

# FreeBSD 系统手册

## 摘要

欢迎到 FreeBSD 系统手册。 本手册在不断由多人编写。 多章是空白，有的章节亟待更新。 如果你对某个主题感兴趣并愿意有所贡献， 请参阅 [FreeBSD 文档文件列表](#)。

本文的最新英文原始版本可从 [FreeBSD Web 站点](#) 得到， 由 <http://www.FreeBSD.org.cn> 的最新版本可以在 <http://www.FreeBSD.org.cn> 快照 Web 站点 和 <http://www.FreeBSD.org.cn> 文档快照 得到， 一本会不断向主站同步。 此外， 也可以从 [FreeBSD FTP 服务器](#) 或更多的 [图像站点](#) 得到本文的各种其他格式以及多种形式的版本。

---

# 目

# 录

I: 内核	5
1. 引程与内核初始化	6
1.1. 概述	6
1.2. 环	6
1.3. BIOS POST	7
1.4. boot0段	7
1.5. boot2段	8
1.6. loader段	11
1.7. 内核初始化	11
2. 内核中的环	21
2.1. Mutex	21
2.2. 共享互斥环	23
2.3. 原子保环量	23
3. 内核对象	24
3.1. 环	24
3.2. Kobj的工作流程	24
3.3. 使用Kobj	24
4. Jail子系统	29
4.1. Jail的系环	29
4.2. 系环被囚禁程序的限制	35
5. SYSINIT框架	41
5.1. 环	41
5.2. SYSINIT操作	41
5.3. 使用SYSINIT	41
6. TrustedBSD MAC 框架	44
6.1. MAC 文版声明	44
6.2. 环解析	45
6.3. 概述	45
6.4. 安全策略背景知	45
6.5. MAC 框架的内核体系环	45
6.6. MAC策略模体环	49
6.7. MAC策略入口函数参考	51
6.8. 环用体环	109
6.9. 小环	110
7. 虚环内存系统	111
7.1. 物理内存的管理- <code>vm_page_t</code>	111
7.2. 环一的环存信息环体- <code>vm_object_t</code>	112
7.3. 文件系环入/环出- <code>buf</code> 环体	112

7.4. 映射表- <code>vm_map_t</code> , <code>vm_entry_t</code>	112
7.5. KVM存映射	112
7.6. 整FreeBSD的虚内存系	113
8. SMPNgin文	114
8.1. 介绍	114
8.2. 基本工具与上层的基础知识	114
8.3. 架构与CPU概况	115
8.4. 特定数据的策略	117
8.5. 会议说明	120
8.6. 其它	122
命令表	122
II: 内核程序	124
I写 FreeBSD 内核程序	125
1. 介绍	125
2. 内核直接工具-KLD	125
3. 宏内核程序	127
4. 字符驱动	127
5. 网络(消亡中)	135
6. 网卡驱动程序	135
9. ISA驱动程序	136
9.1. 概述	136
9.2. 基本信息	136
9.3. Device_t指针	138
9.4. 配置文件与自配置期间和探测的顺序	138
9.5. 源	140
9.6. 内存映射	142
9.7. DMA	148
9.8. xxx_isa_probe	150
9.9. xxx_isa_attach	157
9.10. xxx_isa_detach	160
9.11. xxx_isa_shutdown	161
9.12. xxx_intr	161
10. PCI	163
10.1. 探测与连接	163
10.2. 源	167
11. 通用方法SCSI控制器	171
11.1. 提交	171
11.2. 通用基础	171
11.3. 指令	191
11.4. 事件	192
11.5. 中断	193

11.6. 超I/O . . . . .	200
11.7. 超IO理 . . . . .	201
12. USBIO . . . . .	202
12.1. 简介 . . . . .	202
12.2. 主控器 . . . . .	202
12.3. USBIO信息 . . . . .	204
12.4. IO的探针和连接 . . . . .	205
12.5. USBIO程序的IO信息 . . . . .	206
13. Newbus . . . . .	208
13.1. IO程序 . . . . .	208
13.2. Newbus概 . . . . .	208
13.3. Newbus API . . . . .	210
14. 声音子系 . . . . .	212
14.1. 简介 . . . . .	212
14.2. 文件 . . . . .	212
14.3. 探针, 连接等 . . . . .	212
14.4. 接口 . . . . .	213
15. PC Card . . . . .	219
15.1. 添加IO . . . . .	219
III: 附 . . . . .	224
参考目 . . . . .	225

# 部分 I: 内核

# Chapter 1. 引入与内核初始化

## 1.1. 概述

第一章是关于引入与内核初始化的。这些过程始于BIOS(固件)POST，直到第一个用户程序建立。由于系统的最初阶段是与硬件相关的、是配合的，这里用IA-32(Intel Architecture 32bit)作为例子。

## 1.2. 引入

一台运行FreeBSD的计算机有多引入方法。这里介绍其中最通常的方法，也就是从安装了操作系统的硬盘上引入。引入过程分几步完成：

- BIOS POST
- `boot0`段
- `boot2`段
- loader段
- 内核初始化

`boot0`和`boot2`段在手册 [boot\(8\)](#)中被称为`bootstrap stages 1 and 2`，是FreeBSD的3段引入过程的开始。在一阶段都有各自的信息显示在屏幕上，可以参考下表列出一些。注意显示的内容可能随机器的不同而有一些区别：

不同机器而定	BIOS(固件)消息
F1 FreeBSD F2 BSD F5 Disk 2	<code>boot0</code>
>>FreeBSD/i386 BOOT Default: 1:ad(1,a)/boot/loader boot:	<code>boot2</code>

```
BTX loader 1.0 BTX version is 1.01
BIOS drive A: is disk0
BIOS drive C: is disk1
BIOS 639kB/64512kB available memory
FreeBSD/i386 bootstrap loader, Revision 0.8
Console internal video/keyboard
(jkh@bento.freebsd.org, Mon Nov 20 11:41:23 GMT 2000)
/kernels text=0x1234 data=0x2345 sysms=[0x4+0x3456]
Hit [Enter] to boot immediately, or any other key for command prompt
Booting [kernel] in 9 seconds...
```

loader

```
Copyright (c) 1992-2002 The FreeBSD Project.
Copyright (c) 1979, 1980, 1983, 1986, 1988, 1989, 1991, 1992, 1993, 1994
    The Regents of the University of California. All rights reserved.
FreeBSD 4.6-RC #0: Sat May 4 22:49:02 GMT 2002
    devnull@kukas:/usr/obj/usr/src/sys/DEVNUL
Timecounter "i8254" frequency 1193182 Hz
```

内核

## 1.3. BIOS POST

当PC加电后，处理器的寄存器被置为某些特定值。在这些寄存器中，**cr1**指令指针寄存器被置为32位0xfffffff0。指令指针寄存器指向处理器将要执行的指令代码。**cr1**是一个32位控制寄存器，在地址0被置为0。cr1的PE(Protected Enabled，保护模式使能)位用来指示处理器是处于保护模式还是地址模式。由于高位被清零，处理器在地址模式中引用。在地址模式中，逻辑地址与物理地址是等同的。

0xfffffff0略小于4G，因此计算机没有4G字节物理内存，就不会是一个有效的内存地址。计算机硬件将这个地址指向BIOS存储器。

BIOS表示Basic Input Output System (基本输入输出系统)。在主板上，它被固化在一个相对容量较小的只读存储器(Read-Only Memory, ROM)。BIOS包含各主板硬件定制的底层例程。就CPU而言，处理器首先指向常驻BIOS存储器的地址 0xfffffff0。通常这个位置包含一条跳转指令，指向BIOS的POST例程。

POST表示Power On Self Test(加电自检)。这套程序包括内存测试、系统启动和其他底层工具，从而使得CPU能够初始化整台计算机。第一段中有一个重要部分，就是启动引导。在所有的BIOS都允许手工干预。可以从软盘、光驱、硬盘等引入。

POST的最后一部分是执行INT 0x19指令。这个指令从引导第一个扇区读取512字节，装入地址0x7c00。第一个扇区的方法最早起源于硬盘的分区，硬盘面被分为若干柱形道。每道有数个扇区，同道又将扇区分成一定数目(通常是64)的扇形。0号扇区是硬盘的最外圈，1号扇区，第一个扇区(道、柱面都从0开始)而扇区从10开始)有着特殊的作用，它又被称为主引导记录(Master Boot Record, MBR)。第一剩余的扇区常常不使用。

## 1.4. boot0段

我们看一下文件/boot/boot0。这是一个512字节的小文件。如果在FreeBSD安装过程中选择"bootmanager"，这个文件中的内容将被写入硬盘MBR。

如前所述，`INT 0x19` 指令装`MBR`，也就是`boot0` 的内容至内存地址 `0x7c00`。再看文件 `sys/boot/i386/boot0/boot0.S`，可以猜想`里面`生了什`-`是引`管理器`，一段由 Robert Nordier`写的令人起敬的程序片段。`

MBR里，也就是`boot0`里，从偏移量`0x1be`始有一个特殊的`分区表`。其中有4条`分区`(称`分区`)，`16字`。分区`表示硬`如何被`分`，在FreeBSD的`中`，`被称`slice(d)。`16字`中有一个`志字`决定`个分区是否可引`。有`只能有一个分区可定一志`。否`， boot0的代将拒`行。

一个分区`有如下域`：

- 1字`文件系`型
- 1字`可引`志
- 6字`CHS格式描述符`
- 8字`LBA格式描述符`

一个分区`描述符包含某一分区在硬`上的`切位置信息`。

CHS`描述符指示相同的信息，但是指示方式有所不同：LBA (物理地址， Logical Block Addressing)指示分区的起始扇区和分区长度，而CHS(柱面 磁道 扇区)指示首扇区和末扇区`

引`管理器`描分区表，并在屏幕上`示菜`，以便用`可以`用于引`的磁`和分区。在`上按下相`的`后`，`boot0`行如下`作`：

- `中的分区可引`，清除以前的可引`志`
- `住本次的分区以`下次引`作缺省`
- 装`中分区的第一个扇区，并跳`行之

什`数据会存在`于一个可引`扇区`(`里指FreeBSD扇区`)的第一扇区里`？正如已猜到的，那就是boot2。`

## 1.5. boot2`段`

也`想知道，什`boot2`是在`boot0`之后，而不是在boot1之后。事`上，`也有一个512字`的文件 boot1存放在目`/boot`里，那是用来从一`引`系`的`。从`引`，boot1起着 boot0`硬`引`相同的作用`:它`到`boot2并`行`之。

可能已`看到`有一文件`/boot/mbr`。`是boot0的化版本`。mbr中的代`不会`示菜`用`，而只是`的引`被`志的分区`。

boot2的代`存放在目` sys/boot/i386/boot2`里，`的可`行文件在`/boot`里。在`/boot`里的文件` boot0 和boot2`不会在引`程中使用，只有boot0cfg`的工具才会使用它`。boot0`的内容`在MBR中才能生效。boot2位于可引`的FreeBSD分区的`始。`些位置不受文件系`控制，所以它`不可用ls`之`的命令`看。

boot2的主要任`是装`文件 /boot/loader，那是引`程的第三`段。在boot2中的代`不能使用`如 `open()` 和`read()` 之`的例程函数`，因`内核`没有被加`。而`当`描硬`，`取文件系`，`到文件`/boot/loader，用BIOS的功能将它`入内存`，然后从其入口点`始`行之。

除此之外，boot2`可提示用`行`， loader可以从其它磁`、系`元、分区装`。

boot2`的二制代`用特殊的方式`生`：

```

sys/boot/i386/boot2/Makefile
boot2: boot2.ldr boot2.bin ${BTX}/btv/btx
    btxld -v -E ${ORG2} -f bin -b ${BTX}/btv/btx -l boot2.ldr \
        -o boot2.ld -P 1 boot2.bin

```

这个Makefile片段表明**btxld(8)**被用来连接二进制代码。BTX表示引导扩展器(Boot eXtender)是程序(称作客户(client))提供保护模式环境、并与客户程序相连接的一段代码。所以 **boot2**是一个BTX客户，使用BTX提供的服务。

工具**btxld**是连接器，它将多个二进制代码连接在一起。**btxld(8)**和**ld(1)**的区别是**ld**通常将多个目标文件连接成一个可执行文件，而**btxld**将一个目标文件与BTX连接起来，生成合于放在分区首部的二进制代码，以系统调用。

**boot0**行跳转至BTX的入口点。然后，BTX将处理器切换至保护模式，并准备一个环境，然后调用客户。这个环境包括：

- 虚拟8086模式。这意味着BTX是虚拟8086的程序。模式指令，如pushf, popf, cli, sti, if，均可被客户调用。
- 建立中断描述符表(Interrupt Descriptor Table, IDT)，使得所有的硬件中断可被缺省的BIOS程序处理。建立中断0x30，是系统调用口。
- 一个系统调用**exec**和**exit**的定义如下：

```

sys/boot/i386/btx/lib/btxsys.s:
    .set INT_SYS,0x30      # 中断号
#
# System call: exit
#
__exit:    xorl %eax,%eax      # BTX系统调用0x0
           int $INT_SYS       #
#
# System call: exec
#
__exec:   movl $0x1,%eax      # BTX系统调用0x1
           int $INT_SYS       #

```

BTX建立全局描述符表(Global Descriptor Table, GDT)：

```

sys/boot/i386/btx/btx.s:
gdt:      .word 0x0,0x0,0x0,0x0      # 以空入口
          .word 0xffff,0x0,0x9a00,0xcf    # SEL_SCODE
          .word 0xffff,0x0,0x9200,0xcf    # SEL_SDATA
          .word 0xffff,0x0,0x9a00,0x0 # SEL_RCODE
          .word 0xffff,0x0,0x9200,0x0 # SEL_RDATA
          .word 0xffff,MEM_USR,0xfa00,0xcf# SEL_UCODE
          .word 0xffff,MEM_USR,0xf200,0xcf# SEL_UDATA
          .word _TSSL,0x8900,0x0 # SEL_TSS

```

客口的代口和数据始于地址MEM\_USR(0xa000), 口符(selector) SEL\_UCODE指向客口的数据段。口符 SEL\_UCODE 口有第3口描述符限 (Descriptor Privilege Level, DPL), 口是最低口限。但是 INT 0x30 指令的 理程序存口于口一个段里, 口个段的口符SEL\_SCODE (supervisor code)由有着管理口限。正如代口建立 IDT(中断描述符表)口行的操作那口:

```

    mov $SEL_SCODE,%dh      # 段口符
init.2:   shr %bx          # 是否理个中断?
    jnc init.3            # 否
    mov %ax,(%di)          # 置理程序偏移量
    mov %dh,0x2(%di)        # 置理程序口符
    mov %dl,0x5(%di)        # 置 P:DPL:type
    add $0x4,%ax           # 下一个中断理程序

```

所以, 当客口用 \_exec(), 代口将被以最高口限行。 口使得内核可以修改保口模式数据口, 如分口表(page tables)、全局描述符表(GDT)、中断描述符表(IDT)等。

**boot2** 定口了一个重要的数据口: **struct bootinfo**。口个口由 **boot2** 初始化, 然后被口送到loader, 之后又被口入内核。口个口的部分目由**boot2**定, 其余的由loader定。口个口中的信息包括内核文件名、BIOS提供的硬口柱面/磁口/扇区数目信息、口 BIOS提供的引口的口器号, 可用的物理内存大小, **envp** 指口(境指口)等。定口如下:

```

/usr/include/machine/bootinfo.h
struct bootinfo {
    u_int32_t bi_version;
    u_int32_t bi_kernelname; /* 用一个字表示 * */
    u_int32_t bi_nfs_diskless; /* struct nfs_diskless * */
    /* 以上常 */
#define bi_endcommon bi_n_bios_used
    u_int32_t bi_n_bios_used;
    u_int32_t bi_bios_geom[N_BIOS_GEOM];
    u_int32_t bi_size;
    u_int8_t bi_memsizes_valid;
    u_int8_t bi_bios_dev; /* 引口的BIOS元号 */
    u_int8_t bi_pad[2];
    u_int32_t bi_basemem;
    u_int32_t bi_extmem;
    u_int32_t bi_symtab; /* struct symtab * */
    u_int32_t bi_esymtab; /* struct symtab * */
    /* 以下目高bootloader提供 */
    u_int32_t bi_kernend; /* 内核空末端 */
    u_int32_t bi_envp; /* 境 */
    u_int32_t bi_modulep; /* 装的模 */
};


```

**boot2** 口入一个循口等待用口入, 然后口用 **load()**。如果用口不做任何口入, 循口将在一段口后口束, **load()** 将会装口缺省文件(/boot/loader)。函数 **ino\_t lookup(char \*filename)** 和 **int xfsread(ino\_t inode, void \*buf, size\_t nbytes)** 用来将文件内容口入内存。/boot/loader是一个ELF格式二口制文件, 不口它的部被口成了a.out格式中的**struct exec**口。**load()**描loader的ELF部, 装口/boot/loader 至内存, 然后跳

到入口执行之：

```
sys/boot/i386/boot2/boot2.c:  
    __exec((caddr_t)addr, RB_BOOTINFO | (opts RBX_MASK),  
        MAKEBOOTDEV(dev_maj[dsk.type], 0, dsk.slice, dsk.unit, dsk.part),  
        0, 0, 0, VTOP(bootinfo));
```

## 1.6. loader段

loader也是一个 BTX 客户，在这里不作描述。已有一部内容全面的手册 [loader\(8\)](#)，由Mike Smith编写。比loader更低层的BTX的机理已在前面讲过。

loader的主要任务是引导内核。当内核被装入内存后，即被loader使用：

```
sys/boot/common/boot.c:  
/* 从loader中调用内核中的exec程序 */  
module_formats[km->loader]-l_exec(km);
```

## 1.7. 内核初始化

我们来看一下链接内核的命令。能帮助我了解 loader 确定内核的准确位置。那个位置就是内核真正的入口点。

```
sys/conf/Makefile.i386:  
ld -elf -Bdynamic -T /usr/src/sys/conf/ldscript.i386 -export-dynamic \  
-dynamic-linker /red/herring -o kernel -X locore.o \  
lots of kernel .o files
```

在一行中有一些有趣的東西。首先，内核是一个ELF二进制文件，可是连接器却是 /red/herring，一个莫须有的文件。其次，看一下文件sys/conf/ldscript.i386，可以理解内核ld的有些什么。最前几行，字符串

```
sys/conf/ldscript.i386:  
ENTRY(btext)
```

表示内核的入口点是符号 **btext**。这个符号在locore.s中定义：

```

sys/i386/i386/locore.s:
.text
*****
*
* This is where the bootblocks start us, set the ball rolling...
* 入口
*/
NON_GPROF_ENTRY(bttext)

