.c

```

1  static const TypeInfo pit_common_type = {
2      .name = "pit-common",
3      .parent = "isa-device",
4      .instance_size = sizeof(PITCommonState),
5      .class_size = sizeof(PITCommonClass),
6      .class_init = pit_common_class_init,
7      .abstract = true,
8  };

```

hwi8254.c

```

1  static const TypeInfo pit_info = {
2      .name = "isa-pit",
3      .parent = "pit-common",
4      .instance_size = sizeof(PITCommonState),
5      .class_init = pit_class_initfn,
6  };

```

由于TypeInfo只是注册时临时使用, 而TypeImpl是TypeInfo的超集, 所以, 这层关系也反应了TypeImpl的继承关系。

```

1  struct TypeImpl
2  {
3      const char *name;
4      size_t class_size;
5      size_t instance_size;
6      void (*class_init)(ObjectClass *klass, void *data);

```

```

7 void (*class_base_init)(ObjectClass *klass, void *data);
8 void (*class_finalize)(ObjectClass *klass, void *data);
9 void *class_data;
10 void (*instance_init)(Object *obj);
11 void (*instance_finalize)(Object *obj);
12 bool abstract;
13 const char *parent;
14 TypeImpl *parent_type;
15 ObjectClass *klass;
16 int num_interfaces;
17 InterfaceImpl interfaces[MAX_INTERFACES];
18 };

```

TypeImpl
ObjectClass *class
size_t class_size
size_t instance_size
Callback class_init
Callback class_base_init
Callback class_finalize
Callback instance_init
Callback instance_finalize
TypeImpl *parent_type
InterfaceImpl *interfaces

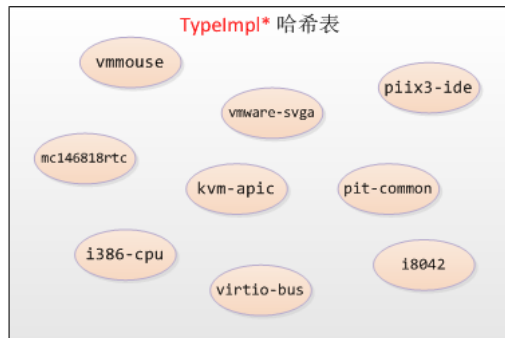


Figure 1 Typelmpl图解

打印查看TypeImpl属性：

```

(gdb) p *obj->class->type //struct TypeImpl * type
$13 = {name = 0x5555566e5e30 "mc146818rtc", class_size = 128, instance_size = 664, class_init = 0x5555555555555555,
class_finalize = 0, class_data = 0x0, instance_init = 0, instance_finalize = 0, abstract = false, parent_type = 0x5555566d8bd0, class = 0x555556a50e50, num_interfaces = 0, interfaces = {{typename = 0x5555566d8bd0, ...}}
其主要包含如下部分：

```

- name/parent/parent\_type 表示自己的，父亲的KEY和TypeImpl指针。
- class/class\_size/class\_init/class\_base\_init/class\_finalize/class\_data 和ObjectClass联系，组成继承关系。
- instance\_init/instance\_finalize和ObjectClass有裙带关系的Object，共同完成继承体系。
- num\_interfaces/interfaces 用于管理接口。

## Object和ObjectClass的关系

还是通过这条继承链来看：

"isa-pit" -> "pit-common" -> "isa-device" -> "device" -> "object"

其中ObjectClass链的定义为：

```

1 struct ObjectClass
2 {
3     /*< private >*/
4     Type type;
5     GSList *interfaces;
6     ObjectUnparent *unparent;
7 };
8 typedef struct DeviceClass {
9     /*< private >*/
10    ObjectClass parent_class;
11    /*< public >*/
12    const char *fw_name;
13    const char *desc;
14    Property *props;
15    int no_user;
16    /* callbacks */
17    void (*reset)(DeviceState *dev);
18    DeviceRealize realize;
19    DeviceUnrealize unrealize;
20    /* device state */
21    const struct VMStateDescription *vmsd;
22    /* Private to qdev / bus. */
23    qdev_initfn init; /* TODO remove, once users are converted to realize */
24    qdev_event unplug;
25    qdev_event exit;
26    const char *bus_type;
27 } DeviceClass;
28 typedef struct ISADeviceClass {
29     DeviceClass parent_class;
30     int (*init)(ISADevice *dev);
31 } ISADeviceClass;
32 typedef struct PITCommonClass {
33     ISADeviceClass parent_class;
34     int (*init)(PITCommonState *s);
35     void (*set_channel_gate)(PITCommonState *s, PITChannelState *sc, int val);
36     void (*get_channel_info)(PITCommonState *s, PITChannelState *sc,
37                             PITChannelInfo *info);
38     void (*pre_save)(PITCommonState *s);
39     void (*post_load)(PITCommonState *s);
40 } PITCommonClass;

```

下层定义包含上层，很明显的继承模型，ObjectClass更像C++的CLASS，而Object链的定义为：

```

1 struct Object
2 {
3     /*< private >*/
4     ObjectClass *class;
5     ObjectFree *free;
6     QTAILQ_HEAD(, ObjectProperty) properties;
7     uint32_t ref;
8     Object *parent;
9 };
10 struct DeviceState {
11     /*< private >*/
12     Object parent_obj;
13     /*< public >*/
14     const char *id;
15     bool realized;
16     QemuOpts *opts;
17     int hotplugged;
18     BusState *parent_bus;
19     int num_gpio_out;
20     qemu_irq *gpio_out;
21     int num_gpio_in;
22     qemu_irq *gpio_in;
23     QLIST_HEAD(, BusState) child_bus;
24     int num_child_bus;
25     int instance_id_alias;
26     int alias_required_for_version;
27 };
28 struct ISADevice {
29     DeviceState qdev;
30     uint32_t isairq[2];
31     int nirqs;
32     int ioport_id;
33 };
34 typedef struct PITCommonState {
35     ISADevice dev;
36     MemoryRegion ioports;
37     uint32_t iobase;
38     PITChannelState channels[3];
39 } PITCommonState;