```

首先将寄存器EFLAGS置一个固定的0x00000002，然后初始化所有段寄存器：

```

sys/i386/i386/locore.s
/* 不要相信BIOS给出的EFLAGS */
    pushl $PSL_KERNEL
    popfl

/*
 * 不要相信BIOS给出的%fs、%gs。相信引导程序中定的%cs、%ds、%es、%ss
 */
    mov %ds, %ax
    mov %ax, %fs
    mov %ax, %gs

```

bttext用例程recover\_bootinfo(), identify\_cpu(), create\_pagetables()。有些例程也定在locore.s之中。有些例程的功能如下：

[recover\\_bootinfo](#)

该例程分析由引导程序送内核的参数。引导内核有3种方式：由loader引导(如前所述)，由老式磁盘引导，无引导方式。该函数决定引导方式，并将struct bootinfo存至内核内存。

[identify\\_cpu](#)

该函数CPU型，结果存放在变量(cpu)中。

[create\\_pagetables](#)

该函数分配表在内核内存空部分配一块空，并填写一定内容

下一节是VME(如果CPU有这个功能)：

```

testl $CPUID_VME, R(_cpu_feature)
jz 1f
movl %cr4, %eax
orl $CR4_VME, %eax
movl %eax, %cr4

```

然后，分模式：

```

/* Now enable paging */
    movl    R_IdlePTD, %eax
    movl    %eax,%cr3           /* load ptd addr into mmu */
    movl    %cr0,%eax          /* get control word */
    orl $CR0_PE|CR0_PG,%eax   /* enable paging */
    movl    %eax,%cr0          /* and let's page NOW! */

```

由于分段模式已启动，原先的段地址寻址方式随即失效。随后三行代码用来跳转至虚地址：

```

pushl  $begin           /* jump to high virtualized address */
ret

/* 在跳转至KERNBASE， 那里是操作系统内核被链接后真正的入口 */
begin:

```

函数`init386()`被调用；随参数传入的是一个指针，指向第一个空物理页。随后执行`mi_startup()`。  
`init386`是一个与硬件系统相关的初始化函数，`mi_startup()`是个与硬件系统无关的函数（前缀'mi'\_表示Machine Independent，不依赖于机器）。内核不再从`mi_startup()`里返回；调用这个函数后，内核完成引导：

```

sys/i386/i386/locore.s:
    movl    physfree, %esi
    pushl  %esi      /* 送给init386()的第一个参数 */
    call    _init386  /* 设置386芯片使之能UNIX工作 */
    call    _mi_startup /* 自己配置硬件，挂接根文件系统，等 */
    hlt     /* 不再返回到这里！ */

```

### 1.7.1. init386()

`init386()`定义在`sys/i386/i386/machdep.c`中，它对Intel 386芯片进行低级初始化。loader已将CPU切换至保护模式。loader已建立了最早的任务。



作者注

每个任务都是与其它“任务”相对独立的运行环境。任务之间可以分割切合，并且程序/线程的提供了必要基础。对于Intel 80x86任务的描述，对Intel公司于80386 CPU及后产品的材料，或者在[清华大学](#)收藏中用“80386”作关键词所找到的系方面的目。

在多个任务中，内核将工作。在其代码前，我将处理器保护模式必须完成的一系列准备工作一并列出：

- 初始化内核的整体参数，有些参数由引导程序传来
- 准备GDT(全局描述符表)
- 准备IDT(中断描述符表)
- 初始化系统控制台
- 初始化DDB(内核的调试器)，如果它被内核的
- 初始化TSS(任务状态段)

- 准备GDT(全局描述符表)
- 建立proc0(0号进程, 即内核的进程)的pcb(进程控制块)

`init386()`首先初始化内核的可调整参数, 有些参数由引导程序来。先设置环境指针(environment pointer, envp)用, 再用`init_param1()`。envp指针已由loader存放在bootinfo中:

```
sys/i386/i386/machdep.c:
kern_envp = (caddr_t)bootinfo.bi_envp + KERNBASE;

/* 初始化基本可调整, 如hz等 */
init_param1();
```

`init_param1()`定在 sys/kern/subr\_param.c之中。这个文件里有一些sysctl, 有10个函数, `init_param1()`和`init_param2()`。10个函数从`init386()`中用:

```
sys/kern/subr_param.c
hz = HZ;
TUNABLE_INT_FETCH("kern.hz", hz);
```

TUNABLE\_typename\_FETCH用来取环境量的:

```
/usr/src/sys/sys/kernel.h
#define TUNABLE_INT_FETCH(path, var) getenv_int((path), (var))
```

Sysctl`kern.hz`是系统频率。同时, 有些sysctl被`init_param1()`定: `kern.maxswzone`, `kern.maxbcache`, `kern.maxtsiz`, `kern.dfldsiz`, `kern.maxdsiz`, `kern.dflssiz`, `kern.maxssiz`, `kern.sgrowsiz`。

然后`init386()`准备全局描述符表 (Global Descriptors Table, GDT)。在x86上每个任务都运行在自己的虚拟地址空间里, 每个空间由"段址:偏移量"的数指定。一个例子, 当前将要由处理器执行的指令在 CS:EIP, 那条指令的物理虚地址就是"代码段虚地址CS" + EIP。为了方便, 段起始于虚地址0, 止于界限4G字节。所以, 在一个例子中, 指令的物理虚地址正是EIP的值。段寄存器, 如CS、DS等是16位, 即全局描述符表中的索引(更精确的说, 索引并非16位的全部, 而是16位中的INDEX部分)。



#### 作者注

对于80386, 16位有16位, INDEX部分是其中的高13位。

FreeBSD的全局描述符表16个CPU保存着15个16位:

```

sys/i386/i386/machdep.c:
union descriptor gdt[NGDT * MAXCPU]; /* 全局描述符表 */

sys/i386/include/segments.h:
/*
 * 全局描述符表(GDT)中的入口
 */
#define GNULL_SEL 0 /* 空描述符 */
#define GCODE_SEL 1 /* 内核代码描述符 */
#define GDATA_SEL 2 /* 内核数据描述符 */
#define GPRIV_SEL 3 /* 称多处理器(SMP)物理器有数据 */
#define GPROC0_SEL 4 /* Task state process slot zero and up, 任务状态 */
#define GLDT_SEL 5 /* 一个程序的局部描述符表 */
#define GUSERLDT_SEL 6 /* 用户自定义的局部描述符表 */
#define GTGATE_SEL 7 /* 程序切换门 */
#define GBIOSLOWMEM_SEL 8 /* BIOS低端内存门(必然是第8个入口) */
#define GPANIC_SEL 9 /* 会致全系常中止工作的任务 */
#define GBIOSCODE32_SEL 10 /* BIOS接口(32位代码) */
#define GBIOSCODE16_SEL 11 /* BIOS接口(16位代码) */
#define GBIOSDATA_SEL 12 /* BIOS接口(数据) */
#define GBIOSUTIL_SEL 13 /* BIOS接口(工具) */
#define GBIOSARGS_SEL 14 /* BIOS接口(参数) */

```

注意，这些#definees并非#符本身，而只是#符中的INDEX域，因此它正是全局描述符表中的索引。例如，内核代码的#符(GCODE\_SEL)的#0x08。

下一节是初始化中断描述符表(Interrupt Descriptor Table, IDT)。该表在发生软件或硬件中断时会被处理器引用。例如，一行系用，用#用程序提交INT 0x80 指令。这是一个软件中断，处理器用索引#0x80在中断描述符表中#。一个#指向处理器的一个中断的例程。在#个特定情形中，#是内核的系用#口。

### 作者注



Intel 80386支持“不用#”，可以使得用#程序只通过一条call指令就#用内核中的例程。可是FreeBSD并未采用#机制，也是因为使用#中断接口可免去#接的麻烦。另外有一个附好的：在#真Linux#，当遇到FreeBSD内核不支持的而又并非#性的系用#，内核只会#示一些出#信息，#使得程序能#行；而不是在真正#行程序之前的初始化#程中就因#接失#而不允#程序#行。

中断描述符表最多可以有256 (0x100)条#。内核分配NIDT条#的内存#中断描述符表，#里NIDT=256，是最大#：

```

sys/i386/i386/machdep.c:
static struct gate_descriptor idt0[NIDT];
struct gate_descriptor *idt = idt0[0]; /* 中断描述符表 */

```

#个中断都被#置一个合#的中断#理程序。系用#口INT 0x80也是如此：

```
sys/i386/i386/machdep.c:  
    setidt(0x80, IDTVEC(int0x80_syscall),  
           SDT_SYS386TGT, SEL_UPL, GSEL(GCODE_SEL, SEL_KPL));
```

所以当一个用应用程序提交INT 0x80指令时，全系统的控制权会会调用函数Xint0x80\_syscall，该函数在内核代码段中，将被以管理权限执行。

然后，控制台和DDB(调试器)被初始化：

```
sys/i386/i386/machdep.c:  
    cninit();  
/* 以下代码可能因未定义宏DDB而被跳过 */  
#ifdef DDB  
    kdb_init();  
    if (boothowto & RB_KDB)  
        Debugger("Boot flags requested debugger");  
#endif
```

任务状态段(TSS)是一个x86保护模式中的数据段。当发生任何切换，任务状态段用来硬件存取任务信息。

局部描述符表(LDT)用来指向用代码和数据。系统定义了几个常量，指向局部描述符表，它是系统用和用代码、用数据常量：

```
/usr/include/machine/segments.h  
#define LSYS5CALLS_SEL 0 /* Intel BCS限制要求的 */  
#define LSYS5SIGR_SEL 1  
#define L43BSDCALLS_SEL 2 /* 尚无 */  
#define LUCODE_SEL 3  
#define LSOL26CALLS_SEL 4 /* Solaris =2.6版系用 */  
#define LUDATA_SEL 5  
/* separate stack, es, fs, gs sels ? 分开的、es、fs、gs常量 ? */  
/* #define LPOSIXCALLS_SEL 5 */ /* notyet, 尚无 */  
#define LBSDICALLS_SEL 16 /* BSDI system call gate, BSDI系用 */  
#define NLDT (LBSDICALLS_SEL + 1)
```

然后，proc0(0号进程，即内核所用的进程)的进程控制块(Process Control Block) (struct pcb)被初始化。proc0是一个 struct proc 常量，描述了一个内核进程。内核运行，进程是存在，所以这个常量在内核中被定为全局常量：

```
sys/kern/kern_init.c:  
    struct proc proc0;
```

struct pcb是proc的一部分，它定义在/usr/include/machine/pcb.h之中，内含i386硬件有的信息，如寄存器的值。

### 1.7.2. mi\_startup()

这个函数用冒泡排序算法，将所有系统初始化对象，然后逐个调用它们的入口：

```
sys/kern/init_main.c:  
for (sipp = sysinit; *sipp; sipp++) {  
  
    /* ... 省略 ... */  
  
    /* 调用函数 */  
    (*((*sipp)-func))((*sipp)-udata);  
    /* ... 省略 ... */  
}
```

尽管sysinit框架已在《FreeBSD开发者手册》中有所描述，我是在这里介绍一下其内部原理。

这个系统初始化对象(sysinit对象)通常用宏建立。我将以announce sysinit对象为例。这个对象打印版权信息：

```
sys/kern/init_main.c:  
static void  
print_caddr_t(void *data __unused)  
{  
    printf("%s", (char *)data);  
}  
SYSINIT(announce, SI_SUB_COPYRIGHT, SI_ORDER_FIRST, print_caddr_t, copyright)
```

这个对象的子系统是SI\_SUB\_COPYRIGHT(0x0800001)，数正好排在SI\_SUB\_CONSOLE(0x0800000)后面。所以，版权信息将在控制台初始化之后就被很早的打印出来。

我们看一看宏SYSINIT()到底做了些什么。它展开成宏C\_SYSINIT()。宏C\_SYSINIT()然后展开成一个静态struct sysinit。里面申明里用了两个宏 DATA\_SET:

```
/usr/include/sys/kernel.h:  
#define C_SYSINIT(uniquifier, subsystem, order, func, ident) \  
static struct sysinit uniquifier ## _sys_init = { \ subsystem, \  
order, \ func, \ ident \ }; \ DATA_SET(sysinit_set,uniquifier ## \  
_sys_init);  
  
#define SYSINIT(uniquifier, subsystem, order, func, ident) \  
C_SYSINIT(uniquifier, subsystem, order,           \  
(sysinit_cfunc_t)(sysinit_nfunc_t)func, (void *)ident)
```

宏DATA\_SET()展开成MAKE\_SET()，宏MAKE\_SET()指向所有包含的sysinit函数：

```

/usr/include/linker_set.h
#define MAKE_SET(set, sym) \
    static void const * const __set_##set##_sym##sym = sym; \
    __asm(".section .set." #set ",aw\""); \
    __asm(".long " #sym); \
    __asm(".previous")
#endif
#define TEXT_SET(set, sym) MAKE_SET(set, sym)
#define DATA_SET(set, sym) MAKE_SET(set, sym)

```

回到我的例子中，`#define`宏的展开过程，将会产生如下声明：

```

static struct sysinit announce_sys_init = {
    SI_SUB_COPYRIGHT,
    SI_ORDER_FIRST,
    (sysinit_cfunc_t)(sysinit_nfunc_t) print_caddr_t,
    (void *) copyright
};

static void const *const __set_sysinit_set_sym_announce_sys_init =
    announce_sys_init;
__asm(".section .set.sysinit_set", "aw\"");
__asm(".long " "announce_sys_init");
__asm(".previous");

```

第一个`__asm`指令在内核可执行文件中建立一个ELF的section。产生在内核直接的时候。它将被命令`.set.sysinit_set`。它的内容是一个32位——`announce_sys_init`的地址，它正是第二个`\__asm`指令所定位的。第三个`\__asm`指令结束。如果前面有名字相同的定义句，它的内容(那个32位)将被填加到已存在的里，就造出了一个32位指针数。

用objdump察看一个内核二进制文件，也会影响到里面有好几个小的：

```

% objdump -h /kernel
7 .set.cons_set 00000014 c03164c0 c03164c0 002154c0 2**2
    CONTENTS, ALLOC, LOAD, DATA
8 .set.kbddriver_set 00000010 c03164d4 c03164d4 002154d4 2**2
    CONTENTS, ALLOC, LOAD, DATA
9 .set.scrndr_set 00000024 c03164e4 c03164e4 002154e4 2**2
    CONTENTS, ALLOC, LOAD, DATA
10 .set.scterm_set 0000000c c0316508 c0316508 00215508 2**2
    CONTENTS, ALLOC, LOAD, DATA
11 .set.sysctl_set 0000097c c0316514 c0316514 00215514 2**2
    CONTENTS, ALLOC, LOAD, DATA
12 .set.sysinit_set 00000664 c0316e90 c0316e90 00215e90 2**2
    CONTENTS, ALLOC, LOAD, DATA

```

一屏信息表示表明`.set.sysinit_set`有0x664字节的大小，所以`0x664/sizeof(void *)`个`sysinit`对象被

除了内核。其它的，如`.set.sysctl_set`表示其它的接器集合。

通常定义一个类型的`struct linker_set`的变量，`.set.sysinit_set`将被"收集"到那个变量里：

```
sys/kern/init_main.c:  
extern struct linker_set sysinit_set; /* XXX */
```

`struct linker_set`定义如下：

```
/usr/include/linker_set.h:  
struct linker_set {  
    int ls_length;  
    void *ls_items[1];      /* ls_length个指向对象的指针，以NULL结尾 */  
};
```



作者注

以上是的，用C语言具体`linker_set`来表示那个ELF的。

第一是`sysinit`象的数量，第二是一个以NULL结尾的数组，数组中是指向那些对象的指针。

回到`mi_startup()`的，我清楚了`sysinit`象是如何被组织起来的。函数`mi_startup()`将它排序，并用一个象。最后一个象是调度器：

```
/usr/include/sys/kernel.h:  
enum sysinit_sub_id {  
    SI_SUB_DUMMY      = 0x0000000,      /* 不被执行，仅供接器使用 */  
    SI_SUB_DONE       = 0x0000001,      /* 已被整理 */  
    SI_SUB_CONSOLE    = 0x0800000,      /* 控制台 */  
    SI_SUB_COPYRIGHT  = 0x0800001,      /* 最早使用控制台的象 */  
    ...  
    SI_SUB_RUN_SCHEDULER = 0xffffffff /* 调度器：不返回 */  
};
```

调度器`sysinit`象定义在文件`sys/vm/vm_glue.c`中，这个象的入口点是`scheduler()`。这个函数上是个无限循环，它表示那个进程(PID)0的进程——swapper进程。前面提到的proc0正是用来描述这个进程。

第一个用的程序是`_init_`，由`sysinit`象`init`建立：

```

sys/kern/init_main.c:
static void
create_init(const void *userdata __unused)
{
    int error;
    int s;

    s = splhigh();
    error = fork1(proc0, RFFDG | RFPROC, initproc);
    if (error)
        panic("cannot fork init: %d\n", error);
    initproc->p_flag |= P_INMEM | P_SYSTEM;
    cpu_set_fork_handler(initproc, start_init, NULL);
    remrunqueue(initproc);
    splx(s);
}
SYSINIT(init, SI_SUB_CREATE_INIT, SI_ORDER_FIRST, create_init, NULL)

```

`create_init()`通常用`fork1()`分配一个新的进程，但并不将其置为可运行。当新进程被调度器调度运行时，`start_init()`将被调用。那个函数定义在`init_main.c`中。它装载并运行二进制代码`init`，先从`/sbin/init`，然后是`/sbin/oinit`，`/sbin/init.bak`，最后是`/stand/sysinstall`：

```

sys/kern/init_main.c:
static char init_path[MAXPATHLEN] =
#endif INIT_PATH
    __XSTRING(INIT_PATH);
#else
    "/sbin/init:/sbin/oinit:/sbin/init.bak:/stand/sysinstall";
#endif

```

# Chapter 2. 内核中的口

本章由 FreeBSD SMP Next Generation Project 完成。将口和建口送口称多口理 (SMP) 口件列表。

本文口提口述了在FreeBSD内核中的口，口些口使得有效的多口理成口可能。口可以用几口方式口得。数据口可以用mutex或lockmgr(9)保口。口于口数不多的若干个口量，假如口是使用原子操作口口它口，口些口量就可以得到保口。



## 口者注

口本章内容，口不足以口出"mutex" 和"共享互斥口"的区口。似乎它口的功能有重口之口，前者比后者的功能口口更多。它口似乎都是lockmgr(9)的子集。

## 2.1. Mutex

Mutex就是一口用来解决共享/排它矛盾的口。一个mutex在一个口刻只可以被一个口体口有。如果口一个口体要口得已口被口有的mutex，就会口入等待，直到口个mutex被口放。在FreeBSD内核中，mutex被口程所口有。

Mutex可以被口口的索要，但是mutex一般只被一个口体口有口短的一段口，因此一个口体不能在持有mutex口睡眠。如果口需要在持有mutex口睡眠，可使用一个 lockmgr(9) 的口。