```

有了ObjectClass为什么还要有个Object？从代码看，ObjectClass只有一份实例，而Object是可以多个实例的，Object引用ObjectClass获得ObjectClass的特征，但是同时又节约了初始化和存放ObjectClass的CPU和空间，相同的ObjectClass可以被多个Object引用，例如scsi-disk.c里面有“scsi-hd”, “scsi-cd”, “scsi-block”, “scsi-disk”四种Object共同引用了“scsi-device”。这里可以想象成C++的虚继承，ObjectClass是virtual class而Object是class。其实两者是可以柔和在一起的，Object也有对应的继承关系，用来保存特定属性。

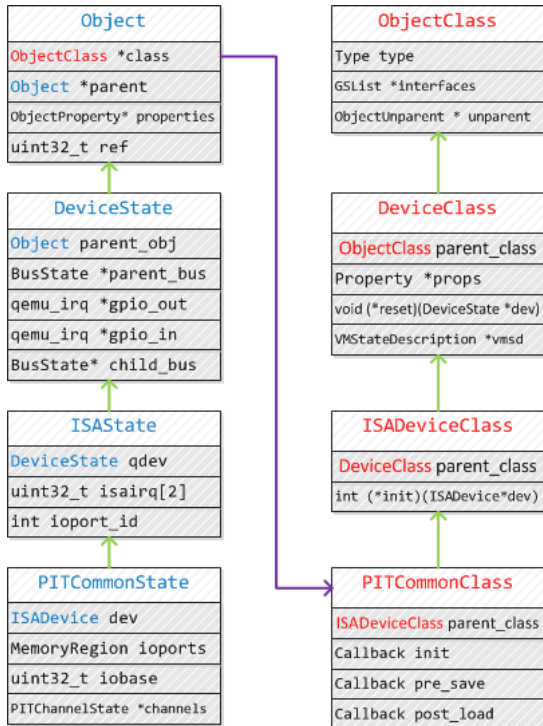


Figure 2 ObjectClass和Object 关系

## Object和ObjectClass的初始化

上面讲的Object和ObjectClass主要是完成一个对象继承模型，从代码看QEMU的这个模型实现并不非常很优雅，封装不够彻底，就像你妈给你做了条裤子，却没有做裤腰带，还得提着上路。

Object和ObjectClass的初始化方式并不一致，需要分别初始化，ObjectClass通常使用`object_class_by_name`获取，此函数会根据提供的KEY去查找TypeImpl并初始化ObjectClass指针；而Object的初始化是使用的`object_new`，通过参数KEY查找TypeImpl然后

malloc 实例。以qdev\_try\_create获取"isa-pit"的Object DeviceState实例来说，其获取DeviceState的函数如此定义：

```

1 DeviceState *qdev_try_create(BusState *bus, const char *type)
2 {
3     DeviceState *dev;
4     //这个type为TypeInfo.name, 例如"isa-pit"
5     if (object_class_by_name(type) == NULL) {
6         return NULL;
7     }
8     //type_initialize完成后, object_new用来实例化一个instance
9     dev = DEVICE(object_new(type)); // = DEVICE(object_new_with_type(type_get_by_name(typename)))
10    if (!dev) {
11        return NULL;
12    }
13    if (!bus) {
14        bus = sysbus_get_default();
15    }
16    qdev_set_parent_bus(dev, bus);
17    object_unref(OBJECT(dev));
18    return dev;
19 }
20
21 ObjectClass *object_class_by_name(const char *typename)
22 {
23     //之前在type_register_static的时候, 注册了TypeInfo.name, 例如"isa-pit"为key的TypeImpl
24     TypeImpl *type = type_get_by_name(typename);
25     if (!type) {
26         return NULL;
27     }
28     type_initialize(type); //这里面, 初始化class,
29     return type->class;
30 }
31
32 //其实这个函数更应该叫做new_TypeInfo_class()
33 static void type_initialize(TypeImpl *ti)
34 {
35     TypeImpl *parent;
36     if (ti->class) {
37         return;
38     }
39     /*
40     type_class_get_size 首先获取自己的class_size变量, 如果没有, 再找parent类型所指的TypeImpl的class_size, 直到找到为止
41     比如"isa-pit"没有设置class_size, 那么获取的是"pit-common"的class_size, 而type_object_get_size也是类似
42     static const TypeInfo pit_common_type = {
43         .name = "pit-common",
44         .parent = "isa-device",
45         .instance_size = sizeof(PITCommonState),
46         .class_size = sizeof(PITCommonClass),
47         .class_init = pit_common_class_init,
48         .abstract = true,
49     };
50     */
51     ti->class_size = type_class_get_size(ti);
52     ti->instance_size = type_object_get_size(ti);
53     ti->class = g_malloc0(ti->class_size);
54     parent = type_get_parent(ti);
55     if (parent) {
56         //1, 保证parent初始化了
57         type_initialize(parent);
58         GSList *e;
59         int i;
60
61         //2, 将parent的class内容memcpy一份给自己的对应的parent区域
62         g_assert(parent->class_size <= ti->class_size);
63         memcpy(ti->class, parent->class, parent->class_size);
64
65         //3, 将parent里面的class的interfaces做一次深度复制, 复制给自己
66         for (e = parent->class->interfaces; e; e = e->next) {
67             ObjectClass *iface = e->data;
68             type_initialize_interface(ti, object_class_get_name(iface));
69         }
70
71         //4. 如果本类型有自己的interfaces, 初始化
72         for (i = 0; i < ti->num_interfaces; i++) {
73             TypeImpl *t = type_get_by_name(ti->interfaces[i].typename);
74             for (e = ti->class->interfaces; e; e = e->next) {
75                 TypeImpl *target_type = OBJECT_CLASS(e->data)->type;
76                 if (type_is_ancestor(target_type, t)) {
77                     break;
78                 }
79             }
80             if (e) {
81                 continue;
82             }
83             type_initialize_interface(ti, ti->interfaces[i].typename);
84         }
85     }
86
87     ti->class->type = ti;
88     while (parent) {
89         if (parent->class_base_init) {
90             //回溯回调parent的class_base_init函数
91             parent->class_base_init(ti->class, ti->class_data);
92         }
93         parent = type_get_parent(parent);
94     }
95
96     if (ti->class_init) {
97         /*
98         如果本类设置了class_init, 回调它, ti->class_data是一个void*的参数
99         比如"isa-pit"我们设置了pit_class_initfn
100        这个函数主要干嘛? 主要填充class里的其他该填充的地方。