口个mutex有几个令人感口趣的属性：

### 口量名

在内核源代口中struct mtx口量的名字

### 口口名

由函数mtx\_init指派的mutex的名字。口个名字口示在KTR跟踪消息和witness出口与警告信息里。口个名字口用于区分口口在witness代口中的各个mutex

### 口型

Mutex的口型，用口志MTX\_表示。口个口志的意口在mutex(9)有所描述。

#### MTX\_DEF

一个睡眠mutex

#### MTX\_SPIN

一个循环mutex

#### MTX\_RECURSE

口个mutex允口口口

### 保口口象

口个入口所要保口的数据口口列表或数据口口成口列表。口于数据口口成口，将按照 口名.成口名的形式命名。

### 依口函数

口当mutex被持有口才可以被口用的函数

表 1. Mutex列表

量名	名	型	保象	依函数
sched_lock	"sched lock"(调度器)	MTX_SPIN   MTX_RECURSE	_gmonparam, cnt.v_swtch, cp_time, curpriority, mtx .mtx_blocked, mtx .mtx_contested, proc.p_procq, proc .p_slpq, proc .p_sflag, proc .p_stat, proc .p_estcpu, proc .p_cpticks proc .p_pctcpu, proc .p_wchan, proc .p_wmesg, proc .p_swtime, proc .p_sptime, proc .p_runtime, proc .p_uu, proc.p_su, proc.p_iu, proc .p_uticks, proc .p_sticks, proc .p_iticks, proc .p_oncpu, proc .p_lastcpu, proc .p_rqindex, proc .p_heldmtx, proc .p_blocked, proc .p_mtxname, proc .p_contested, proc .p_priority, proc .p_usrpri, proc .p_nativepri, proc .p_nice, proc .p_rt prio, psnt, slpque, itqueuebits, itqueues, rtqueuebits, rtqueues, queuebits, queues, idqueuebits, idqueues, switchtime, switchticks	setrunqueue, remrunqueue, mi_switch, chooseproc, schedclock, resetpriority, updatepri, maybe_resched, cpu_switch, cpu_throw, need_resched, resched_wanted, clear_resched, aston, astoff, astpending, calcru, proc_compare

变量名	名	型	保对象	依函数
vm86pcb_lock	"vm86pcb lock"(虚8086模式编程控制)	MTX_DEF	vm86pcb	vm86_bioscall
Giant	"Giant"(巨)	MTX_DEF   MTX_RECURSE	几乎可以是任何东西	很多
callout_lock	"callout lock"(延用)	MTX_SPIN   MTX_RECURSE	callfree, callwheel, nextsoftcheck, proc.p_itcallout, proc.p_slpcallout, softticks, ticks	

## 2.2. 共享互斥

一些提供基本的读写型的功能，可以被一个正在睡眠的线程持有。它在它被移到lockmgr(9)之中。

表 2. 共享互斥列表

变量名	保对象
allproc_lock	allproc zombproc pidhashtbl proc.p_list proc.p_hash nextpid
proctree_lock	proc.p_children proc.p_sibling

## 2.3. 原子保对象

原子保对象并非由一个线程持有的特殊对象，而是：一些对象的所有数据都要使用特殊的原子操作([atomic\(9\)](#))。尽管其它的基本同步机制(例如mutex)就是用原子保对象做的，但是很少有对象直接使用物理方式。

- `mtx mtx_lock`

# Chapter 3. 内核对象

内核对象，也就是Kobj，内核提供了一种面向对象的C语言编程方式。被操作的数据也通过操作它的方法。使得在不破坏二进制兼容性的前提下，某一个接口能同时相互通作。

## 3.1. 对象

对象

数据集合-数据对象-数据分配的集合

方法

某一对象-函数

□

一个或多个方法

接口

一个或多个方法的一个标准集合

## 3.2. Kobj的工作流程



作者注

一小段落中原作者的用词有些含混，参考我在括号中的注释。

Kobj工作，生成方法的描述。每个描述有一个唯一的对象和一个缺省函数。某个描述的地址被用来在一个对象的方法表里唯一的对象方法。

建立一个，就是要建立一个方法表，并将对象表映射到一个或多个函数(方法)；  
有些函数(方法)都有方法描述。使用前，对象要被映射。对象要对象分配一些内存。  
在方法表中的每个方法描述都会被指派一个唯一的对象，除非已被其它引用它的对象指派了对象。由于  
每个将要被使用的方法，都会由脚本生成一个函数(方法对象函数)，以解析外来参数，并在被  
映射出方法描述的地址。被生成的函数(方法对象函数)凭着那个方法描述的唯一对象按Hash的方法对象象的内存。  
如果一个方法不在内存中，函数会直接使用对象的方法表。如果一个方法被映射到了，对象里的相关函数  
(也就是某个方法的代理)就会被使用。否则，一个方法描述的缺省函数将被使用。

流程可被表示如下：

对象-内存-对象

## 3.3. 使用Kobj

### 3.3.1. 对象

```
struct kobj_method
```

### 3.3.2. 函数

```
void kobj_class_compile(kobj_class_t cls);
void kobj_class_compile_static(kobj_class_t cls, kobj_ops_t ops);
void kobj_class_free(kobj_class_t cls);
kobj_t kobj_create(kobj_class_t cls, struct malloc_type *mtype, int mflags);
void kobj_init(kobj_t obj, kobj_class_t cls);
void kobj_delete(kobj_t obj, struct malloc_type *mtype);
```

### 3.3.3. 宏

```
KOBJ_CLASS_FIELDS
KOBJ_FIELDS
DEFINE_CLASS(name, methods, size)
KOBJMETHOD(NAME, FUNC)
```

### 3.3.4. 头文件

```
sys/param.h
sys/kobj.h
```

### 3.3.5. 建立一个接口的模板

使用Kobj的第一步是建立一个接口。建立接口包括建立模板的工作。

建立模板可用脚本src/sys/kern/makeobjops.pl完成，它会生成申明方法的头文件和代码，脚本会生成方法的函数。

在头个模板中如下部分会被使用: #include, INTERFACE, CODE, METHOD, STATICMETHOD, 和 DEFAULT.

#include句的整行内容将被一字不差的复制到被生成的代码文件的头部。

例如:

```
#include sys/foo.h
```

INTERFACE用来定接口名。一个名字将与一个方法名接合在一起，形成 [interface name]\_[method name]。语法是：INTERFACE [接口名];

例如:

```
INTERFACE foo;
```

CODE会将它的参数一字不差的复制到代码文件中。语法是CODE { [任何代码] };

例如：

```
CODE {
    struct foo * foo_alloc_null(struct bar *)
    {
        return NULL;
    }
};
```

METHOD用来描述一个方法。语法是：METHOD [返回类型] [方法名] { [对象 [, 参数若干]] };

例如：

```
METHOD int bar {
    struct object *;
    struct foo *;
    struct bar;
};
```

DEFAULT跟在METHOD之后，是METHOD的补充。它是个方法补充上缺省函数。语法是：METHOD [返回类型] [方法名] { [对象; [其它参数]] } DEFAULT [缺省函数];

例如：

```
METHOD int bar {
    struct object *;
    struct foo *;
    int bar;
} DEFAULT foo_hack;
```

STATICMETHOD类似于METHOD。对于一个Kobj对象，一般其内部都有一些Kobj特有的数据。METHOD定义的方法就假定有些数据位于对象内部；假如对象内部没有这些数据，这些方法对这个对象的调用就可能出错。而STATICMETHOD定义的对象可以不受这个限制：它描述出的方法，其操作的数据不由某个对象例例出，而是全都由用这个方法的操作数(或者注：即参数)给出。也在于在某个对象的方法表之外用这个方法有用。



者注

一段的语与原文相比整整很大。静态方法是不依于对象实例的方法。参看C++中的“静态函数”的概念。

其它完整的例子：

```
src/sys/kern/bus_if.m  
src/sys/kern/device_if.m
```

### 3.3.6. 建立一个对象

使用Kobj的第二步是建立一个对象。一个对象有名字、方法表；假如使用了Kobj的“对象管理工具”(Object Handling Facilities)，对象中包含对象的大小。建立对象使用宏`DEFINE_CLASS()`。建立方法表时，对象建立一个`kobj_method_t`数组，用NULL结尾。一个非NULL对象可用宏`KOBJMETHOD()`建立。

例如：

```
DEFINE_CLASS(fooClass, foomethods, sizeof(struct foodata));  
  
kobj_method_t foomethods[] = {  
    KOBJMETHOD(bar_doo, foo_doo),  
    KOBJMETHOD(bar_foo, foo_foo),  
    { NULL, NULL}  
};
```

对象被“对象”。根据对象被初始化对象的状况，将要用到一个静态分配的对象和“操作数表”(ops table，后者注：即“参数表”)。有些操作可通过声明一个对象体`struct kobj_ops`并使用`kobj_class_compile_static()`，或是只使用`kobj_class_compile()`来完成。

### 3.3.7. 建立一个对象

使用Kobj的第三步是定对象。Kobj对象建立程序假定Kobj对象有数据在一个对象的头部。如果不是如此，应当先自行分配对象，再使用`kobj_init()`初始化对象中的Kobj对象有数据；其后可以使用`kobj_create()`分配对象，并自对象初始化对象中的Kobj对象有内容。`kobj_init()`也可以用来改变一个对象所使用的。

将Kobj的数据集成到对象中要使用宏`KOBJ_FIELDS`。

例如

```
struct foo_data {  
    KOBJ_FIELDS;  
    foo_foo;  
    foo_bar;  
};
```

### 3.3.8. 对象用方法

使用Kobj的最后一部就是通过生成的函数对象用对象中的方法。对象用时，接口名与方法名用下划线接合，而且全部使用大写字母。

例如，接口名对象，方法对象，对象用就是：

```
[返回值 = ] FOO_BAR(对象 [, 其它参数]);
```

### 3.3.9. 善后整理

当一个用 `kobj_create()` 不再需要被使用时，可对这个对象用 `kobj_delete()`。当一个类不再需要被使用时，可对这个类用 `kobj_class_free()`。

# Chapter 4. Jail子系口

在大多数UNIX®系口中，用root是万能的。也就加了多危。

root，就可以在他的指尖掌握系口中所有的功能。

就可以将攻口者造成的口害口小到最低限度。

些安全功能中，有一口叫安全口。

4.0及以后版本中提供的安全功能，就是[jail\(8\)](#)。

如果一个攻口者得了一个系口中的

在FreeBSD里，有一些sysctl口削弱了root的口限，

一在FreeBSD

将一个口行口境的文件口根切口到某一特定位置，并且

一个被口禁的口程不能影响口个jail之外的

口程、不能使用一些特定的系口口用，也就不能口主口算机造成破坏。



口者注

英文口口"jail"的中文意思是"囚禁、口禁"。

Jail已口成口一新型的安全模型。人口可以在jail中口行各口可能很脆弱的服口器程序，如Apache、BIND和sendmail。口一来，即使有攻口者取得了jail中的root，最多口人口口眉口，而不会使人口口慌失措。本文主要口注jail的内部原理(源代口)。如果口正在口口口置Jail的指南性文口，我建口口口我的口一篇文章，口表在Sys Admin Magazine, May 2001, 《Securing FreeBSD using Jail》。

## 4.1. Jail的系口口

Jail由口部分口成：用口口程序，也就是[jail\(8\)](#)；口有在内核中Jail的口口代口：[jail\(2\)](#) 系口口用和相口的口束。我将口口用口口程序和jail在内核中的口口原理。

### 4.1.1. 用口口代口

Jail的用口口源代口在/usr/src/usr.sbin/jail，  
主机名，IP地址，口有需要口行的命令。

由一个文件jail.c口成。口个程序有口些参数：jail的路径，

#### 4.1.1.1. 数据口口

在jail.c中，我将最先注解的是一个重要口体  
`struct jail j;`的声明，口个口口型的声明包含在 /usr/include/sys/jail.h之中。

jail口口的定口是：

```
/usr/include/sys/jail.h:  
  
struct jail {  
    u_int32_t      version;  
    char          *path;  
    char          *hostname;  
    u_int32_t      ip_number;  
};
```

正如口所口，口送口命令[jail\(8\)](#)的口个参数都在口里有口口的一口。

口些参数才由命令行真正口入：

事口上，当命令[jail\(8\)](#)被口行口，

```

/usr/src/usr.sbin/jail.c
char path[PATH_MAX];
...
if(realpath(argv[0], path) == NULL)
    err(1, "realpath: %s", argv[0]);
if (chdir(path) != 0)
    err(1, "chdir: %s", path);
memset(j, 0, sizeof(j));
j.version = 0;
j.path = path;
j.hostname = argv[1];

```

#### 4.1.1.2. 网口

`jail(8)`的参数中有一个是IP地址。是在网上`jail`的地址。`jail(8)`将IP地址翻成网口字序，并存入`(jail型的)体`。

```

/usr/src/usr.sbin/jail/jail.c:
struct in_addr in;
...
if (inet_aton(argv[2], in) == 0)
    errx(1, "Could not make sense of ip-number: %s", argv[2]);
j.ip_number = ntohs(in.s_addr);

```

函数`inet_aton(3)`“将指定的字符串解为一个Internet地址，并将其存到指定的体中”。`inet_aton(3)`定了`体in`，之后`in`中的内容再用`ntohs(3)`成主机字序，并置入`jail`体的`ip_number`成。

#### 4.1.1.3. 囚禁程序

最后，用程序囚禁程序。在Jail自身成了一个被囚禁的程序，并使用`execv(3)`行用指定的命令。

```

/usr/src/usr.sbin/jail/jail.c
i = jail(j);
...
if (execv(argv[3], argv + 3) != 0)
    err(1, "execv: %s", argv[3]);

```

正如所，函数`jail()`被用，参数是`体jail`中被填入数据，而如前所述，这些数据又来自`jail(8)`的命令行参数。最后，行了用指定的命令。下面我将始`jail`在内核中的。

### 4.1.2. 相的内核源代码

在我来看文件`/usr/src/sys/kern/kern_jail.c`。在里定了`jail(2)`的系用、相的`sysctl`，有网函数。

#### 4.1.2.1. sysctl

在`kern_jail.c`里定了如下`sysctl`：

```
/usr/src/sys/kern/kern_jail.c:
```

```
int      jail_set_hostname_allowed = 1;
SYSCTL_INT(_security_jail, OID_AUTO, set_hostname_allowed, CTLFLAG_RW,
           jail_set_hostname_allowed, 0,
           "Processes in jail can set their hostnames");
/* Jail中的进程可以定自身的主机名 */

int      jail_socket_unixiproute_only = 1;
SYSCTL_INT(_security_jail, OID_AUTO, socket_unixiproute_only, CTLFLAG_RW,
           jail_socket_unixiproute_only, 0,
           "Processes in jail are limited to creating UNIX/IPv4/route sockets only");
/* Jail中的进程被限制只能建立UNIX套接字、IPv4套接字、路由套接字 */

int      jail_sysv ipc_allowed = 0;
SYSCTL_INT(_security_jail, OID_AUTO, sysv ipc_allowed, CTLFLAG_RW,
           jail_sysv ipc_allowed, 0,
           "Processes in jail can use System V IPC primitives");
/* Jail中的进程可以使用System V进程通信原 */

static int jail_enforce_statfs = 2;
SYSCTL_INT(_security_jail, OID_AUTO, enforce_statfs, CTLFLAG_RW,
           jail_enforce_statfs, 0,
           "Processes in jail cannot see all mounted file systems");
/* jail 中的进程看系统中挂接的文件系统受到何限制 */

int      jail_allow_raw_sockets = 0;
SYSCTL_INT(_security_jail, OID_AUTO, allow_raw_sockets, CTLFLAG_RW,
           jail_allow_raw_sockets, 0,
           "Prison root can create raw sockets");
/* jail 中的 root 用是否可以建 raw socket */

int      jail_chflags_allowed = 0;
SYSCTL_INT(_security_jail, OID_AUTO, chflags_allowed, CTLFLAG_RW,
           jail_chflags_allowed, 0,
           "Processes in jail can alter system file flags");
/* jail 中的进程是否可以修改系统文件 */

int      jail_mount_allowed = 0;
SYSCTL_INT(_security_jail, OID_AUTO, mount_allowed, CTLFLAG_RW,
           jail_mount_allowed, 0,
           "Processes in jail can mount/unmount jail-friendly file systems");
/* jail 中的进程是否可以挂或卸jail友好的文件系统 */
```

这些sysctl中的一个都可以用命令[sysctl\(8\)](#)。在整个内核中，  
。例如，上述第一个sysctl的名字是 [security.jail.set\\_hostname\\_allowed](#)。

这些sysctl按名称

#### 4.1.2.2. jail(2)系用

像所有的系用一， 系用jail(2)有0个参数， struct thread \*td和struct jail\_args \*uap。 td是一个指向thread体的指， 指用于描述用jail(2)的程。 在0个上下文中， uap指向一个体， 0个体中包含了一个指向从用jail送来的jail体的指。 在前面我述用程序， 已看到一个jail体被作参数送系用jail(2)。

```
/usr/src/sys/kern/kern_jail.c:  
/*  
 * struct jail_args {  
 *     struct jail *jail;  
 * };  
 */  
int  
jail(struct thread *td, struct jail_args *uap)
```

于是uap-jail可以用于被jail(2)的jail体。

然后， jail(2)使用copyin(9)将jail体制到内核内存空中。 copyin(9)需要三个参数：要制内核内存空的数据的地址 uap-jail，在内核内存空存放数据的j， 以及数据的大小。 uap-jail指向的Jail体被制内核内存空，并被存放在一个jail体j里。

```
/usr/src/sys/kern/kern_jail.c:  
error = copyin(uap-jail, j, sizeof(j));
```

在jail.h中定0了一个重要的体型prison。 体prison只被用在内核空中。 下面是prison体的定0。

```
/usr/include/sys/jail.h:  
struct prison {  
    LIST_ENTRY(prison) pr_list;           /* (a) all prisons */  
    int             pr_id;                /* (c) prison id */  
    int             pr_ref;               /* (p) refcount */  
    char            pr_path[MAXPATHLEN];  /* (c) chroot path */  
    struct vnode   *pr_root;              /* (c) vnode to rdir */  
    char            pr_host[MAXHOSTNAMELEN]; /* (p) jail hostname */  
    u_int32_t       pr_ip;                /* (c) ip addr host */  
    void            *pr_linux;              /* (p) linux abi */  
    int             pr_securelevel;        /* (p) securelevel */  
    struct task    *pr_task;               /* (d) destroy task */  
    struct mtx     pr_mtx;  
    void            **pr_slots;              /* (p) additional data */  
};
```