```

```

101     malloc之后你总得调用构造函数吧, 调用构造函数的第一句都是super(xxx)
102     这工作前面, 2, 3步骤已经做了, 然后干你自己的活。见pit_class_initfn定义
103     */
104     ti->class_init(ti->class, ti->class_data);
105 }
106 }

```

上述代码把object\_class\_by\_name的流程说完了, 再看看object\_new(type) = object\_new\_with\_type(type\_get\_by\_name(typename))的流程:

```

1  Object *object_new_with_type(Type type)
2  {
3      Object *obj;
4      g_assert(type != NULL);
5      type_initialize(type);
6      obj = g_malloc(type->instance_size); //这个instance_size是初始化TypeInfo的时候设置的sizeof(PITCommonState)
7      object_initialize_with_type(obj, type);
8      obj->free = g_free;
9      return obj;
10 }
11 void object_initialize_with_type(void *data, TypeImpl *type)
12 {
13     Object *obj = data;
14     g_assert(type != NULL);
15     type_initialize(type);
16     g_assert(type->instance_size >= sizeof(Object));
17     g_assert(type->abstract == false);
18     memset(obj, 0, type->instance_size);
19     obj->class = type->class; //instance的类型通过class指针指定
20     object_ref(obj);
21     QTAILQ_INIT(&obj->properties);
22     object_init_with_type(obj, type); //深度递归调用TypeImpl及其parent的instance_init函数指针, 相当于new instance的构造
23 }

```

qdev\_try\_create->object\_class\_by\_name->type\_initialize的调用流程, 如果父ObjectClass没初始化, 会初始化父ObjectClass, 此时调用到父ObjectClass对应name的TypeImpl的class\_init函数, 例如"pit-common"的class\_init回调pit\_common\_class\_init, 此回调会设置ISADeviceClass的init回调为pit\_init\_common。

## Object类型转换

再解释下Object里常用的一个宏: OBJECT\_CHECK, 以ISA\_DEVICE这段代码为例:

```

1  static const TypeInfo mc146818rtc_info = {
2      .name = "mc146818rtc",
3      .parent = "isa-device",
4      .instance_size = sizeof(RTCState),
5      .class_init = rtc_class_initfn,
6  };
7  ISADevice *isa_create(ISABus *bus, const char *name)
8  {
9      DeviceState *dev;
10     if (!bus) {
11         hw_error("Tried to create isa device %s with no isa bus present.", name);
12     }
13     dev = qdev_create(&bus->qbus, name);
14     return ISA_DEVICE(dev); //毫无疑问, "mc146818rtc"肯定可以转换为"isa-device", why? 见mc146818rtc_info定义
15 }
16 ISADevice *dev = isa_create(bus, "mc146818rtc");

```

这里的ISA\_DEVICE体现了OBJECT的类型转换功能, 宏定义为:

```

1  #define OBJECT(obj)
2      ((Object *) (obj))
3  #define OBJECT_CHECK(type, obj, name)
4      ((type *) object_dynamic_cast_assert(OBJECT(obj), (name)))
5  #define ISA_DEVICE(obj)
6      OBJECT_CHECK(ISADevice, (obj), TYPE_ISA_DEVICE)

```

展开后为:

```

1  #define ISA_DEVICE(dev) (ISADevice *) object_class_dynamic_cast(((Object *) dev)->class, "isa-device")

```

object.c里面定义了object\_class\_dynamic\_cast函数, 其实此函数功能比较简单, 就是通过遍历parent看当前class是否有一个祖先是typename, 其定义如下:

```

1  ObjectClass *object_class_dynamic_cast(ObjectClass *class,
2      const char *typename)
3  {
4      TypeImpl *target_type = type_get_by_name(typename); //找到"isa-device"对应的TypeImpl*
5      TypeImpl *type = class->type; //本ObjectClass真实的TypeImpl*, 其实这里是"mc146818rtc"
6      ObjectClass *ret = NULL;
7      /*
8       (gdb) p *target_type
9      $31 = {name = 0x5555556d0b20 "isa-device", class_size = 128, instance_size = 160, class_init = 0x55555568ddd0 <isa
10         class_finalize = 0, class_data = 0x0, instance_init = 0, instance_finalize = 0, abstract = true, parent = 0x5555