然后， 系用jail(2)一个prison体分配一内存，并在jail和prison体之制数据。

```

/usr/src/sys/kern/kern_jail.c:
MALLOC(pr, struct prison *, sizeof(*pr), M_PRISON, M_WAITOK | M_ZERO);
...
error = copyinstr(j.path, pr->pr_path, sizeof(pr->pr_path), 0);
if (error)
    goto e_killmtx;
...
error = copyinstr(j.hostname, pr->pr_host, sizeof(pr->pr_host), 0);
if (error)
    goto e_dropvnref;
pr->pr_ip = j.ip_number;

```

下面，我将介绍另一个重要的系统调用[jail\\_attach\(2\)](#)，它实现了将进程禁的功能。

```

/usr/src/sys/kern/kern_jail.c
/*
 * struct jail_attach_args {
 *     int jid;
 * };
 */
int
jail_attach(struct thread *td, struct jail_attach_args *uap)

```

这个系统调用做出一些可以用于区分被禁和未被禁的进程的改动。

要理解[jail\\_attach\(2\)](#)做了什么，我们首先要理解一些背景信息。

在FreeBSD中，一个内核可运行的进程是通过其thread实体来实现的，同时，进程都由它自己的proc实体描述。可以在/usr/include/sys/proc.h中找到thread和proc实体的定义。例如，在任何系统调用中，参数td上是个指向可用进程的thread实体的指针，正如前面所讲的那样。td所指向的thread实体中的td\_proc成员是一个指针，该指针指向td所表示的进程所属进程的proc实体。proc实体proc包含的成员可以描述所有者的身份(p\_ucred)，进程限制(p\_limit)，等等。在由proc实体的p\_ucred成员所指向的ucred实体的定义中，还有一个指向prison实体的指针(cr\_prison)。

```

/usr/include/sys/proc.h:
struct thread {
    ...
    struct proc *td_proc;
    ...
};

struct proc {
    ...
    struct ucred *p_ucred;
    ...
};

/usr/include/sys/ucred.h
struct ucred {
    ...
    struct prison *cr_prison;
    ...
};

```

在kern\_jail.c中，函数**jail()**以固定的**jid**用函数**jail\_attach()**。随后**jail\_attach()**用函数**change\_root()**以改用程序的根目录。接下来，**jail\_attach()**创建一个新的**ucred**对象，并在成功地将**prison**对象接到一个**ucred**对象后，将这个**ucred**对象连接到该程序上。从此日起，这个程序就会被限制被禁的。当我以新建的一个**ucred**对象参数用内核路径**jailed()**，它将返回1来表明这个用户身是和一个jail相同的。在jail中叉分出来的所有程序的公共祖先程序就是那个执行了**jail(2)**的程序，因为正是它用了**jail(2)**系统调用。当一个程序通过**execve(2)**而被执行，它将从其父程序的**ucred**对象继承被禁的属性，因而它也会有一个被禁的**ucred**对象。

```

/usr/src/sys/kern/kern_jail.c
int
jail(struct thread *td, struct jail_args *uap)
{
...
    struct jail_attach_args jaa;
...
    error = jail_attach(td, jaa);
    if (error)
        goto e_droppref;
...
}

int
jail_attach(struct thread *td, struct jail_attach_args *uap)
{
    struct proc *p;
    struct ucred *newcred, *oldcred;
    struct prison *pr;
...
    p = td->td_proc;
...
    pr = prison_find(uap->jid);
...
    change_root(pr->pr_root, td);
...
    newcred->cr_prison = pr;
    p->p_ucred = newcred;
...
}

```

当一个进程被从其父进程叉分来的時候，系统用fork(2)将用crhold()来共享其身凭。如果  
，很自然的就保持了子进程的身凭于其父进程一致，所以子进程也是被禁的。

```

/usr/src/sys/kern/kern_fork.c:
p2->p_ucred = crhold(td->td_ucred);
...
td2->td_ucred = crhold(p2->p_ucred);

```

## 4.2. 系统被囚禁程序的限制

在整个内核中，有一系列被囚禁程序的约束措施。  
突破这些约束，相的函数将出返回。例如：

```

if (jailed(td->td_ucred))
    return EPERM;

```

### 4.2.1. SysV 程序通信(IPC)

System V 程序通信 (IPC) 是通过消息的。一个程序都可以向其它程序发送消息，告诉对方做什么。管理消息的函数是：`msgctl(3)`、`msgget(3)`、`msgsnd(3)` 和 `msgrecv(3)`。前面已提到，一些 `sysctl` 可以影响 jail 的行为，其中有一个是 `security.jail.sysvipc_allowed`。在大多数系统上，一个 `sysctl` 会成为 0。如果将它设为 1，程序会完全失去 jail 的意图：因为那时在 jail 中特权限就可以影响被禁的境外的程序了。消息与信号的区别是：消息由一个信号号组成。

/usr/src/sys/kern/sysv\_msg.c:

- `msgget(key, msgflg)`: `msgget`返回(也可能创建)一个消息描述符，以指派一个在其它函数中使用的消息队列。
- `msgctl(msqid, cmd, buf)`: 通常一个函数，一个程序可以拥有一个消息描述符的状态。
- `msgsnd(msqid, msgp, msgs, msgflg)`: `msgsnd`向一个程序送一条消息。
- `msgrecv(msqid, msgp, msgs, msgtyp, msgflg)`: 程序用一个函数接收消息。

在这些函数的实现代码中，都有一个条件判断：

```
/usr/src/sys/kern/sysv_msg.c:  
if (!jail_sysvipc_allowed jailed(td->td_ucred))  
    return (ENOSYS);
```

信号量系统调用使得程序可以通过一系列原子操作同步。信号量系统调用提供了另一种途径。然而，程序将正在被使用的信号量放入等待状态，一直休眠到资源被释放。在 jail 中如下的信号量系统调用将会失效：`semget(2)`, `semctl(2)` 和 `semop(2)`。

/usr/src/sys/kern/sysv\_sem.c:

- `semctl(semid, num, cmd, ...)`: `semctl`在信号量队列中用 `semid` 的信号量执行 `cmd` 指定的命令。
  - `semget(key, nsems, flag)`: `semget`建立一个基于 `key` 的信号量数组。
- 参数 `key` 和 `flag` 与它们在 `msgget()` 的意图相同。
- `semop(semid, array, nops)`: `semop` 对 `semid` 的信号量完成一个由 `array` 所指定的操作。

System V IPC 使程序可以共享内存。程序之间可以通过它们虚拟地址空间的共享部分以及相写操作直接通信。某些系统调用在被禁的境内将失效：`shmdt(2)`、`shmat(2)`、`shmctl(2)` 和 `shmget(2)`

/usr/src/sys/kern/sysv\_shm.c:

- `shmctl(shmid, cmd, buf)`: `shmctl` 对共享内存区域做各自的控制。
- `shmget(key, size, flag)`: `shmget` 建立一个 `size` 字节的共享内存区域。
- `shmat(shmid, addr, flag)`: `shmat` 将 `shmid` 的共享内存区域指派到程序的地址空间里。
- `shmdt(addr)`: `shmdt` 取消共享内存区域的地址指派。

## 4.2.2. 套接字

Jail以一种特殊的方式处理socket(2)系统调用和相关的低级套接字函数。为了决定一个套接字是否允许被创建，它首先通过sysctl security.jail.socket\_unixiproute\_only来检查是否被置为1。如果被置为1，套接字建立将只能指定某些族：PF\_LOCAL, PF\_INET, PF\_ROUTE。否则，socket(2)将会返回出口。

```
/usr/src/sys/kern/uipc_socket.c:  
int  
socreate(int dom, struct socket **aso, int type, int proto,  
        struct ucred *cred, struct thread *td)  
{  
    struct protosw *prp;  
    ...  
    if (jailed(cred) && jail_socket_unixiproute_only  
        prp->pr_domain->dom_family != PF_LOCAL  
        prp->pr_domain->dom_family != PF_INET  
        prp->pr_domain->dom_family != PF_ROUTE) {  
        return (EPROTONOSUPPORT);  
    }  
    ...  
}
```

## 4.2.3. Berkeley包过滤器

Berkeley包过滤器提供了一个与内核的，直接通向数据路的低级接口。在BPF是否可以在禁的环境中被使用是通过devfs(8)来控制的。

## 4.2.4. 网络

网络TCP, UDP, IP和ICMP很常见。IP和ICMP位于同一层次：第二层，网层。当参数nam被置为，有一些限制措施会防止被囚禁的程序绑定到一些网口接口上。nam是一个指向sockaddr结构体的指针，描述可以绑定服务的地址。一个更确切的定义：sockaddr"是一个模板，包含了地址的字符和地址的长度"。

在函数in\_pcbbind\_setup()中sin是一个指向sockaddr\_in结构体的指针，该结构体包含了套接字可以绑定的端口、地址、长度、族。这就禁止了在jail中的进程指定不属于该进程所存在于的jail的IP地址。

```

/usr/src/sys/kern/netinet/in_pcbs.c:
int
in_pcbbind_setup(struct inpcb *inp, struct sockaddr *nam, in_addr_t *laddrp,
    u_short *lportp, struct ucred *cred)
{
    ...
    struct sockaddr_in *sin;
    ...
    if (nam) {
        sin = (struct sockaddr_in *)nam;
        ...
        if (sin->sin_addr.s_addr != INADDR_ANY)
            if (prison_ip(cred, 0, sin->sin_addr.s_addr))
                return(EINVAL);
        ...
        if (lport) {
            ...
            if (prison_ip(cred, 0, sin->sin_addr.s_addr))
                return (EADDRNOTAVAIL);
            ...
        }
    }
    if (lport == 0) {
        ...
        if (laddr.s_addr != INADDR_ANY)
            if (prison_ip(cred, 0, laddr.s_addr))
                return (EINVAL);
        ...
    }
    ...
    if (prison_ip(cred, 0, laddr.s_addr))
        return (EINVAL);
}

```

想知道函数`prison_ip()`做什么。`prison_ip()`有三个参数，一个指向身份凭证的指针（用`cred`表示），一些标志和一个IP地址。当该IP地址不属于某个jail时，返回1；否则返回0。正如从代码中看到的，如果，那个IP地址不属于某个jail，就不再允许向该网络地址绑定。

```
/usr/src/sys/kern/kern_jail.c:  
int  
prison_ip(struct ucred *cred, int flag, u_int32_t *ip)  
{  
    u_int32_t tmp;  
  
    if (!jailed(cred))  
        return (0);  
    if (flag)  
        tmp = *ip;  
    else  
        tmp = ntohl(*ip);  
    if (tmp == INADDR_ANY) {  
        if (flag)  
            *ip = cred->cr_prison->pr_ip;  
        else  
            *ip = htonl(cred->cr_prison->pr_ip);  
        return (0);  
    }  
    if (tmp == INADDR_LOOPBACK) {  
        if (flag)  
            *ip = cred->cr_prison->pr_ip;  
        else  
            *ip = htonl(cred->cr_prison->pr_ip);  
        return (0);  
    }  
    if (cred->cr_prison->pr_ip != tmp)  
        return (1);  
    return (0);  
}
```

#### 4.2.5. 文件系统

如果完全大于0，即便是jail里面的root，也不允许在Jail中取消或更改文件标志，如“不可修改”、“只可添加”、“不可删除”标志。

```
/usr/src/sys/ufs/ufs/ufs_vnops.c:  
static int  
ufs_setattr(ap)  
{  
    ...  
    if (!priv_check_cred(cred, PRIV_VFS_SYSFLAGS, 0)) {  
        if ((ip-i_flags  
            (SF_NOUNLINK | SF_IMMUTABLE | SF_APPEND)) {  
            error = securelevel_gt(cred, 0);  
            if (error)  
                return (error);  
        }  
        ...  
    }  
}/usr/src/sys/kern/kern_priv.c  
int  
priv_check_cred(struct ucred *cred, int priv, int flags)  
{  
    ...  
    error = prison_priv_check(cred, priv);  
    if (error)  
        return (error);  
    ...  
}  
/usr/src/sys/kern/kern_jail.c  
int  
prison_priv_check(struct ucred *cred, int priv)  
{  
    ...  
    switch (priv) {  
        ...  
        case PRIV_VFS_SYSFLAGS:  
            if (jail_chflags_allowed)  
                return (0);  
            else  
                return (EPERM);  
        ...  
    }  
    ...  
}
```

# Chapter 5. SYSINIT框架

SYSINIT是一个通用的用排序与分行机制的框架。

SYSINIT使得FreeBSD的内核各子系可以在内核或模块上加接被重整、添加、删除、替换，内核和模块加就不必去修改一个静态的有序初始化安排表甚至重新内核。个体系也使得内核模块（在称为KLD）可以与内核不同、接、在引入系统加载，甚至在系统运行加载。这些操作是通过“内核接器”（kernel linker）和“接器集合”（linker set）完成的。

FreeBSD目前使用它来执行内核的初始化。

## 5.1. 接器集合

### 接器集合(Linker Set)

一种方法。这种方法将整个程序源文件中静态申明的数据收集到一个可近址的数据元中。

## 5.2. SYSINIT操作

SYSINIT要依赖接器取遍布整个程序源代码申明的静态数据，并把它变成一个彼此相的数据。这种方法被称为“接器集合”（linker set）。SYSINIT使用一个接器集合以多个数据集合，包含多个数据条目的顺序、函数、一个会被提交函数的数据指。

SYSINIT按照先全局函数排序以便执行。第一先的是子系统的，其次是SYSINIT分行子系统的函数的全局排序，定在sys/kernel.h中的枚举 `sysinit_sub_id` 内。第二先在子系统中的元素的排序，定在sys/kernel.h中的枚举 `sysinit_elem_order` 内。

有些时刻需要使用SYSINIT：系统或内核模块加载，系统析或内核模块卸载。内核子系统通常在系统使用SYSINIT的定义以初始化数据。例如，程度子系统使用一个SYSINIT 定义来初始化行列数据。程度程序避免直接使用 `SYSINIT()`，对于上的物理真部分使用 `DRIVER_MODULE()` 用的函数先部分的存在，如果存在，再执行的初始化。一进程中，会做一些的事情，然后用 `SYSINIT()` 本身。对于非一部分的虚，改用 `DEV_MODULE()`。

## 5.3. 使用SYSINIT

### 5.3.1. 接口

#### 5.3.1.1. 文件

```
sys/kernel.h
```

#### 5.3.1.2. 宏

```
SYSINIT(uniquifier, subsystem, order, func, ident)  
SYSUNINIT(uniquifier, subsystem, order, func, ident)
```

### 5.3.2. 宏

宏SYSINIT()在SYSINIT数据集合中 建立一个SYSINIT数据，以便SYSINIT在系□□或模□加□排序 并□行其中的函数。SYSINIT()有一个参数uniquifier， SYSINIT用它来□□数据，随后是子系□序号、子系□元素□序号、待□用函数、□□函数的数据。所有的函数必□有一个恒量指□参数。

例 1. SYSINIT()的例子

```
#include sys/kernel.h

void foo_null(void *unused)
{
    foo_doo();
}

SYSINIT(foo, SI_SUB_FOO, SI_ORDER_FOO, foo_null, NULL);

struct foo foo_voodoo = {
    FOO_VOODOO;
}

void foo_arg(void *vdata)
{
    struct foo *foo = (struct foo *)vdata;
    foo_data(foo);
}

SYSINIT(bar, SI_SUB_FOO, SI_ORDER_FOO, foo_arg, foo_voodoo);
```

注意，SI\_SUB\_FOO和SI\_ORDER\_FOO

□当分□在上面提到的枚□sysinit\_sub\_id和sysinit\_elem\_order之中。既可以使用已有的枚□，也可以将自己的枚□添加到□□个枚□的定□之中。□可以使用数学表□式微□SYSINIT的□行□序。 以下的例子示例了一个需要□好要在内核参数□整的SYSINIT之前□行的SYSINIT。

例 2. □整SYSINIT()□序的例子

```
static void
mptable_register(void *dummy __unused)
{
    apic_register_enumerator(mptable_enumerator);
}

SYSINIT(mptable_register, SI_SUB_TUNABLES - 1, SI_ORDER_FIRST,
        mptable_register, NULL);
```

### 5.3.3. 析口

宏SYSUNINIT()的行口与SYSINIT()的相当，只是它将数据口填加至SYSINIT的数据集合。

例3. SYSINIT()的例子

```
#include sys/kernel.h

void foo_cleanup(void *unused)
{
    foo_kill();
}
SYSINIT(foobar, SI_SUB_FOO, SI_ORDER_FOO, foo_cleanup, NULL);

struct foo_stack foo_stack = {
    FOO_STACK_VOODOO;
}

void foo_flush(void *vdata)
{
}
SYSINIT(barfoo, SI_SUB_FOO, SI_ORDER_FOO, foo_flush, foo_stack);
```

# Chapter 6. TrustedBSD MAC 框架

## 6.1. MAC 文档声明

本文档是作为 DARPA CHATS 研究的一部分，由供于 Security Research Division of Network Associates 公司Safeport Network Services and Network Associates Laboratories 的Chris Costello依据 DARPA/SPAWAR 合同 N66001-01-C-8035 ("CBOSS")，由 FreeBSD 团队写的。

Redistribution and use in source (SGML DocBook) and 'compiled' forms (SGML, HTML, PDF, PostScript, RTF and so forth) with or without modification, are permitted provided that the following conditions are met:

1. Redistributions of source code (SGML DocBook) must retain the above copyright notice, this list of conditions and the following disclaimer as the first lines of this file unmodified.
2. Redistributions in compiled form (transformed to other DTDs, converted to PDF, PostScript, RTF and other formats) must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

THIS DOCUMENTATION IS PROVIDED BY THE NETWORKS ASSOCIATES TECHNOLOGY, INC "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL NETWORKS ASSOCIATES TECHNOLOGY, INC BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS DOCUMENTATION, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.