```



```

11  class = 0x555556a14750, num_interfaces = 0, interfaces = {{typename = 0x0} <repeats 32 times>}}
12  (gdb) p *type_interface
13  $33 = {name = 0x5555566e40b0 "interface", class_size = 32, instance_size = 0, class_init = 0, class_base_init = 0,
14  instance_finalize = 0, instance_finalize = 0, abstract = true, parent = 0x0, parent_type = 0x0, class = 0x0, num_int
15  typename = 0x0} <repeats 32 times>}}
16  type_is_ancestor用于判断, "interface"是否是"isa-device"的祖先 (判断方法是递归遍历"isa-device"的parent, 比较是否有"i
17  */
18  if (type->num_interfaces && type_is_ancestor(target_type, type_interface)) {
19      int found = 0;
20      GSList *i;
21      for (i = class->interfaces; i; i = i->next) {
22          ObjectClass *target_class = i->data;
23          if (type_is_ancestor(target_class->type, target_type)) {
24              ret = target_class;
25              found++;
26          }
27      }
28      /* The match was ambiguous, don't allow a cast */
29      if (found > 1) {
30          ret = NULL;
31      }
32  } else if (type_is_ancestor(type, target_type)) {
33      /*
34       判断type="mc146818rtc"的祖先是否是target_type="isa-device",
35       如果是, 这里表示子类class="mc146818rtc"能成功转换为父类typename="isa-device"所指的ObjectClass
36       */
37      ret = class;
38  }
39  return ret;
40 }

```

GDB显示其内部数据为：

```

(gdb) p *obj //Object *obj
$11 = {class = 0x555556a50e50, free = 0x7ffff7424020, properties = {tqh_first = 0x555556a17da0, tqh_las
(gdb) p *obj->class //ObjectClass * class
$12 = {type = 0x5555566e5cb0, interfaces = 0x0, unparent = 0x555556b1ac0 }
(gdb) p *obj->class->type //struct TypeImpl * type
$13 = {name = 0x5555566e5e30 "mc146818rtc", class_size = 128, instance_size = 664, class_init = 0x555555
class_finalize = 0, class_data = 0x0, instance_init = 0, instance_finalize = 0, abstract = false, pare
parent_type = 0x5555566d8bd0, class = 0x555556a50e50, num_interfaces = 0, interfaces = {{typename = 0x

```

## QOM设备初始化

基于Object和ObjectClass实现的QOM设备，何时触发他的初始化，以PIT为例，将之前的Object和ObjectClass想象成C++，那么PIT对应的PitCommonState定义应该类似如下所示：

```

1  class PITCommonClass : public ISADeviceClass {
2  public:
3      virtual int init(PITCommonState *s) = 0;
4  };
5
6  class ISADevice : public DeviceState {
7  public:
8      int nirqs;
9      int ioport_id;
10 };
11
12 class PITCommonState : public ISADevice, public PITCommonClass {
13     int init(PITCommonState *s);
14 };

```

看吧，QEMU绕了这么大一个圈子，就想实现这样一个结构，所以有的时候用C++还是有好处的（虽然本人生理周期现正处于不太喜欢C++时间）。

那么，何处调用了new PITCommonState()？

这得从main函数开始看，main函数里面，有machine->init(&args);函数调用，这是对注册的machine的初始化，而默认的machine是在pc\_piix.c里面pc\_machine\_init函数注册的第一个machine，即：

```

1  static QEMUMachine pc_i440fx_machine_v1_4 = {
2      .name = "pc-i440fx-1.4",
3      .alias = "pc",
4      .desc = "Standard PC (i440FX + PIIX, 1996)",
5      .init = pc_init_pci,
6      .max_cpus = 255,
7      .is_default = 1,
8      .default_machine_opts = KVM_MACHINE_OPTIONS,
9      DEFAULT_MACHINE_OPTIONS,
10 };
11 qemu_register_machine(&pc_i440fx_machine_v1_4);

```

当main函数调用machine->init时，我的实验环境默认情况其实就是调用的pc\_i440fx\_machine\_v1\_4的初始化回调pc\_init\_pci -> pc\_init1，这个函数主要初始化相关PC硬件：

```

1  static void pc_init1(MemoryRegion *system_memory,
2                      MemoryRegion *system_io,
3                      ram_addr_t ram_size,
4                      const char *boot_device,
5                      const char *kernel_filename,