本文档可作为非官方中文翻译供参考，  
不作为判定任何  
任的依据。如与英文原文有出入，以英文原文为准。

在满足下列可条件的前提下，允再分或以源代 (SGML DocBook) 或 "XXX" (SGML, HTML, PDF, PostScript, RTF 等) 的修改或未修改的形式：

1. 再分源代 (SGML DocBook) 必不加修改的保留上述版告示、本条件清和下述作  
文件的最先若干行。
2. 再分XXX的形式 (XXX其它DTD、PDF、PostScript、RTF 或其它形式)，必将上述版告示、本条件清  
和下述XXX制到与分品一同提供的文件，以及其它材料中。

本文由 NETWORKS ASSOCIATES TECHNOLOGY, INC "按**现状**条件提供，并在此明示不提供任何明示或暗示的保障，包括但不限于**商****品性**、**适****用性**、**特****定目的的****用性**的**暗示保障**。任何情况下，NETWORKS ASSOCIATES TECHNOLOGY, INC 均不**任何直接、间接、偶然、特殊、**口**性的**，**或必然的****失**（包括但不限于**替代商品或服务的采****用、使用、数据或利益的****失或****中断**）**或**，**无论**是**如何****致的并以任何有****任****的**，**无论**是否是在本文**使用以外以任何方式****生的契****格****任**或**民事侵****行**（包括疏忽或其它）中的，即使已被告知**生****失的可能性**。



## 6.2. **口****解****析**

FreeBSD 以一个内核安全**扩展**性框架(TrustedBSD MAC 框架)的方式，**若干****制****控****策****略**（也称"**集****式****控制策略**"）**提****供****性****支****持**。MAC 框架是一个**入****式****控****框****架**，**允****新****安****全****策****略****更****便****融****入****内****核**：**安****全****策****略****可****静****入****内****核**，**也****可****在****引****加****加**，**甚****至****在****行****加****加**。**框****架****所****提****供****的****准****化****接****口**，**使****得****行****在其****上****的****安****全****策****略****模****能****系****象****的****安****全****属****性****行****如****等****一****系****列****操****作**。MAC 框架的存在，**化****了****些****操****作****在****策****略****模****中****的****，****从****而****著****降****低****了****新****安****全****策****略****模****的****度**。

本章将介**MAC** 策略框架，**者****提****供****一个****示****例****性****的****MAC** 策略模**文****。**

## 6.3. **概****述**

TrustedBSD MAC 框架提供的机制，**允****在其****上****行****的****内****核****模****在****内****核****或****者****行****，****内****核****的****控****模****行****展**。**新****的****系****安****全****策****略****作****一****内****核****模****并****被****接****到****内****核****中**；**如****果****系****中****同****存****多****个****安****全****策****略****模****，****它****的****决****策****果****将****以****某****定****方****式****合**。**了****化****新****安****全****策****略****的****，**MAC 向上提供了大量用于**控****的****基****施**，**特****是**，**的****或****久****、****策****略****无****象****安****全****的****支****持**。**支持****目****前****仍****是****性****的**。

本章所提供的信息不**将****使****在** MAC 使能**境****下****工****的****潜****在****用****受益**，**也****可****以****需****要****解****MAC** 框架是如何支持**内****核****控****行****展****的****策****略****模****人****所****用****。**

## 6.4. 安全策略背景知**口**

**制****控****策****略**（**称** MAC），**是****指****由****操****作****系****制****施****的一****用****的****控****策****略**。**在****某****些****情****况****下**，**制****控****策****略****可****能****会****与****自****主****控****（****称** DAC）**所****提****供****保****措****施****冲****突**，**后****者****是****用****向****非****管****理****用****数****据****采****取****保****措****施****提****供****支****持****的**。在**的** UNIX 系**中**，DAC 保**措****施****包****括****文****件****模****和****控****列****表**；而 MAC 提供**程****控****和****火****等**。操作系**者****和****安****全****机****制****研****究****人****多****典****的** MAC 安全策略作了形式化的表述，比如，**多****安****全** (MLS) 机密性策略，Biba 完整性策略，基于角色的**控****策****略** (RBAC)，域和型裁决策策(DTE),以及型裁决策策(TE)。安全策略的形式化表述被称**安****全****模****型**。**个****模****型****根****据****一****系****列****条****件****做****出****安****全****相****的****决****策**，**些****条****件****包****括****用****的****身****、****角****和****安****全****信****状****，****以****象****的****安****全****（****用****来****代****表****象****数****据****的****机****密****性****/****完****整****性****）**。

TrustedBSD MAC 框架所提供的**策****略****模****的****支****持**，**不****可****以****用****来****上述****所****有****策****略**，**能****用****于****其****他****利****用****已****有****安****全****属****性****（****如****用****ID**、**文件****属****性****等**）**决****策****的****安****全****化****策****略**。**此****外****，****因****具****体****策****略****模****在****授****方****面****所****有****的****高****度****活****性和****自****主****性**，**所****以** MAC 框架同**可****以****用****来****完全****自****主****的****安****全****策****略**。

## 6.5. MAC 框架的内核体系**口**

TrustedBSD MAC 框架**大****部****分****的****控****模****提****供****基****施**，**允****它****以****内****核****模****的****形****式****活****地****展****系****中****施****的****安****全****策****略**。**如****果****系****中****同****加****了****多****个****策****略**，MAC 框架将**将****各****个****策****略****的****授****果****以****一****（****某**

程度上)有意的方式结合,形成最后的决策。

### 6.5.1. 内核元素

MAC 框架由下列内核元素组成：

- 框架管理接口
- 并与同原
- 策略注册
- 内核象的扩展性安全
- 策略入口函数的结合操作
- 管理原
- 由内核服务用的入口函数 API
- 策略模块的入口函数 API
- 入口函数的 (包括策略生命周期管理、管理和控制三部分)
- 管理策略无的系用
- 用的 `mac_syscall()` 系用
- 以 MAC 的策略加模形式的各安全策略

### 6.5.2. 框架管理接口

TrustedBSD MAC 框架行直接管理的方式有三种:通过 sysctl 子系、通过 loader 配置,或者使用系用。

多数情况下,与同一个内核内部量相同的 sysctl 值和 loader 参数的名字是相同的,通过它,可以控制保措施的实施,比如,某个策略在各个内核子系中的实施与否等等。此外,如果在内核中支持 MAC ,内核将若干数器以跟踪的分配使用情况。通常不建议在用户空间下通过在不同子系上设置不同的量或参数来实施控制,因为方法将会作用于系中所有的活策略。如果希望具体策略实施管理而不影响其他活策略,当使用策略的控制,因为方法的控制粒度更细,并能更好地保持策略模块的功能一致性。

与其他内核模块一样,系管理可以通过系的模块管理系用和其他系接口,包括 boot loader 值,策略模块行加与卸操作;策略模块可以在加,置加志,来指示系其加、卸操作行相控制,比如阻止非期望的卸操作。

### 6.5.3. 策略表的并同

在行,系中活的策略集合可能生化,然而策略入口函数的使用操作并不是原子性的,因此,当某一个入口函数正被使用,系需要提供外的同机制来阻止策略模块的加与卸,以保当前活的策略集合不会在此过程中生改。通过使用"框架忙"数器,就可以做到一点:一旦某个入口函数被用,数器的被加1;而当一个入口函数用完,数器的被少1。数器的,如果其正,框架将阻止策略表的修改操作,求操作的过程将被迫入睡眠,直到数器的重新少到0为止。数器本身由一个互斥保,同合一个条件量(用于醒等待策略表行修改操作的睡眠程序)。采用同一模型的一个副作用是,在同一个策略模块内部,允嵌套地用框架,不情况其很少出口。

了少由于采用数器引入的外,后者采用了各量化措施。其中包括,当策略表空或者其中含有静态表

(那些只能在系统运行之前加载而且不能卸载的策略)。框架不负责器运行操作，其原因是0，从而将此0的同0到0。一个0端的0法是，使用一个0000来禁止在0行0加0的策略0表0行修改，此0不再需要0策略0表的使用0行同0保0。

因0 MAC 框架不允许在某些入口函数之内阻塞，所以不能使用普通的睡眠0。故而，加载或卸载操作可能会0等待框架空0而被阻塞相当0的一段0。

#### 6.5.4. 00同0

MAC 框架必须为其0000的安全属性00的存000提供同0保0。下列00情形，可能导致0安全属性00的不一致0：第一，作0安全属性00的持有者，内核0象本身可能同时被多个0程00；第二，MAC 框架代码是可重入的，即允0多个0程同时在框架内0行。通常，MAC 框架使用内核0象数据上已有的内核同0机制来保0其上附加的MAC 安全00。例如，套接字上的 MAC 00由已有的套接字互斥0保0。0似的，0于安全00的并000的0程与0其所在0象0行的并000在00上是一0的，例如，信任状安全00，将保持与0数据00中其他内容一致的"写00制"的更新0程。MAC 框架在引用一个内核0象0，将首先00000象上的00需要用到的00行断言。策略模0的0写者必0了解0些同000，因0它0可能会限制0安全00所能0行的000型。0个例子，如果通0入口函数00策略模0的是0某个信任状的只0引用，那0在策略内部，只能000000的00状0。

#### 6.5.5. 策略0的同0与并0

FreeBSD 内核是一个可0占式的内核，因此，作0内核一部分的策略模0也必0是可重入的，也就是0，在00策略模0必0假0多个内核0程可以同时通0不同的入口函数0入0模0。如果策略模0使用可被修改的内核状0，那00需要在策略内部使用恰当的同0原0，0保在策略内部的多个0程不会因此0察到不一致的内核状0，从而避免由此0生的策略0操作。0此，策略可以使用 FreeBSD 0有的同0原0，包括互斥0、睡眠0、条件0量和0数信号量。0些同0原0的使用必0慎重，需要特0注意0点：第一，保持0有的内核上0次序；第二，在非睡眠的入口函数之内不要使用互斥0和0醒操作。

0避免0反内核上0次序或造成00上0，策略模0在0用其他内核子系0之前，通常要0放所有在策略内部申0的0。00做的0果是，在全局上0次序形成的拓朴00中，策略内部的00是作0叶子0点，从而保0了0些0的使用不会0致由于上0次序混乱造成的死0。

#### 6.5.6. 策略注册

0了0当前使用的策略模0集合，MAC 框架000个0表：一个静00表和一个000表。0个0表的数据00和操作基本相同，只是000表00外使用了一个"引用0数"以同0其的00操作。当包含 MAC 框架策略的内核模0被加00，0策略模0会通0 **SYSINIT** 0用一个注册函数；相00的，0当一个策略模0被卸0，**SYSINIT** 也会0用一个注0函数。只有当遇到下列情况之一0，注册0程才会失0：一个策略模0被加0多次，或者系00源不足不能0足注册0程的需要（例如，策略模0需要0内核0象添加00而可用0源不足），或者其他0策略加0前提条件不0足（有些策略要求只能在系0引0之前加0）。0似的，如果一个策略被000不可卸0的，0其0用注00程将会失0。

#### 6.5.7. 入口函数

内核服0与 MAC 框架之00行交互有00途径：一是，内核服0用一系列 API 通知 MAC 框架安全事件的0生；二是，内核服0向 MAC 框架提供一个指向安全0象的策略无0安全00数据00的指0。00指0由 MAC 框架0由00管理入口函数0行00，并且，只要0管理相00象的内核子系00作修改，就可以允0 MAC 框架向策略模0提供00服0。例如，在0程、0程信任状、套接字、管道、Mbuf、网0接口、IP 重00列和其他各0安全相0的数据00中均0加了指向安全00的指0。0外，当需要做出重要的安全决策0，内核服0也会0用 MAC 框架，以便各个策略模0根据其自己的0准（可以使用存0在安全00中的数据）完善0些决策。0大多数安全相

的决策是式的控制；也有少数及更加一般的决策函数，比如，套接字的数据包匹配和程序行时刻的。

### 6.5.8. 策略组合

如果内核中同时加入了多个策略模块，一些策略的决策结果将由框架使用一个合成算子来执行组合，得出最终的结果。

目前，算子是硬的，并且只有当所有的活动策略均要求表示同意才会返回成功。由于各个策略返回的出口条件可能并不相同（成功、被拒、求象不存在等等），需要使用一个子先从各个策略返回的出口条件集合中取出一个作为最终返回结果。一般情况下，与“被拒”相比，将更倾向于“求象不存在”。尽管不能从理论上保证合成结果的有效性与安全性，但结果表明，对于多用的策略集合来说，事情的如此。例如，可信系统常常采用类似的方法用多个安全策略执行组合。

### 6.5.9. 支持

与多对象添加安全相关的控制扩展，MAC 框架各可用对象提供了一个用于管理策略无的系统。常用的类型有，partition符、机密性、完整性、区（非等）、域、角色和型。“策略无”的意思是指，它的方法与使用它的具体策略模块无关，而同一个策略模块能完全独立地定义和使用与对象相同的元数据的。

用应用程序提供统一格式的基于字符串的，由使用它的策略模块解析其内在含义并决定其外在表示。如果需要，应用程序可以使用多重元素。

内存中的例被存放在由 slab 分配的 `struct label` 数据中。这是一个固定宽度的数，个元素是由一个 `void *` 指和一个 `long` 成的组合。申存的策略模块向 MAC 注册，将被分配一个“slot”，作框架分配其使用的策略元素在整个存中的位置索引。而所分配的存空的完全由策略模块来决定：MAC 框架向策略模块提供了一系列入口函数用于内核对象生命周期的各事件进行控制，包括，对象的初始化、的创建和对象的注销。使用这些接口，可以如数等存模型。MAC 框架是入口函数入一个指向对象的指和一个指向对象的指，因此，策略模块能直接而无需知悉对象的内部。唯一的例外是进程信任状态，指向其的指必须由策略模块解析。今后的 MAC 框架可能会对此行改动。

初始化入口函数通常有一个睡眠标志位，用来表明一个初始化操作是否允许中途睡眠等待；如果不允，可能会失败返回，并要求撤此次分配操作（乃至对象分配操作）。例如，如果在网络管理中断因不允睡眠或者用户持有一个互斥，就可能出现情况。考虑到在管理中的网数据包（Mbufs）上性能的损失太大，策略必须就自己 Mbuf 行的要求向 MAC 框架做出特别声明。附加到系中而又使用的策略必须管理未被其初始化函数管理的对象作好准备，有些对象在策略附加之前就已存在，故而无法在初始化用策略的相函数行管理。MAC 框架向策略保证，没有被初始化的 slot 的必须或者 NULL，策略可以借此附加到未初始化的。需要注意的是，因为 Mbuf 的存分配是有条件的，因此需要使用其的附加策略可能需要管理 Mbuf 中的 NULL 的指。

于文件系统的，MAC 框架在文件的扩展属性中为其分配永久存储。只要可能，扩展属性的原子化的事操作就被用于保 vnode 上安全的组合更新操作的一致性——目前，特性只被 UFS2 文件系统支持。除了粒度的文件系统对象（即一个文件系统一个），策略写者可能使用一个（或者若干）扩展属性。为了提高性能，vnode 数据中有一个 `(v_label)` 字段，用作磁的缓冲；vnode 例化后，策略可以将其装入缓冲，并在需要时对其进行更新。如此，不必在每次行控制，均无条件地磁上的扩展属性。



目前，如果一个使用的策略允被卸，卸模块之后，其状态 slot 尚无法被系统回收重用，由此导致了 MAC 框架策略卸载操作数目上的严格限制。

## 6.5.10. 相互通用

MAC 框架向应用程序提供了一系列公用：其中大多数用于向执行和修改策略无操作的公用 API 提供支持。

一些公用管理系用，接受一个描述，`struct mac`，作为参数。这个的主体是一个数，其中个元素包含了一个公用的 MAC 形式。个元素又由部分组成：一个字符串名字，和其的。一个策略可以向系声明一个特定的元素名字，一来，如果需要，就可以将若干个相互独立的元素作为一个整体执行。策略模由入口函数，在内核和用户提供的之间工作，提供了元素上的高度活性。公用管理系用通常有公用的函数包装，些包装函数可以提供内存分配和清理功能，从而简化了公用程序的管理工作。

目前的 FreeBSD 内核提供了下列 MAC 相的系用：

- `mac_get_proc()` 用于当前程的安全。
- `mac_set_proc()` 用于求改当前程的安全。
- `mac_get_fd()` 用于由文件描述符所引用的对象（文件、套接字、管道文件等等）的安全。
- `mac_get_file()` 用于由文件系路径所描述的对象的安全。
- `mac_set_fd()` 用于求改由文件描述符所引用的对象（文件、套接字、管道文件等等）的安全。
- `mac_set_file()` 用于求改由文件系路径所描述的对象的安全。
- `mac_syscall()` 通用公用，策略模能在不修改系用表的前提下创建新的系用；公用参数包括：目标策略名字、操作号和将被策略内部使用的参数。
- `mac_get_pid()` 用于由程号指定的一个程的安全。
- `mac_get_link()` 与 `mac_get_file()` 功能相同，只是当路径参数的最后一符号接，前者将返回符号接的安全，而后者将返回其所指文件的安全。
- `mac_set_link()` 与 `mac_set_file()` 功能相同，只是当路径参数的最后一符号接，前者将置符号接的安全，而后者将置其所指文件的安全。
- `mac_execve()` 与 `execve()` 功能相似，只是前者可以在始执行一个新程序，根据入的求参数，置执行程的安全。由于执行一个新程序而致的程安全的改，被称“”。
- `mac_get_peer()`，通过一个套接字自，用于一个程套接字等体的安全。

除了上述系用之外，也可以通过 `SIOCSIGMAC` 和 `SIOCSIFMAC` 网口接口的 ioctl 系用用来和置网口接口的安全。

## 6.6. MAC策略模体系

安全策略可以直接入内核，也可以成独立的内核模，在系引或者执行使用模加命令加。策略模通一先定好的入口函数与系交互。通过它，策略模能掌握某些系事件的生，并且在必要的时候影系的控制决策。一个策略模包含下列部分：

- 可：策略配置参数
- 策略和参数的集中
- 可：策略生命周期事件的，比如，策略的初始化和
- 可：所内核对象的安全行初始化、和的支持

- 可以：对象的使用、执行控制以及修改对象安全性的支持
- 策略相关的控制入口函数的API
- 策略标志、模块入口函数和策略特性的声明

### 6.6.1. 策略注入

策略模块可以使用 **MAC\_POLICY\_SET()** 宏来声明。该宏完成以下工作：为策略命名（向系统声明策略提供的名字）；提交策略定义的 MAC 入口函数向量的地址；按照策略的要求配置策略的加权位，保证 MAC 框架将以策略所期望的方式执行操作；此外，可能要求框架将策略分配到状态 slot 中。

```
static struct mac_policy_ops mac_policy_ops =
{
    .mpo_destroy = mac_policy_destroy,
    .mpo_init = mac_policy_init,
    .mpo_init_bpfdesc_label = mac_policy_init_bpfdesc_label,
    .mpo_init_cred_label = mac_policy_init_label,
/* ... */
    .mpo_check vnode_setutimes = mac_policy_check vnode_setutimes,
    .mpo_check vnode_stat = mac_policy_check vnode_stat,
    .mpo_check vnode_write = mac_policy_check vnode_write,
};
```

如上所示，MAC 策略入口函数向量，**macpolicyops**，将策略模块中定义的功能函数挂接到特定的入口函数地址上。  
在“[入口函数参考](#)”小节中，将提供可用入口函数功能描述和原型的完整列表。与模块注册相关的入口函数有两个：**.mpo\_destroy** 和 **.mpo\_init**。当某个策略向模块框架注册操作成功后，**.mpo\_init** 将被调用，此后其他的入口函数才能被使用。  
通过特定的分配和初始化操作，比如特殊数据或的初始化。卸载一个策略模块，将调用 **.mpo\_destroy** 用来释放策略分配的内存空间或注销其申请的。目前，为了防止其他入口函数被同调用，调用上述两个入口函数的编程必须持有 MAC 策略表的互斥锁：限制将被放，但与此同时，将要求策略必须谨慎使用内核原语，以避免由于上次序或睡眠造成死锁。