```



```

6      const char *kernel_cmdline,
7      const char *initrd_filename,
8      const char *cpu_model,
9      int pci_enabled,
10     int kvmclock_enabled)
11 {
12     //CPU类型初始化-> cpu_x86_init -> mce_init/qemu_init_vcpu, 初始化VCPU
13     pc_cpus_init(cpu_model);
14     //初始化acpi_tables
15     pc_acpi_init("acpi-dsdt.aml");
16     if (!xen_enabled()) {
17         //ROM, BIOS, RAM相关初始化
18         fw_cfg = pc_memory_init(system_memory,
19                                 kernel_filename, kernel_cmdline, initrd_filename,
20                                 below_4g_mem_size, above_4g_mem_size,
21                                 rom_memory, &ram_memory);
22     }
23     //IRQ, 初始化
24     //VGA初始化
25     pc_vga_init(isa_bus, pci_enabled ? pci_bus : NULL);
26     /* init basic PC hardware */
27     pc_basic_device_init(isa_bus, gsi, &rtc_state, &floppy, xen_enabled()); //这里调用pit_init
28     //初始化网卡
29     pc_nic_init(isa_bus, pci_bus);
30     //初始化硬盘, 音频设备
31     //初始化cmos数据, 比如设置cmos rtc时钟, 是否提供PS/2设备
32     pc_cmos_init(below_4g_mem_size, above_4g_mem_size, boot_device,
33                 floppy, idebus[0], idebus[1], rtc_state);
34     //初始化USB
35     if (pci_enabled && usb_enabled(false)) {
36         pci_create_simple(pci_bus, piix3_devfn + 2, "piix3-usb-uhci");
37     }
38 }
39
40 void pc_basic_device_init(ISABus *isa_bus, qemu_irq *gsi,
41                           ISADevice **rtc_state,
42                           ISADevice **floppy,
43                           bool no_vmport)
44 {
45     //初始化HPET
46     //初始化mc146818 rtc
47     //初始化i8042 PIT
48     pit = pit_init(isa_bus, 0x40, pit_isa_irq, pit_alt_irq);
49     //初始化串口, 并口
50     //初始化vmmouse ps2_mouse
51 }

```

接下来的流程是pit\_init -> isa\_create(bus, "isa-pit") -> qdev\_create -> qdev\_try\_create, qdev\_try\_create的实现在前面已经讲了, 如上节所述, 它分别使用object\_class\_by\_name和object\_new来初始化ObjectClass和Object。

## 属性设置

Object同ObjectClass的显著区别就是Object提供了属性的概念, 以MC146818为例, 其定义时设置了"base\_year"和"lost\_tick\_policy" :

```

1 static Property mc146818rtc_properties[] = {
2     DEFINE_PROP_INT32("base_year", RTCState, base_year, 1980),
3     DEFINE_PROP_LOSTTICKPOLICY("lost_tick_policy", RTCState,
4     lost_tick_policy, LOST_TICK_DISCARD),
5     DEFINE_PROP_END_OF_LIST(),
6 };

```

但是用GDB一看 :

```

(gdb) p *obj->properties.tqh_first
$15 = {name = 0x555556a17df0 "type", type = 0x555556a17e10 "string", get = 0x555555727ae0
, release = 0x555555727980
, opaque = 0x555556a17d80, node = {tqe_next = 0x555556a17e50, tqe_prev = 0x555556a17af0}}
(gdb) p *obj->properties.tqh_first.node.tqe_next
$21 = {name = 0x555556a17ea0 "realized", type = 0x555556a17ec0 "bool", get = 0x5555557279c0
, release = 0x555555727940
, opaque = 0x555556a17e30, node = {tqe_next = 0x555556a17ee0, tqe_prev = 0x555556a17dd0}}
(gdb) p *obj->properties.tqh_first.node.tqe_next.node.tqe_next
$22 = {name = 0x555556a17f30 "base_year", type = 0x555556a1be10 "int32", get = 0x5555556af370 , set = 0x5555555d4e1a0, node = {tqe_next = 0x555556a50ee0, tqe_prev = 0x555556a17e80}}
(gdb) p *obj->properties.tqh_first.node.tqe_next.node.tqe_next.node.tqe_next
$23 = {name = 0x555556a1be30 "lost_tick_policy", type = 0x555556a1be50 "LostTickPolicy", get = 0x5555556af370 , set = 0x5555555d4e1c0, node = {tqe_next = 0x555556a50fa0, tqe_prev = 0x555556a17f10}}
(gdb) p *obj->properties.tqh_first.node.tqe_next.node.tqe_next.node.tqe_next.node.tqe_next
$24 = {name = 0x555556a50f80 "parent_bus", type = 0x55555680b1d0 "link", get = 0x555555729e50 , set = 0x55555572a320 , release = 0, opaque = 0x555556a17b30, node = {tqe_next = 0x0, tqe_prev = 0x555555729e50}}
(gdb) p *obj->properties.tqh_first.node.tqe_next.node.tqe_next.node.tqe_next.node.tqe_next
Cannot access memory at address 0x0
事实上却多了"type" "realized" "parent_bus", 这些属性都是动态添加的。

```

在"object"类型的, instance\_init = object\_instance\_init回调处, 添加了

```

1 object_property_add_str(obj, "type", qdev_get_type, NULL, NULL);

```

在"device"类型的instance\_init = device\_initfn回调处，添加了

```
1 object_property_add_bool(obj, "realized",
2                          device_get_realized, device_set_realized, NULL);
3 object_property_add_link(OBJECT(dev), "parent_bus", TYPE_BUS,
4                          (Object **)&dev->parent_bus, NULL);
```

设置属性的时候，调用类似qdev\_prop\_set\_int32(&dev->qdev, "base\_year", base\_year);进行设置，这里，注意第一个参数为什么是DeviceState\* dev->qdev而不是ISADevice \*dev？

因为rtc\_class\_initfn里初始化props是给 DeviceClass \*dc初始化的，所以对应的应该是DeviceState而不是子类ISADevice。