之所以向策略声明提供模块名字域，是为了能唯一地标识模块，以便解析模块依赖关系。使用恰当的字符串作名字。  
在策略加载和卸载，策略的完整字符串名字将由内核日志显示。此外，当向用进程报告状态信息时也会包含字符串。

### 6.6.2. 策略标志

在声明提供标志参数域的机制，允许策略模块在加载模块时，就自身特性向 MAC 框架提供声明。目前，已定义的标志有三个：

#### **MPC\_LOADTIME\_FLAG\_UNLOADOK**

表示策略模块可以被卸载。如果未提供标志，表示策略模块拒绝被卸载。那些使用安全模式的策略，而又不能在执行状态下运行的模块可能会被置为标志。

#### **MPC\_LOADTIME\_FLAG\_NOTLATE**

表示策略模块必须在系统引导过程中执行加载和初始化。如果标志被设置，那么在系统引导之后注册模块的需求将被忽略。

MAC 框架所拒。那些需要大量安全的系统象进行安全初始化工作，而又不能理解含有未被正确初始化安全的象的策略模可能会置志。

## MPC\_LOADTIME\_FLAG\_LABELMBUFS

表示策略模要求 Mbuf 指定安全，并且存其所需的内存空是提前分配好的。缺省情况下，MAC 框架并不会 Mbuf 分配存，除非系中注册的策略模中至少有一个置了志。做法在没有策略需要 Mbuf 做，著地提升了系网性能。外，在某些特殊境下，可以通过内核，`MAC_ALWAYS_LABEL_MBUF`，抵制 MAC 框架 Mbuf 的安全分配存，而不上述志如何置。

 那些使用了 `MPC_LOADTIME_FLAG_LABELMBUFS` 志但没有置 `MPC_LOADTIME_FLAG_NOTLABEL` 志的策略模必能正地通过入口函数入的 `NULL` 的 Mbuf 安全指。是因那些没有分配存的理中的 Mbuf 在一个需要 Mbuf 安全的策略模加之后，其安全的指将仍然空。如果策略在网子系活之前被加（即，策略不是被推加的），那所有的 Mbuf 的存的分配就可以得到保。

### 6.6.3. 策略入口函数

MAC 框架注册的策略提供四型的入口函数：策略注册和管理入口函数；用于理内核象声明周期事件，如初始化、建和，的入口函数；理策略模感兴趣的控制决策事件的入口函数；以及用于管理象安全的用入口函数。此外，有一个 `mac_syscall()` 入口函数，被策略模用于在不注册新的系用的前提下，展内核接口。

策略模的写人除了必清楚在入特定入口函数之后，些象是可用的之外，熟知内核所采用的加策略。程人在入口函数之内避免使用非叶点，并且遵循和修改象的加程，以降低致死的可能性。特地，程序清楚，然在通常情况下，入入口函数之后，已上了一些，可以安全地象及其安全，但是并不能保证它行修改（包括象本身和其安全）也是安全的。相的上信息，可以参考 MAC 框架入口函数的相文。

策略入口函数把个分指向象本身和其安全的指策略模。一来，即使策略并不熟悉象内部，也能基于作出正决策。只有程信任状个象例外：MAC 框架是假所有的策略模是理解其内部的。

## 6.7. MAC策略入口函数参考

### 6.7.1. 通用的模入口函数

#### 6.7.1.1. `mpo_init`

```
void  
mpo_init (struct mac_policy_conf);
```

参数	明	定
conf	MAC 策略定	

策略加事件。当前程正持有策略表上的互斥，因此是非睡眠的，其他内核子系的用也慎重。

如果需要在策略初始化段执行可能造成睡眠阻塞的内存分配操作，可以将它放在一个独立的模块中集中执行。 SYSINIT()

### 6.7.1.2. mpo\_destroy

```
void  
mpo_destroy (struct mac_policy_conf);
```

参数	说明	约定
conf	MAC 策略定义	

策略加入事件。必须持有策略表互斥锁，因此需要慎重行事。

### 6.7.1.3. mpo\_syscall

```
int  
mpo_syscall (struct thread,  
              int call,  
              void *arg);
```

参数	说明	约定
td	线程	
call	策略特有的系统调用号	
arg	系统调用参数的指针	

入口函数提供策略调用的系统调用，策略模块不需要向应用程序提供的一个外服而注册调用的系统调用。由应用程序提供的策略注册名字来决定提供其所申请的特定策略，所有参数将通过入口函数被调用的策略。当新服务，安全模块必须在必要通过 MAC 框架调用相应的控制机制。比方，假如一个策略调用了某些外的信号功能，那么它调用相应的信号控制，以接受 MAC 框架中注册的其他策略的。



不同的模块需要异地手动执行 copyin() 拷贝系统调用数据。

### 6.7.1.4. mpo\_thread\_userret

```
void  
mpo_thread_userret (struct thread);
```

参数	说明	约定
td	返回线程	

使用入口函数，策略模块能在线程返回用空（系统调用返回、异常返回等等）执行 MAC 相应的清理工作。使用线程的策略需要使用入口函数，因为在清理系统调用的线程中，并不是在任意时刻都能申请到线程的；线程的可能表示线程信息、线程历史或者其他数据。使用入口函数，线程信任状态所作的修改

可能被存放在 `p_label`，且受一个线程自旋的保护；接下来，置线程的 `TDF_ASTPENDING` 状态位和线程的 `PS_MACPENDM` 状态位，表明将调度一个 userret 入口函数的功用。通过入口函数，策略可以在相同上下文中创建信任状的替代品。策略工程师必须清楚，需要保证与调度一个 AST 相同的事件平行次序，同行所执行的 AST 可能很晚，而且在处理多线程应用程序时可能被重入。

## 6.7.2. 操作〇〇

### 6.7.2.1. `mpo_init_bpfdesc_label`

```
void  
mpo_init_bpfdesc_label (struct label);
```

参数	〇明	〇定
label	将被〇用的新〇〇	

一个新近例化的 bpfdesc (BPF 描述子) 初始化〇〇。可以睡眠。

### 6.7.2.2. `mpo_init_cred_label`

```
void  
mpo_init_cred_label (struct label);
```

参数	〇明	〇定
label	将被初始化的新〇〇	

一个新近例化的用〇信任状初始化〇〇。可以睡眠。

### 6.7.2.3. `mpo_init_devfsdirent_label`

```
void  
mpo_init_devfsdirent_label (struct label);
```

参数	〇明	〇定
label	将被〇用的新〇〇	

一个新近例化的 devfs 表〇初始化〇〇。可以睡眠。

### 6.7.2.4. `mpo_init_ifnet_label`

```
void  
mpo_init_ifnet_label (struct label);
```

参数	□明	□定
label	将被□用的新□□	

□一个新近□例化的网□接口初始化□□。可以睡眠。

#### 6.7.2.5. `mpo_init_ipq_label`

```
void
    mpo_init_ipq_label (struct label,
                        int flag);
```

参数	□明	□定
label	将被□用的新□□	
flag	睡眠/不睡眠 <code>malloc(9)</code> ; 参□下文	

□一个新近□例化的 IP 分片重□□列初始化□□。其中的flag域可能取M\_WAITOK 或 M\_NOWAIT之一，用来避免在□初始化□用中因□ `malloc(9)` 而□入睡眠。IP 分片重□□列的分配操作通常是在□性能有□格要求的□境下□行的，因此□代□必□小心地避免睡眠和□□□的操作。IP 分片重□□列分配操作失□□上述入口函数将失□返回。

#### 6.7.2.6. `mpo_init_mbuf_label`

```
void
    mpo_init_mbuf_label (int flag,
                          struct label);
```

参数	□明	□定
flag	睡眠/不睡眠 <code>malloc(9)</code> ; 参□下文	
label	将被初始化的策略□□	

□一个新近□例化的 mbuf 数据包□部 (mbuf) 初始化□□。其中的flag的□可能取M\_WAITOK和 M\_NOWAIT之一，用来避免在□初始化□用中因□ `malloc(9)` 而□入睡眠。Mbuf □部的分配操作常常在□性能有□格要求的□境下被□繁□行，因此□代□必□小心地避免睡眠和□□□的操作。上述入口函数在 Mbuf □部分配操作失□□将失□返回。

#### 6.7.2.7. `mpo_init_mount_label`

```
void
    mpo_init_mount_label (struct label,
                           struct label);
```

参数	□明	□定
mntlabel	将被初始化的mount □□策略□□	

参数	□明	□定
fslabel	将被初始化的文件系□策略□□	

□一个新近□例化的 mount 点初始化□□。可以睡眠。

#### 6.7.2.8. `mpo_init_mount_fs_label`

```
void
    mpo_init_mount_fs_label (struct label);
```

参数	□明	□定
label	将被初始化的□□	

□一个新近加□的文件系□初始化□□。可以睡眠。

#### 6.7.2.9. `mpo_init_pipe_label`

```
void
    mpo_init_pipe_label (struct);
```

参数	□明	□定
label	将被填写的□□	

□一个□□□例化的管道初始化安全□□。可以睡眠。

#### 6.7.2.10. `mpo_init_socket_label`

```
void
    mpo_init_socket_label (struct label,
                           int flag);
```

参数	□明	□定
label	将被初始化的新□□	
flag	<code>malloc(9)</code> flags	

□一个□□□例化的套接字初始化安全□□。其中的 `flag` 域的□必□被指定□ `M_WAITOK` 和 `M_NOWAIT` 之一，以避免在□初始化程中使用可能睡眠的`malloc(9)`。

#### 6.7.2.11. `mpo_init_socket_peer_label`

```
void
    mpo_init_socket_peer_label (struct label,
                                int flag);
```

参数	□明	□定
label	将被初始化的新□	
flag	malloc(9) flags	

□□□例化的套接字□等体□行□□的初始化。其中的 flag 域的□必□被指定□ M\_WAITOK 和 M\_NOWAIT 之一，以避免在□初始化程中使用可能睡眠的 malloc(9)。

#### 6.7.2.12. mpo\_init\_proc\_label

```
void
    mpo_init_proc_label (struct label);
```

参数	□明	□定
label	将被初始化的新□	

□一个□□□例化的□程初始化安全□□。可以睡眠。

#### 6.7.2.13. mpo\_init\_vnode\_label

```
void
    mpo_init_vnode_label (struct label);
```

参数	□明	□定
label	将被初始化的新□	

□一个□□□例化的 vnode 初始化安全□□。可以睡眠。

#### 6.7.2.14. mpo\_destroy\_bpfdesc\_label

```
void
    mpo_destroy_bpfdesc_label (struct label);
```

参数	□明	□定
label	bpfdesc □□	

□一个 BPF 描述子上的□□。在□入口函数中，策略□当□放所有在内部分配与 label 相□□的存□空□，以便□□□□□。

### 6.7.2.15. mpo\_destroy\_cred\_label

```
void  
    mpo_destroy_cred_label (struct label);
```

参数	□明	□定
label	将被□□的□□	

□□一个信任状上的□□。在□□入口函数中，策略□□当□□放所有在内部分配的与 label 相□□的存□□空□□，以便□□□□□□。

### 6.7.2.16. mpo\_destroy\_devfs dirent\_label

```
void  
    mpo_destroy_devfs dirent_label (struct label);
```

参数	□明	□定
label	将被□□的□□	

□□一个 devfs 表□□上的□□。在□□入口函数中，策略□□当□□放所有在内部分配的与 label 相□□的存□□空□□，以便□□□□□□。

### 6.7.2.17. mpo\_destroy\_ifnet\_label

```
void  
    mpo_destroy_ifnet_label (struct label);
```

参数	□明	□定
label	将被□□的□□	

□□与一个已□□除接口相□□的□□。在□□入口函数中，策略□□当□□放所有在内部分配的与 label 相□□的存□□空□□，以便□□□□□□。

### 6.7.2.18. mpo\_destroy\_ipq\_label

```
void  
    mpo_destroy_ipq_label (struct label);
```

参数	□明	□定
label	将被□□的□□	

□□与一个 IP 分片□□列相□□的□□。在□□入口函数中，策略□□当□□放所有在内部分配的与 label 相□□的存□□空□□，以便□□□□□□。

### 6.7.2.19. `mpo_destroy_mbuf_label`

```
void  
    mpo_destroy_mbuf_label (struct label);
```

参数	□明	□定
label	将被□的□	

□与一个 Mbuf 相□的□。在□入口函数中，策略□当□放所有在内部分配的与 label 相□的存□空□，以便□□□□。

### 6.7.2.20. `mpo_destroy_mount_label`

```
void  
    mpo_destroy_mount_label (struct label);
```

参数	□明	□定
label	将被□的 Mount 点□	

□与一个 mount 点相□的□。在□入口函数中，策略□当□放所有在内部分配的与 mntlabel 相□的存□空□，以便□□□□。

### 6.7.2.21. `mpo_destroy_mount_label`

```
void  
    mpo_destroy_mount_label (struct label,  
                            struct label);
```

参数	□明	□定
mntlabel	将被□的 Mount 点□	
fslabel	File system label being destroyed	

□与一个 mount 点相□的□。在□入口函数中，策略□当□放所有在内部分配的，与 mntlabel 和fslabel 相□的存□空□，以便□□□□。

### 6.7.2.22. `mpo_destroy_socket_label`

```
void  
    mpo_destroy_socket_label (struct label);
```

参数	□明	□定
label	将被□的套接字□	

□与一个套接字相□的□。在□入口函数中，策略□当□放所有在内部分配的，与 label 相□的存□空□，以便□□□□。

#### 6.7.2.23. `mpo_destroy_socket_peer_label`

```
void
    mpo_destroy_socket_peer_label (struct label);
```

参数	□明	□定
peerlabel	将被□的套接字□等□体□	

□与一个套接字相□的□等□体□。在□入口函数中，策略□当□放所有在内部分配的，与 label 相□的存□空□，以便□□□□。

#### 6.7.2.24. `mpo_destroy_pipe_label`

```
void
    mpo_destroy_pipe_label (struct label);
```

参数	□明	□定
label	管道□	

□一个管道的□。在□入口函数中，策略□当□放所有在内部分配的，与 label 相□的存□空□，以便□□□□。

#### 6.7.2.25. `mpo_destroy_proc_label`

```
void
    mpo_destroy_proc_label (struct label);
```

参数	□明	□定
label	□程□	

□一个□程的□。在□入口函数中，策略□当□放所有在内部分配的，与 label 相□的存□空□，以便□□□□。

#### 6.7.2.26. `mpo_destroy_vnode_label`

```
void
    mpo_destroy_vnode_label (struct label);
```

参数	□明	□定
label	□程□□	

□一个 vnode 的□□。在□入口函数中，策略□当□放所有在内部分配的，与 label 相□□的存□空□，以便□□□□。

#### 6.7.2.27. `mpo_copy_mbuf_label`

```
void
mpo_copy_mbuf_label (struct label,
                      struct label);
```

参数	□明	□定
src	源□□	
dest	目□□□	

将 src 中的□□信息拷□到 dest中。

#### 6.7.2.28. `mpo_copy_pipe_label`

```
void
mpo_copy_pipe_label (struct label,
                      struct label);
```

参数	□明	□定
src	源□□	
dest	目□□□	

将 src 中的□□信息拷□至 dest。

#### 6.7.2.29. `mpo_copy_vnode_label`

```
void
mpo_copy_vnode_label (struct label,
                      struct label);
```

参数	□明	□定
src	源□□	
dest	目□□□	

将 src 中的□□信息拷□至 dest。

### 6.7.2.30. mpo\_externalize\_cred\_label

```
int  
mpo_externalize_cred_label (struct label *label,  
                           char *element_name,  
                           struct sbuf *sb,  
                           int *claimed);
```

参数	说明	约定
label	将用外部形式表示的 <b>label</b>	
element_name	需要外部表示 <b>label</b> 的策略的名字	
sb	用来存放 <b>label</b> 的文本表示形式的字符 buffer	
claimed	如果可以填充 <b>element_data</b> 域, 其数 <b>sb</b>	

根据**label**入的**element\_name**, 生成一个以外部形式表示的**label**。一个外部形式**label**, 是**label**内容的文本表示, 它由**label**的程序使用, 是用**label**可**label**的。目前的MAC**label**方案将依次**label**用策略的相**label**入口函数, 因此, 具体策略的**label**代**label**, 需要在填写**sb**之前, 先**label****element\_name**中指定的名字。如果**element\_name**中的内容与**label**的策略名字不相符, 直接返回0。**label**当**sb**中**label**数据的**label**程中出**label**时, 才返回非0。一旦策略决定填写**element\_data**, **label**\***claim**的数**label**。

### 6.7.2.31. mpo\_externalize\_ifnet\_label

```
int  
mpo_externalize_ifnet_label (struct label *label,  
                           char *element_name,  
                           struct sbuf *sb,  
                           int *claimed);
```

参数	说明	约定
label	将用外部形式表示的 <b>label</b>	
element_name	需要外部表示 <b>label</b> 的策略的名字	
sb	用来存放 <b>label</b> 的文本表示形式的字符 buffer	
claimed	如果可以填充 <b>element_data</b> 域, 其数 <b>sb</b>	

根据**label**入的**element\_name**, 生成一个以外部形式表示的**label**。一个外部形式**label**, 是**label**内容的文本表示, 它由**label**的程序使用, 是用**label**可**label**的。目前的MAC**label**方案将依次**label**用策略的相**label**入口函数, 因此, 具体策略的**label**代**label**, 需要在填写**sb**之前, 先**label****element\_name**中指定的名字。如果**element\_name**中的内容与**label**的策略名字不相符, 直接返回0。**label**当**sb**中**label**数据的**label**程中出**label**时, 才返回非0。一旦策略决定填写**element\_data**, **label**\***claim**的数**label**。

### 6.7.2.32. mpo\_externalize\_pipe\_label

```
int  
mpo_externalize_pipe_label (struct label *label,  
                           char *element_name,  
                           struct sbuf *sb,  
                           int *claimed);
```

参数	□明	□定
label	将用外部形式表示的□□	
element_name	需要外部表示□□的策略的名字	
sb	用来存放□□的文本表示形式的字符 buffer	
claimed	如果可以填充element_data 域, □其数□□□	

根据□入的□□□□， □生一个以外部形式表示的□□。 一个外部形式□□， 是□□内容的文本表示， 它由□□的□用程序使用， 是用□可□的。 目前的MAC□□方案将依次□用策略的相□入口函数， 因此， 具体策略的□□代□， 需要在填写sb之前， 先□□element\_name中指定的名字。 如果element\_name中的内容与□的策略名字不相符， □直接返回0。 □当□□□□数据的□程中出□□□□， 才返回非0□。 一旦策略决定填写element\_data， □□\*claim的数□。

### 6.7.2.33. mpo\_externalize\_socket\_label

```
int  
mpo_externalize_socket_label (struct label *label,  
                               char *element_name,  
                               struct sbuf *sb,  
                               int *claimed);
```

参数	□明	□定
label	将用外部形式表示的□□	
element_name	需要外部表示□□的策略的名字	
sb	用来存放□□的文本表示形式的字符 buffer	
claimed	如果可以填充element_data 域, □其数□□□	

根据□入的□□□□， □生一个以外部形式表示的□□。 一个外部形式□□， 是□□内容的文本表示， 它由□□的□用程序使用， 是用□可□的。 目前的MAC□□方案将依次□用策略的相□入口函数， 因此， 具体策略的□□代□， 需要在填写sb之前， 先□□element\_name中指定的名字。 如果element\_name中的内容与□的策略名字不相符， □直接返回0。 □当□□□□数据的□程中出□□□□， 才返回非0□。 一旦策略决定填写element\_data， □□\*claim的数□。

#### 6.7.2.34. mpo\_externalize\_socket\_peer\_label

```
int  
    mpo_externalize_socket_peer_label (struct label *label,  
                                         char *element_name,  
                                         struct sbuf *sb,  
                                         int *claimed);
```

参数	□明	□定
label	将用外部形式表示的□□	
element_name	需要外部表示□□的策略的名字	
sb	用来存放□□的文本表示形式的字符 buffer	
claimed	如果可以填充element_data 域, □其数□□□	

根据□入的□□□□， □生一个以外部形式表示的□□。 一个外部形式□□， 是□□内容的文本表示， 它由□□的□用程序使用， 是用□可□的。 目前的MAC□□方案将依次□用策略的相□入口函数， 因此， 具体策略的□□代□， 需要在填写sb之前， 先□□element\_name中指定的名字。 如果element\_name中的内容与□的策略名字不相符， □直接返回0。 □当□□□□数据的□程中出□□□□， 才返回非0□。 一旦策略决定填写element\_data， □□\*claim的数□。

#### 6.7.2.35. mpo\_externalize\_vnode\_label

```
int  
    mpo_externalize_vnode_label (struct label *label,  
                                 char *element_name,  
                                 struct sbuf *sb,  
                                 int *claimed);
```

参数	□明	□定
label	将用外部形式表示的□□	
element_name	需要外部表示□□的策略的名字	
sb	用来存放□□的文本表示形式的字符 buffer	
claimed	如果可以填充element_data 域, □其数□□□	

根据□入的□□□□， □生一个以外部形式表示的□□。 一个外部形式□□， 是□□内容的文本表示， 它由□□的□用程序使用， 是用□可□的。 目前的MAC□□方案将依次□用策略的相□入口函数， 因此， 具体策略的□□代□， 需要在填写sb之前， 先□□element\_name中指定的名字。 如果element\_name中的内容与□的策略名字不相符， □直接返回0。 □当□□□□数据的□程中出□□□□， 才返回非0□。 一旦策略决定填写element\_data， □□\*claim的数□。

#### 6.7.2.36. mpo\_internalize\_cred\_label

```
int  
mpo_internalize_cred_label (struct label *label,  
                           char *element_name,  
                           char *element_data,  
                           int *claimed);
```

参数	□明	□定
label	将被填充的□□	
element_name	需要内部表示□□的策略的名字	
element_data	需要被□□的文本数据	
claimed	如果数据被正□□□， □其数□□□	

根据一个文本形式的外部表示□□数据， □建一个内部形式的□□□□。 目前的MAC方案将依次□用所有策略的相□入口函数， 来□□□□的内部□□□求， 因此， □□代□必□首先通□比□element\_name中的内容和自己的策略名字， 来□定是否需要□□element\_data中存放的数据。 □似的， 如果名字不匹配或者数据□□操作成功， □函数返回0， 并□□\*claimed的□。

#### 6.7.2.37. mpo\_internalize\_ifnet\_label

```
int  
mpo_internalize_ifnet_label (struct label *label,  
                           char *element_name,  
                           char *element_data,  
                           int *claimed);
```

参数	□明	□定
label	将被填充的□□	
element_name	需要内部表示□□的策略的名字	
element_data	需要被□□的文本数据	
claimed	如果数据被正□□□， □其数□□□	

根据一个文本形式的外部表示□□数据， □建一个内部形式的□□□□。 目前的MAC方案将依次□用所有策略的相□入口函数， 来□□□□的内部□□□求， 因此， □□代□必□首先通□比□element\_name中的内容和自己的策略名字， 来□定是否需要□□element\_data中存放的数据。 □似的， 如果名字不匹配或者数据□□操作成功， □函数返回0， 并□□\*claimed的□。

#### 6.7.2.38. mpo\_internalize\_pipe\_label

```

int
    mpo_internalize_pipe_label (struct label *label,
                                char *element_name,
                                char *element_data,
                                int *claimed);

```

参数	□明	□定
label	将被填充的□	
element_name	需要内部表示□的策略的名字	
element_data	需要被□的文本数据	
claimed	如果数据被正□□， □其数□□	

根据一个文本形式的外部表示□数据， □建一个内部形式的□□□。 目前的MAC方案将依次□用所有策略的相□入口函数， 来□□□的内部□□□求， 因此， □□代□必□首先通□比□element\_name中的内容和自己的策略名字， 来□定是否需要□□element\_data中存放的数据。 □似的， 如果名字不匹配或者数据□□操作成功， □函数返回0， 并□□\*claimed的□。

#### 6.7.2.39. mpo\_internalize\_socket\_label

```

int
    mpo_internalize_socket_label (struct label *label,
                                   char *element_name,
                                   char *element_data,
                                   int *claimed);

```

参数	□明	□定
label	将被填充的□	
element_name	需要内部表示□的策略的名字	
element_data	需要被□的文本数据	
claimed	如果数据被正□□， □其数□□	

根据一个文本形式的外部表示□数据， □建一个内部形式的□□□。 目前的MAC方案将依次□用所有策略的相□入口函数， 来□□□的内部□□□求， 因此， □□代□必□首先通□比□element\_name中的内容和自己的策略名字， 来□定是否需要□□element\_data中存放的数据。 □似的， 如果名字不匹配或者数据□□操作成功， □函数返回0， 并□□\*claimed的□。

#### 6.7.2.40. mpo\_internalize\_vnode\_label

```

int
mpo_initialize_vnode_label (struct label *label,
                           char *element_name,
                           char *element_data,
                           int *claimed);

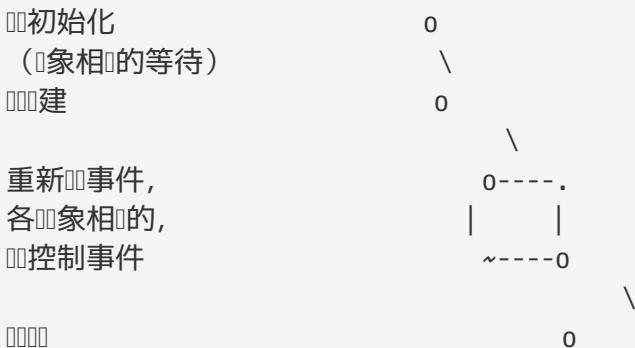
```

参数	说明	判定
label	将被填充的	
element_name	需要内部表示的策略的名字	
element_data	需要被的文本数据	
claimed	如果数据被正，其数	

根据一个文本形式的外部表示的数据，创建一个内部形式的。目前的MAC方案将依次用所有策略的相入口函数，来内部的内部要求。因此，代必须首先通过比element\_name中的内容和自己的策略名字，来判定是否需要element\_data中存放的数据。似的，如果名字不匹配或者数据操作成功，函数返回0，并\*claimed的。

### 6.7.3. 事件

策略模块使用MAC 框架提供的“事件”入口函数，内核对象的进行操作。策略模块将感兴趣的被内核对象的相生命周期事件 注册在恰当的入口点上。对象的初始化、创建和事件均提供了子点。在某些对象上可以重新，即，允许编程改变对象上的。某些对象可以其特定的对象事件，比如与IP 重相关联的事件。一个典型的被对象在其生命周期中将有下列入口函数：



使用初始化入口函数，策略可以以一一的、与对象使用境无的方式置对象的初始。分配一个策略的缺省 slot 0，不使用对象的策略可能并不需要执行对象的初始化操作。

的创建事件发生在将一个内核数据同一个真的内核对象相（内核对象实例化）的时刻。例如，在真正被使用之前，在一个缓冲池内已分配的 mbuf 数据，将保持“未使用”状态。因此，mbuf 的分配操作将导致 mbuf 的初始化操作，而 mbuf 的创建操作被推到 mbuf 真正与一个数据相的时刻。

通常，用者将会提供创建事件的上下文，包括创建境、创建程序中及的其他对象的等。例如，当一个套接字创建一个 mbuf 时，除了新建的 mbuf 及其之外，作新建者的套接字与其也被提交策略。

不提倡在对象就其分配内存的原因有二个：创建操作可能发生在性能有严格要求的内核接口上；而且，因创建用不允许失败，所以无法宣告内存分配失败。

对象特有的事件一般不会引起其他的事件，但是在对象上下文发生改变时，策略使用它可以对相对象进行修改或更新操作。例如，在MAC\_UPDATE\_IPQ 入口函数之内，某个 IP 分片重排列的可能会因排列中接收了新的mbuf 而被更新。

控制事件将在后章中。

策略通常执行操作，释放其分配的内存或的状况，之后内核才可以重用或者释放对象的内核数据。

除了与特定内核对象定的普通之外，有一类外的类型：。这些用于存放由用户程序提交的更新信息。它的初始化和操作与其他一样，只是创建事件，MAC\_INTERNALIZE，略有不同：函数接受用户提交的，将其转化为内核表示形式。

### 6.7.3.1. 文件系统对象事件操作

#### 6.7.3.1.1. mpo\_associate\_vnode\_devfs

```
void  
mpo_associate_vnode_devfs (struct mount,  
                           struct label,  
                           struct devfs_dirent,  
                           struct label,  
                           struct vnode,  
                           struct label);
```

参数	说明	约定
mp	Devfs 挂点	
fslabel	Devfs 文件系统 (mp- mnt_fslabel)	
de	Devfs 目录	
delabel	与 de 相同的策略	
vp	与 de 相同的 vnode	
vlabel	与 vp 相同的策略	

根据参数 de 传入的 devfs 目录及其信息，一个新近创建的 devfs vnode 填充 (vlabel)。

#### 6.7.3.1.2. mpo\_associate\_vnode\_extattr

```
int  
mpo_associate_vnode_extattr (struct mount,  
                           struct label,  
                           struct vnode,  
                           struct label);
```

参数	说明	约定
mp	文件系统挂点	

参数	□明	□定
fslabel	文件系□□□	
vp	将被□□的 vnode	
vlabel	与 vp 相□□的策略□□	

从文件系□□□展属性中□取 vp 的□□。成功，返回 0。不成功，□在 errno 指定的相□的□□□□。如果文件系□不支持□□□展属性的□□□操作，□可以考□将 fslabel 拷□至 vlabel。

#### 6.7.3.1.3. `mpo_associate_vnode_singlelabel`

```
void
mpo_associate_vnode_singlelabel (struct mount,
                                 struct label,
                                 struct vnode,
                                 struct label);
```

参数	□明	□定
mp	文件系□挂□点	
fslabel	文件系□□□	
vp	将被□□的 vnode	
vlabel	与 vp 相□□的策略□□	

在非多重□□□上，使用□入口函数，根据文件系□□□， fslabel， □ vp □置策略□□。

#### 6.7.3.1.4. `mpo_create_devfs_device`

```
void
mpo_create_devfs_device (dev_t dev,
                         struct devfs_dirent,
                         struct label);
```

参数	□明	□定
dev	devfs_dirent □□的□□	
devfs_dirent	将被□□的 Devfs 目□□	
label	将被填写的 devfs_dirent □□	

□□入□□新建的 devfs\_dirent 填写□□。□函数将在□□□文件系□加□、重□或添加新□□□被□用。

#### 6.7.3.1.5. `mpo_create_devfs_directory`

```

void
    mpo_create_devfs_directory (char *dirname,
                                int dirnamelen,
                                struct devfs_dirent,
                                struct label);

```

参数	□明	□定
dirname	新建目□的名字	
namelen	字符串 dirname 的□度	
devfs_dirent	新建目□在 Devfs 中□的目□	

□入目□参数的新建 devfs\_dirent 填写□□。□函数将在加□、重□□□文件系□，或者添加一个需要指定目□□□的新□□□被□用。

#### 6.7.3.1.6. mpo\_create\_devfs\_symlink

```

void
    mpo_create_devfs_symlink (struct ucred,
                               struct mount,
                               struct devfs_dirent,
                               struct label,
                               struct devfs_dirent,
                               struct label);

```

参数	□明	□定
cred	主体信任状	
mp	devfs 挂□点	
dd	□接目□	
ddlabel	与 dd 相□□的□□	
de	符号□接□	
delabel	与 de 相□□的策略□□	

□新近□建的 [devfs\(5\)](#) 符号□接□填写□□ (delabel) 。

#### 6.7.3.1.7. mpo\_create\_vnode\_extattr

```

int
mpo_create_vnode_extattr (struct ucred,
                          struct mount,
                          struct label,
                          struct vnode,
                          struct label,
                          struct vnode,
                          struct label,
                          struct componentname);

```

参数	□明	□定
cred	主体信任状	
mount	文件系□挂□点	
label	文件系□□	
dvp	父目□ vnode	
dlabel	与 dvp 相□□的策略□□	
vp	新□建的 vnode	
vlabel	与 vp 相□□的策略□□	
cnp	vp中的子域名字	

将 vp 的□□写入文件□展属性。成功，将□□填入 vlabel， 并返回 0。否□， 返回□□的□□□□。

#### 6.7.3.1.8. mpo\_create\_mount

```

void
mpo_create_mount (struct ucred,
                   struct mount,
                   struct label,
                   struct label);

```

参数	□明	□定
cred	主体信任状	
mp	客体；将被挂□的文件系□	
mntlabel	将被填写的 mp 的策略□□	
fslabel	将被挂□到 mp 的文件系□的策略□□。	

□□入的主体信任状所□建的挂□点填写□□。□函数将在文件系□挂□□被□用。

#### 6.7.3.1.9. mpo\_create\_root\_mount

```

void
    mpo_create_root_mount (struct ucred,
                           struct mount,
                           struct label,
                           struct label);

```

参数	□明	□定
□ mpo_create_mount.		

□入的主体信任状所建的挂点填写□。□函数将在挂根文件系□， mpo\_create\_mount; 之后被□用。

#### 6.7.3.1.10. mpo\_relabel\_vnode

```

void
    mpo_relabel vnode (struct ucred,
                       struct vnode,
                       struct label,
                       struct label);

```

参数	□明	□定
cred	主体信任状	
vp	将被重新□的 vnode	
vnode_label	vp □有的策略□	
newlabel	将取代vnode_label的新（可能只是部分）□	

根据□入的新□和主体信任状，更新参数 vnode 的□。

#### 6.7.3.1.11. mpo\_setlabel vnode\_extattr

```

int
    mpo_setlabel vnode_extattr (struct ucred,
                                struct vnode,
                                struct label,
                                struct label);

```

参数	□明	□定
cred	主体信任状	
vp	写出□所□的 vnode	
vlabel	vp的策略□	
intlabel	将被写入磁□的□	

将参数 intlabel 传出的策略信息写入指定 vnode 的扩展属性。该函数被 [vop\\_stdcreatevnode\\_ea](#) 所调用。

#### 6.7.3.1.12. [mpo\\_update\\_devfsdirent](#)

```
void  
mpo_update_devfsdirent (struct devfs_dirent,  
                         struct label,  
                         struct vnode,  
                         struct label);
```

参数	说明	约定
devfs_dirent	客体；devfs 目录	
direntlabel	将被更新的devfs_dirent的策略	
vp	父 vnode	已约定
vnode_label	vp的策略	

根据所传入的 devfs vnode，对 devfs\_dirent 的策略进行更新。重新分配一个 devfs vnode 的操作成功之后，将调用该函数来更新它的策略，如此，即使相同的 vnode 数据被内核回收重用，也不会丢失它的新状态。另外，在 devfs 中新建一个符号链接，接着 [mac\\_vnode\\_create\\_from\\_vnode](#)，也将调用该函数，对 vnode 进行初始化操作。

#### 6.7.3.2. IPC 对象事件操作

##### 6.7.3.2.1. [mpo\\_create\\_mbuf\\_from\\_socket](#)

```
void  
mpo_create_mbuf_from_socket (struct socket,  
                             struct label,  
                             struct mbuf *m,  
                             struct label);
```

参数	说明	约定
socket	套接字	套接字约定 WIP
socketlabel	socket 的策略	
m	客体；mbuf	
mbuflabel	将被填写的 m 的策略	

根据传入的套接字新建的 mbuf 部分设置。当套接字产生一个新的数据或者消息，并将其存储在参数 mbuf 中，将调用该函数。

##### 6.7.3.2.2. [mpo\\_create\\_pipe](#)

```

void
    mpo_create_pipe (struct ucred,
                     struct pipe,
                     struct label);

```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe 的策略□□	

根据□入的主体信任状参数，□置新建管道的□□。□当一个新管道被□建，□函数将被□用。

#### 6.7.3.2.3. mpo\_create\_socket

```

void
    mpo_create_socket (struct ucred,
                       struct socket,
                       struct label);

```

参数	□明	□定
cred	主体信任状	不可改□
so	客体；将被□□的套接字	
socketlabel	将被填写的 so 的□□	

根据□入的主体信任状参数，□置新建套接字的□□。□当新建一个套接字，□函数将被□用。

#### 6.7.3.2.4. mpo\_create\_socket\_from\_socket

```

void
    mpo_create_socket_from_socket (struct socket,
                                   struct label,
                                   struct socket,
                                   struct label);

```

参数	□明	□定
oldsocket	□听套接字	
oldsocketlabel	oldsocket 的策略□□	
newsocket	新建套接字	
newsocketlabel	newsocket 的策略□□	

根据 [listen\(2\)](#) 套接字 oldsocket，□新建 [accept\(2\)](#) 的套接字 newsocket，□置□□。

#### 6.7.3.2.5. mpo\_relabel\_pipe

```
void  
mpo_relabel_pipe (struct ucred,  
                   struct pipe,  
                   struct label,  
                   struct label);
```

参数	□明	□定
cred	主体信任状	
pipe	管道	
oldlabel	pipe 的当前策略□□	
newlabel	将□pipe □置的新的策略□□	

□pipe□置新□□newlabel。

#### 6.7.3.2.6. mpo\_relabel\_socket

```
void  
mpo_relabel_socket (struct ucred,  
                     struct socket,  
                     struct label,  
                     struct label);
```