设置属性是如何实现的？

以"realized"的bool属性设置为例，调用顺序为object\_property\_set\_bool -> object\_property\_set\_qobject -> object\_property\_set，此函数定义：

```
1 void object_property_set(Object *obj, Visitor *v, const char *name,
2                          Error **errp)
3 {
4     //obj还是"mc146818rtc"的实例，name为"realized"，object_property_find其实就是查找obj的properties链表里是否存在名字为
5     ObjectProperty *prop = object_property_find(obj, name, errp);
6     if (prop == NULL) {
7         return;
8     }
9     if (!prop->set) { //如果存在，且没有设置过set handler，错误
10        error_set(errp, QERR_PERMISSION_DENIED);
11    } else { // "realized"的set函数为property_set_bool
12        prop->set(obj, v, prop->opaque, name, errp);
13    }
14 }
15
16 static void property_set_bool(Object *obj, Visitor *v, void *opaque,
17                               const char *name, Error **errp)
18 {
19     BoolProperty *prop = opaque;
20     bool value;
21     Error *local_err = NULL;
22     visit_type_bool(v, &value, name, &local_err);
23     if (local_err) {
24         error_propagate(errp, local_err);
25         return;
26     }
27     prop->set(obj, value, errp); //对于realized来说，其实就是调用device_set_realized
28 }
```

而CALLBACK=device\_set\_realized 又会调用CALLBACK=device\_realize。

## 模块PIO回调流程

### 注册回调

以mc146818 rtc为例，在rtc\_initfn的时候，注册回调的代码如下：

```
1 static const MemoryRegionOps cmos_ops = {
2     .read = cmos_ioport_read,
3     .write = cmos_ioport_write,
4     .impl = {
5         .min_access_size = 1,
6         .max_access_size = 1,
7     },
8     .endianness = DEVICE_LITTLE_ENDIAN,
9 };
10 void isa_register_ioport(ISADevice *dev, MemoryRegion *io, uint16_t start)
11 {
12     memory_region_add_subregion(isabus->address_space_io, start, io);
13     isa_init_ioport(dev, start);
14 }
15
16 memory_region_init_io(&s->io, &cmos_ops, s, "rtc", 2);
17 isa_register_ioport(dev, &s->io, base);
```

其中s->io是MemoryRegion类型，MemoryRegion是可以像树一样，多级挂载，比如，现在将rtc的MemoryRegion挂载在isabus的address\_space\_io这个MemoryRegion下，其start参数为offset在整个isabus->address\_space\_io MemoryRegion中的偏移，即0x70，那么END呢？END在memory\_region\_init\_io的时候已经存储到MemoryRegion的size里面了。

再看看isabus内容，有个更深入的性感的认识：

```
(gdb) p *isabus
$8 = {qbus = {obj = {class = 0x5555556a14e70, free = 0x7ffff7424020, properties = {tqh_first = 0x5555556a14e70, tqh_last = 0x5555556a16450, sibling = {le_next = 0x0, le_prev = 0x555555698d7b8}}, address_space_io = {ops = 0x0, opaque = 0x0, parent = 0x0, size = {lo = 65536, hi = 0}, addr = 0, destructor = 0x5555556a14e70, terminated = false, readable = true, w...
(gdb) p *isabus->address_space_io
$9 = {ops = 0x0, opaque = 0x0, parent = 0x0, size = {lo = 65536, hi = 0}, addr = 0, destructor = 0x5555556a14e70, terminated = false, readable = true, w...}
```

```
suppage = false, terminates = false, readable = true, ram = false, readonly = false, enabled = true, flush_coalesced_mmio = false, alias = 0x0, alias_offset = 0, priority = 0, may_overlap = false, subregions_link = {tqh_last = 0x55555698bf60}, subregions_link = {tqh_next = 0x0, tqe_prev = 0x0}, coalesced = {tqh_first = 0x5555566f7220 "io", dirty_log_mask = 0 '00', ioeventfd_nb = 0, ioeventfds = 0x0, updateaddr =
```

## 回调流程

QEMU通过kvm\_cpu\_exec -> kvm\_vcpu\_ioctl(cpu, KVM\_RUN, 0) 执行GUEST机CODE，当GUEST遇到IO等操作需要退出，会先在KVM里处理，KVM不能处理，kvm\_vcpu\_ioctl就返回，给QEMU处理，QEMU根据返回的run->exit\_reason进行分派，比如PROT\_READ 0x71操作，退出时其exit\_reason为KVM\_EXIT\_IO，kvm\_handle\_io里，根据direction判断是read还是write，根据read的长度，判断该回调哪个函数。比如0x71 read 1字节的时候，调用的是：

```
stb_p(ptr, cpu_inb(port));
```

stb\_p是将第二个参数cpu\_inb(port)的结果转换为一个字节大小赋值给第一个参数ptr所指内存。

```
cpu_inb (addr=addr@entry=113)
```

cpu\_inb和cpu\_inw和cpu\_inl是一家人的三兄弟，长得极其神似，我们看cpu\_inb，在他调用的ioport\_read的时候，第一个参数叫index=0，这是和他的兄弟cpu\_inw, cpu\_inl区别开来的关键特征。本来这里应该没有index啥事的，但总有人偷懒不设置对应addr的handler，index就是对专门为这种懒人擦屁股，没有handler的时候，给选择一个默认handler，你大概也看明白了，就index的话，三兄弟的差别在于inb=0, inw=1, inl=2。

```
1 uint8_t cpu_inb(pio_addr_t addr)
2 {
3     uint8_t val;
4     val = ioport_read(0, addr);
5     return val;
6 }
```

read的地址比较重要，例如rtc的0x71，index其实是用来选择默认handler的，当在对应的ioport\_read\_table里面没有注册函数的时候就根据index的值，分别选择readb, readw, readl来做默认操作。

```
1 static uint32_t ioport_read(int index, uint32_t address)
2 {
3     static IOPortReadFunc * const default_func[3] = {
4         default_ioport_readb,
5         default_ioport_readw,
6         default_ioport_readl
7     };
8     IOPortReadFunc *func = ioport_read_table[index][address];
9     if (!func)
10         func = default_func[index];
11     /*
12      * func一般都是ioport_readb_thunk 关键在于ioport_opaque[address]这里面存放的是不同端口的IORange*
13      * 这个ioport_opaque的每个值，都存储的该端口对应的IORange*
14      * 当读写此端口的时候，就会找到之前注册的IORange回调，比如mc146818的IORange为
15      * (gdb) p *(IORange *)ioport_opaque[0x71]
16      * $22 = {ops = 0x7fca4a51c380, base = 112, len = 2}
17      * (gdb) p *(IORangeOps *)0x7fca4a51c380
18      * $25 = {read = 0x7fca49fe1190 <memory_region_iorange_read>, write = 0x7fca49fe1050 <memory_region_iorange_write>,
19      * destructor = 0x7fca49fdcb0 <memory_region_iorange_destructor>}
20      *
21      * 由于x70, 0x71都是readb，所以在mc146818设备的时候，这个func其实是ioport_readb_thunk
22      * (gdb) p ioport_read_table[0][0x70]
23      * $67 = (IOPortReadFunc *) 0x5555557c6980 <ioport_readb_thunk>
24      * (gdb) p ioport_read_table[0][0x71]
25      * $68 = (IOPortReadFunc *) 0x5555557c6980 <ioport_readb_thunk>
26      */
27     return func(ioport_opaque[address], address);
28 }
```

```
ioport_readb_thunk (opaque=, addr=)
```

ioport\_register里面，注册了对一个字节的ioport read handler为ioport\_readb\_thunk，其实这个函数非常简单

就是调用了ops->read的时候，将width设置为1，和ioport\_readw\_thunk，ioport\_readl\_thunk之类的就差一个width的区别

为什么要搞这么复杂？这是为了64K的read空间设计的回调，因为不同的offset位置，我们需要知道是应该调用readb还是readw还是readl。

```
1 static IOPortReadFunc *ioport_read_table[3][64 * 1024]
2 static uint32_t ioport_readb_thunk(void *opaque, uint32_t addr)
3 {
4     IORange *iport = opaque;
5     uint64_t data;
6     //read is memory_region_iorange_read when input char to ps/2 keyboard
7     //比如，mc146818的时候，addr为x71，iport->base为x71，iport->len=2
8     iport->ops->read(iport, addr - iport->base, 1, &data);
9     return data;
10 }
```

上面的回调为memory\_region\_iorange\_read (iorange=0x55555680b1a0, offset=1, width=1, data=0x7ffec1fdc00)

```
1 static const MemoryRegionOps cmos_ops = {
2     .read = cmos_ioport_read,
3     .write = cmos_ioport_write,
4     .impl = {
```

```

5     .min_access_size = 1,
6     .max_access_size = 1,
7 },
8     .endianness = DEVICE_LITTLE_ENDIAN,
9 };
10
11 static void memory_region_iorange_read(IORange *iorange,
12                                         uint64_t offset,
13                                         unsigned width,
14                                         uint64_t *data)
15 {
16     MemoryRegionIORange *mr
17     = container_of(iorange, MemoryRegionIORange, iorange);
18     /*
19      * 一段MemoryRegionIORange里包含了IORange iorange和MemoryRegion* mr
20      (gdb) p *mr
21 $58 = {iorange = {ops = 0x555555d16200, base = 0x70, len = 2}, mr = 0x555556a17b80, offset = 0}
22 (gdb) p *mr->iorange.ops
23 $59 = {read = 0x5555557ccf50 <memory_region_iorange_read>, write = 0x5555557cce10 <memory_region_iorange_write>,
24        destructor = 0x5555557cba70 <memory_region_iorange_destructor>}
25 (gdb) p *mr->mr
26 $60 = {ops = 0x555555d14ba0, opaque = 0x555556a17ae0, parent = 0x555556708930, size = {lo = 2, hi = 0}, addr = 112
27        destructor = 0x5555557cb910 <memory_region_destructor_iomem>, ram_addr = 18446744073709551615, subpage = false,
28        readonly = false, enabled = true, rom_device = false, warning_printed = false, flush_coalesced_mmio = false, ali
29        may_overlap = false, subregions = {tqh_first = 0x0, tqh_last = 0x555556a17be8}, subregions_link = {tqe_next = 0x
30        coalesced = {tqh_first = 0x0, tqh_last = 0x555556a17c08}, name = 0x55555680b300 "rtc", dirty_log_mask = 0 '00',
31        updateaddr = 0, updateopaque = 0x0}
32 (gdb) p *mr->mr->ops
33 $61 = {read = 0x55555579ea20 <cmos_ioport_read>, write = 0x55555579e040 <cmos_ioport_write>, endianness = DEVICE_L
34        max_access_size = 0, unaligned = false, accepts = 0}, impl = {min_access_size = 1, max_access_size = 1, unalign
35        read = {0, 0, 0}, write = {0, 0, 0}}
36     */
37     MemoryRegion *mr = mr->mr;
38
39     //如果mr还有offset, 要加上这个偏移, 这个offset其实是当成地址来用的, 比如, 我认为read 0x71应该在已有offset=1的基础上
40     offset += mr->offset;
41     if (mr->ops->old_portio) { //对"mc146818rtc"已经没有old_portio的CALLBACK了, 跳过
42         const MemoryRegionPortio *mrp = find_portio(mr, offset - mr->offset,
43                                                     width, false);
44
45         *data = ((uint64_t)1 << (width * 8)) - 1;
46         if (mrp) {
47             *data = mrp->read(mr->opaque, offset);
48         } else if (width == 2) {
49             mrp = find_portio(mr, offset - mr->offset, 1, false);
50             assert(mrp);
51             *data = mrp->read(mr->opaque, offset) |
52                   (mrp->read(mr->opaque, offset + 1) << 8);
53         }
54     }
55     return;
56     *data = 0; //这是read后的返回值存储区域, 提前清零
57     access_with_adjusted_size(offset, data, width,
58                               mr->ops->impl.min_access_size, //这个min_access_size和max_access_size是在设置ops的时候
59                               mr->ops->impl.max_access_size,
60                               memory_region_read_accessor, mr);
61 }

```

从参数看, 这里的access参数为memory\_region\_read\_accessor, 而这个value参数, 用来存放read的返回值。接下来进入access\_with\_adjusted\_size (addr=addr@entry=1, value=value@entry=0x7ffec1fdc00, size=1, access\_size\_min=, access\_size\_max=, access=access@entry=0x5555557cbf70, opaque=opaque@entry=0x555556a17b80), 其access参数非常重要, 继续回调的就是access。

```

1 static void access_with_adjusted_size(hwaddr addr,
2                                       uint64_t *value,
3                                       unsigned size,
4                                       unsigned access_size_min,
5                                       unsigned access_size_max,
6                                       void (*access)(void *opaque,
7                                                       hwaddr addr,
8                                                       uint64_t *value,
9                                                       unsigned size,
10                                                      unsigned shift,
11                                                      uint64_t mask),
12                                       void *opaque)
13 {
14     uint64_t access_mask;
15     unsigned access_size;
16     unsigned i;
17
18     if (!access_size_min) {
19         access_size_min = 1;
20     }
21     if (!access_size_max) {
22         access_size_max = 4;
23     }
24     access_size = MAX(MIN(size, access_size_max), access_size_min); //size其实在参数里面已经指定了, 但是为了安全, 要确保
25     access_mask = -1ULL >> (64 - access_size * 8); //作为mask, 对Read的几个mask, 确保结果大小为预期大小
26     for (i = 0; i < size; i += access_size) {
27         /*
28          * 最大返回结果其实只有sizeof(value) = 64bit, 这里的设计是, 每次取一个字节的返回结果
29          * 但是access_size可以不为bit, 比如read 0x100, 假设read范围为bit, 就是x100-0x104, access_size可以为,
30          * 这样就分两步走, 第一步Read 0x100-0x102返回个字节的结果, 存储到value的低字节
31          * 第二步Read 0x103-0x104返回个字节的结果, 存储到value的高字节
32          * 最后返回的value就存储了两次Read值, 只占用了bit, 不会超过bit
33          */
34         access(opaque, addr + i, value, access_size, i * 8, access_mask);
35     }
36 }

```

上面的access\_with\_adjusted\_size的access参数其实就是memory\_region\_read\_accessor，可以通过GDB打印出来：

```
memory_region_read_accessor (opaque=0x555556a17b80, addr=, value=0x7fffec1fdc00, size=1, shift=0, mask=2
(gdb) p *mr
$52 = {ops = 0x555555d14ba0, opaque = 0x555556a17ae0, parent = 0x555556708930, size = {lo = 2, hi = 0},
  destructor = 0x5555557cb910, ram_addr = 18446744073709551615, subpage = false, terminates = true, read
  readonly = false, enabled = true, rom_device = false, warning_printed = false, flush_coalesced_mmio =
  may_overlap = false, subregions = {tqh_first = 0x0, tqh_last = 0x555556a17be8}, subregions_link = {tqh
  coalesced = {tqh_first = 0x0, tqh_last = 0x555556a17c08}, name = 0x55555680b300 "rtc", dirty_log_mask
  updateaddr = 0, updateopaque = 0x0}
(gdb) p *mr->ops
$53 = {read = 0x55555579ea20, write = 0x55555579e040, endianness = DEVICE_LITTLE_ENDIAN, valid = {min_
  max_access_size = 0, unaligned = false, accepts = 0}, impl = {min_access_size = 1, max_access_size =
  read = {0, 0, 0}, write = {0, 0, 0}}}
绕了这么多，memory_region_read_accessor里的mr->ops->read终于到了我们注册的函数，如mc146818的cmos_ioport_read
(opaque=0x555556a17ae0, addr=1, size=1)
```

```
1 static void memory_region_read_accessor(void *opaque,
2                                         hwaddr addr,
3                                         uint64_t *value,
4                                         unsigned size,
5                                         unsigned shift,
6                                         uint64_t mask)
7 {
8     MemoryRegion *mr = opaque;
9     uint64_t tmp;
10
11     if (mr->flush_coalesced_mmio) {
12         qemu_flush_coalesced_mmio_buffer();
13     }
14     tmp = mr->ops->read(mr->opaque, addr, size);
15     *value |= (tmp & mask) << shift;
16 }
```

read/write回调函数，就是纯功能逻辑，比如mc146818主要干注入时钟，写入寄存器，读取日期等。

[PDF下载](#)

此条目是由 admin 发表在 未分类 分类目录的。将固定链接 <http://mnstory.net/2014/10/qemu-device-simulation/> 加入收藏夹。