参数	□明	□定
cred	主体信任状	不可改□
so	客体；套接字	
oldlabel	so 的当前□□	
newlabel	将□socket □置的新□□	

根据□入的□□参数， □套接字的当前□□□行更新。

#### 6.7.3.2.7. mpo\_set\_socket\_peer\_from\_mbuf

```
void  
mpo_set_socket_peer_from_mbuf (struct mbuf,  
                                struct label,  
                                struct label,  
                                struct label);
```

参数	□明	□定
mbuf	从套接字接收到的第一个数据□	

参数	□明	□定
mbuflabel	mbuf 的□□	
oldlabel	套接字的当前□□	
newlabel	将□套接字填写的策略□□	

根据□入的 mbuf □□, □置某个 stream 套接字的□等□志。除Unix域的套接字之外, □当一个 stream 套接字接收到第一个数据□□, □函数将被□用。

#### 6.7.3.2.8. `mpo_set_socket_peer_from_socket`

```
void
mpo_set_socket_peer_from_socket (struct socket,
                                 struct label,
                                 struct socket,
                                 struct label);
```

参数	□明	□定
oldsocket	本地套接字	
oldsocketlabel	oldsocket 的策略□□	
newsocket	□等套接字	
newsocketpeerlabel	将□newsocket填写的策略□□	

根据□入的□程套接字端点, □一个 stream UNIX 与套接字□置□等□□。□当相□的套接字□之□□行□接□, □函数将在□端分□被□用。

#### 6.7.3.3. Network Object Labeling Event Operations

##### 6.7.3.3.1. `mpo_create_bpfdesc`

```
void
mpo_create_bpfdesc (struct ucred,
                     struct bpf_d,
                     struct label);
```

参数	□明	□定
cred	主体信任状	不可改□
bpf_d	客体 ; bpf 描述子	
bpf	将□bpf_d填写的策略□□	

根据□入的主体信任状参数, □新建的 BPF 描述子□置□□。当□程打□ BPF □□□点□, □函数将被□用。

#### 6.7.3.3.2. mpo\_create\_ifnet

```
void  
mpo_create_ifnet (struct ifnet,  
                   struct label);
```

参数	□明	□定
ifnet	网□接口	
ifnetlabel	将□ifnet填写的策略□□	

□新建的网□接口□置□□。□函数在以下情况下被□用：当一个新的物理接口□□可用□，或者当一个□接口在引□□或由于某个用□操作而□例化□。

#### 6.7.3.3.3. mpo\_create\_ipq

```
void  
mpo_create_ipq (struct mbuf,  
                  struct label,  
                  struct ipq,  
                  struct label);
```

参数	□明	□定
fragment	第一个被接收的 IP 分片	
fragmentlabel	fragment 的策略□□	
ipq	将被□□的 IP 重□□列	
ipqlabel	将□ipq填写的策略□□	

根据第一个接收到的分片的 mbuf □部信息，□新建的 IP 分片重□□列□置□□。

#### 6.7.3.3.4. mpo\_create\_datagram\_from\_ipq

```
void  
mpo_create_create_datagram_from_ipq (struct ipq,  
                                       struct label,  
                                       struct mbuf,  
                                       struct label);
```

参数	□明	□定
ipq	IP 重□□列	
ipqlabel	ipq 的策略□□	
datagram	将被□□的数据□	
datagramlabel	将□datagramlabel填写的策略□□	

根据 IP 分片重~~复~~列，~~复~~重~~复~~完~~复~~的 IP 数据~~复~~置~~复~~。

#### 6.7.3.3.5. `mpo_create_fragment`

```
void
    mpo_create_fragment (struct mbuf,
                          struct label,
                          struct mbuf,
                          struct label);
```

参数	说明	约定
datagram	数据 <del>复</del>	
datagramlabel	datagram 的策略 <del>复</del>	
fragment	将被 <del>复</del> 的分片	
fragmentlabel	将 <del>复</del> datagram 填写的策略 <del>复</del>	

根据数据~~复~~所~~复~~的 mbuf ~~复~~部信息，~~复~~其新建的分片的 mbuf ~~复~~部~~复~~置~~复~~。

#### 6.7.3.3.6. `mpo_create_mbuf_from_mbuf`

```
void
    mpo_create_mbuf_from_mbuf (struct mbuf,
                               struct label,
                               struct mbuf,
                               struct label);
```

参数	说明	约定
oldmbuf	已有的（源） mbuf	
oldmbuflabel	oldmbuf 的策略 <del>复</del>	
newmbuf	将被 <del>复</del> 的新建 mbuf	
newmbuflabel	将 <del>复</del> newmbuf 填写的策略 <del>复</del>	

根据某个~~复~~有数据~~复~~的 mbuf ~~复~~部信息，~~复~~新建数据~~复~~的 mbuf ~~复~~部~~复~~置~~复~~。在~~复~~多条件下将会~~复~~用~~复~~函数，比如，由于~~复~~要求而重新分配某个 mbuf ~~复~~。

#### 6.7.3.3.7. `mpo_create_mbuf_linklayer`

```
void
    mpo_create_mbuf_linklayer (struct ifnet,
                               struct label,
                               struct mbuf,
                               struct label);
```

参数	□明	□定
ifnet	网口接口	
ifnetlabel	ifnet 的策略□□	
mbuf	新建数据□的 mbuf □部	
mbuflabel	将□mbuf填写的策略□□	

□在□定接口上由于某个□路□□而新建的数据□的mbuf□部□置□□。 □函数将在若干条件下被□用， 比如当IPv4和IPv6□□□在□□ARP或者ND6□。

#### 6.7.3.3.8. `mpo_create_mbuf_from_bpfdesc`

```
void
    mpo_create_mbuf_from_bpfdesc (struct bpf_d,
                                struct label,
                                struct mbuf,
                                struct label);
```

参数	□明	□定
bpf_d	BPF 描述子	
bpflabel	bpflabel 的策略□□	
mbuf	将被□□的新建 mbuf	
mbuflabel	将□mbuf填写的策略□□	

□使用参数 BPF 描述子□建的新数据□的 mbuf □部□置□□。 当□参数 BPF 描述子所□□的 BPF □□□行写操作□， □函数将被□用。

#### 6.7.3.3.9. `mpo_create_mbuf_from_ifnet`

```
void
    mpo_create_mbuf_from_ifnet (struct ifnet,
                                struct label,
                                struct mbuf,
                                struct label);
```

参数	□明	□定
ifnet	网口接口	
ifnetlabel	ifnetlabel 的策略□□	
mbuf	新建数据□的 mbuf □部	
mbuflabel	将□mbuf填写的策略□□	

□从网口接口参数□建的数据□的 mbuf □部□置□□。

#### 6.7.3.3.10. `mpo_create_mbuf_multicast_encap`

```
void
    mpo_create_mbuf_multicast_encap (struct mbuf,
                                      struct label,
                                      struct ifnet,
                                      struct label,
                                      struct mbuf,
                                      struct label);
```

参数	□明	□定
oldmbuf	□有数据□的 mbuf □部	
oldmbuflabel	oldmbuf 的策略□□	
ifnet	网□接口	
ifnetlabel	ifnet 的策略□□	
newmbuf	将被□□的新建数据□ mbuf □部	
newmbuflabel	将□newmbuf填写的策略□□	

当□入的已有数据□被□定多播封装接口 (multicast encapsulation interface) □理□被□用， □新□建的数据□所在 mbuf □部□置□□。 □当使用□虚□接口□□一个mbuf□， 将□用□函数。

#### 6.7.3.3.11. `mpo_create_mbuf_netlayer`

```
void
    mpo_create_mbuf_netlayer (struct mbuf,
                               struct label,
                               struct mbuf,
                               struct label);
```

参数	□明	□定
oldmbuf	接收的数据□	
oldmbuflabel	oldmbuf 的策略□□	
newmbuf	新建数据□	
newmbuflabel	newmbuf 的策略□□	

□由 IP 堆□因□□□接收数据□ (oldmbuf) 而新建的数据□□置其 mbuf □部的□□。 □多情况下需要□用□函数， 比如， □□ ICMP □求数据□□。

#### 6.7.3.3.12. `mpo_fragment_match`

```

int
mpo_fragment_match (struct mbuf,
                     struct label,
                     struct ipq,
                     struct label);

```

参数	说明	约定
fragment	IP 数据分片	
fragmentlabel	fragment 的策略	
ipq	IP 分片重叠列	
ipqlabel	ipq 的策略	

根据所入的 IP 分片重叠列 (ipq) 的部，部包含一个 IP 数据 (fragment) 的 mbuf 的部是否符合其要求。如果符合，返回1。否则，返回0。当接收到的分片放入某个已有的分片重叠列中时，将用函数进行安全检查；如果失败，将分片重新实例化一个新的分片重叠列。策略可以利用入口函数，根据或者其他信息阻止不期望的 IP 分片重叠。

#### 6.7.3.3.13. mpo\_relabel\_ifnet

```

void
mpo_relabel_ifnet (struct ucred,
                    struct ifnet,
                    struct label,
                    struct label);

```

参数	说明	约定
cred	主体信任状	
ifnet	客体；网口接口	
ifnetlabel	ifnet 的策略	
newlabel	将ifnet置的新	

根据所入的新， newlabel， 以及主体信任状， cred， 网口接口的进行更新。

#### 6.7.3.3.14. mpo\_update\_ipq

```

void
mpo_update_ipq (struct mbuf,
                 struct label,
                 struct ipq,
                 struct label);

```

参数	□明	□定
mbuf	IP 分片	
mbuflabel	mbuf 的策略□□	
ipq	IP 分片重□□列	
ipqlabel	将被更新的ipq的当前策略□□	

根据所□入的 IP 分片 mbuf □部 (mbuf) □接收 它的 IP 分片重□□列 (ipq) 的□□□行更新。

#### 6.7.3.4. □程□□事件操作

##### 6.7.3.4.1. `mpo_create_cred`

```
void
mpo_create_cred (struct ucred,
                  struct ucred);
```

参数	□明	□定
parent_cred	父主体信任状	
child_cred	子主体信任状	

根据所□入的主体信任状， □新建的主体信任状□置□□。 □当□一个新建的 struct ucred□用 `crcopy(9)` □， 将□用此函数。 □函数不□与□程□制 (forking) 或者□建事件混□一□。

##### 6.7.3.4.2. `mpo_execve_transition`

```
void
mpo_execve_transition (struct ucred,
                       struct ucred,
                       struct vnode,
                       struct label);
```

参数	□明	□定
old	已有的主体信任状	不可改□
new	将被□□的新主体信任状	
vp	将被□行的文件	已被□定
vnode_label	vp 的策略□□	

一个□有信任状old的主体由于□行(vp文件而□致□□□□， □函数根据vnode□□□□主体重新□□□new。 □当一个□程□求□行vnode文件， 而通□ 入口函数`mpo_execve_will_transition` 有成功返回的策略□， 将□用□函数。 策略模□可以通□□入□个主体信任状和□□地□用 `mpo_create_cred` 来□□□入口函数， so as not to implement a transitioning event. 一旦策略□□了`mpo_create_cred`函数， 即使没有□□ `mpo_execve_will_transition`， 也□□□□函数。

#### 6.7.3.4.3. mpo\_execve\_will\_transition

```
int  
mpo_execve_will_transition (struct ucred,  
                           struct vnode,  
                           struct label);
```

参数	□明	□定
old	在□行execve(2)之前的主体信任状	不可改□
vp	将被□行的文件	
vnode_label	vp 的策略□□	

由策略决定，当参数主体信任状□行参数 vnode □，是否需要□行一个□□□操作。如果需要，返回1；否□，返回0。即使一个策略返回0，它也必□□自己不期望的□ mpo\_execve\_transition的□用作好准□，因□只要有其他任何一个策略要求□□，就将□行此函数。

#### 6.7.3.4.4. mpo\_create\_proc0

```
void  
mpo_create_proc0 (struct ucred);
```

参数	□明	□定
cred	将被填写的主体信任状	

□程0，所有内核□程的祖先，□建主体信任状。

#### 6.7.3.4.5. mpo\_create\_proc1

```
void  
mpo_create_proc1 (struct ucred);
```

参数	□明	□定
cred	将被填写的主体信任状	

□程1，所有用□程的祖先，□建主体信任状。

#### 6.7.3.4.6. mpo\_relabel\_cred

```
void  
mpo_relabel_cred (struct ucred,  
                   struct label);
```

参数	□明	□定
cred	主体信任状	
newlabel	将被用到 cred 上的新□□	

根据□入的新□□， □主体信任状上的□□□行更新。

#### 6.7.4. □□控制□□

通□□□控制入口函数，策略模□能影□内核的□□控制决策。

□□控制入口函数的参数有，一个或者若干个授□信任状，和相□操作□及的其他任何□象的信息（其中可能包含□□）。 □□控制入口函数返回0，表示允□□操作；否□，返回一个 `errno(2)` □□□□。 □用□入口函数，将遍□所有系□注册的策略模□，逐一□行 策略相□的□□□和决策，之后按照下述方法□合不同策略的返回□□。 只有当所有的模□均允□□操作□，才成功返回。 否□，如果有□或者若干模□失□返回，□整个□□不通□。 如果有多个模□的□□出□返回，将由定□在kern\_mac.c 中的 `error_select()` 函数从它□返回的□□□□中，□□一个合□的，返回□用□。

最高□先□	EDEADLK
	EINVAL
	ESRCH
	EACCES
最低□先□	EPERM

如果所有策略模□返回的□□□□均没有出□在上述□先□序列表中，□任意□□一个返回。  
：内核□□，无效的参数，□象不存在，□□被拒□，和其他□□。

##### 6.7.4.1. `mpo_check_bpfdesc_receive`

```
int
mpo_check_bpfdesc_receive (struct bpf_d,
                           struct label,
                           struct ifnet,
                           struct label);
```

参数	□明	□定
bpf_d	主体；BPF 描述子	
bpflabel	bpf_d 的策略□□	
ifnet	客体；网□接口	
ifnetlabel	ifnet 的策略□□	

决定 MAC 框架是否□□允□将由参数接口接收到的数据□□□□由 BPF 描述子所□□的□冲区。成功，□返回0；否□，返回□□□□信息`errno`。建□使用的□□□□有：EACCES，用于□□不符的情况；EPERM，用于缺少特□的情况。

#### 6.7.4.2. mpo\_check\_kenv\_dump

```
int  
mpo_check_kenv_dump (struct ucred);
```

参数	□明	□定
cred	主体信任状	

决定相□主体是否□□被允□□内核□境状□（参考 [kenv\(2\)](#)）。

#### 6.7.4.3. mpo\_check\_kenv\_get

```
int  
mpo_check_kenv_get (struct ucred,  
                     char *name);
```

参数	□明	□定
cred	主体信任状	
name	内核的□境□量名字	

决定相□主体是否可以□□内核中□定□境□量的状□。

#### 6.7.4.4. mpo\_check\_kenv\_set

```
int  
mpo_check_kenv_set (struct ucred,  
                     char *name);
```

参数	□明	□定
cred	主体信任状	
name	内核的□境□量名字	

决定相□主体是否有□□置□定内核□境□量的□。

#### 6.7.4.5. mpo\_check\_kenv\_unset

```
int  
mpo_check_kenv_unset (struct ucred,  
                      char *name);
```

参数	□明	□定
cred	主体信任状	

参数	□明 □定
name	内核的环境变量名字Kernel environment variable name

决定相□主体是否有□清除□定的内核环境变量的□置。

#### 6.7.4.6. `mpo_check_kld_load`

```
int
mpo_check_kld_load (struct ucred,
                     struct vnode,
                     struct label);
```

参数	□明 □定
cred	主体信任状
vp	内核模块的 vnode
vlabel	vp的策略□

决定相□主体是否有□加载□定的模块文件。

#### 6.7.4.7. `mpo_check_kld_stat`

```
int
mpo_check_kld_stat (struct ucred);
```

参数	□明 □定
cred	主体信任状

决定相□主体是否有□内核的加载模块文件□表以及相□的□数据。

#### 6.7.4.8. `mpo_check_kld_unload`

```
int
mpo_check_kld_unload (struct ucred);
```

参数	□明 □定
cred	主体信任状

决定相□主体是否有□卸□一个内核模块。

#### 6.7.4.9. `mpo_check_pipe_ioctl`

```

int
mpo_check_pipe_ioctl (struct ucred,
                      struct pipe,
                      struct label,
                      unsigned long,
                      void *data);

```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe的策略□□	
cmd	ioctl(2) 命令	
data	ioctl(2) 数据	

决定相□主体是否有□□用指定的 [ioctl\(2\)](#) 系□□用。

#### 6.7.4.10. mpo\_check\_pipe\_poll

```

int
mpo_check_pipe_poll (struct ucred,
                      struct pipe,
                      struct label);

```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe的策略□□	

决定相□主体是否有□□管道pipe□行poll操作。

#### 6.7.4.11. mpo\_check\_pipe\_read

```

int
mpo_check_pipe_read (struct ucred,
                      struct pipe,
                      struct label);

```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe的策略□□	

决定主体是否有权取pipe。

#### 6.7.4.12. mpo\_check\_pipe\_relabel

```
int  
mpo_check_pipe_relabel (struct ucred,  
                         struct pipe,  
                         struct label,  
                         struct label);
```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe的当前策略□□	
newlabel	将□pipelabel□置的新□□	

决定主体是否有权pipe重新□置□□。

#### 6.7.4.13. mpo\_check\_pipe\_stat

```
int  
mpo_check_pipe_stat (struct ucred,  
                      struct pipe,  
                      struct label);
```

参数	□明	□定
cred	主体信任状	
pipe	管道	
pipelabel	pipe的策略□□	

决定主体是否有权与pipe相□的□□信息。

#### 6.7.4.14. mpo\_check\_pipe\_write

```
int  
mpo_check_pipe_write (struct ucred,  
                      struct pipe,  
                      struct label);
```

参数	□明	□定
cred	主体信任状	
pipe	管道	

参数	□明	□定
pipelabel	pipe的策略□□	

决定□主体是否有□写pipe。

#### 6.7.4.15. `mpo_check_socket_bind`

```
int
mpo_check_socket_bind (struct ucred,
                      struct socket,
                      struct label,
                      struct sockaddr);
```

参数	□明	□定
cred	主体信任状	
socket	将被□定的套接字	
socketlabel	socket的策略□□	
sockaddr	socket的地址	

#### 6.7.4.16. `mpo_check_socket_connect`

```
int
mpo_check_socket_connect (struct ucred,
                          struct socket,
                          struct label,
                          struct sockaddr);
```

参数	□明	□定
cred	主体信任状	
socket	将被□接的套接字	
socketlabel	socket的策略□□	
sockaddr	socket的地址	

决定□主体（cred）是否有□将套接字（socket）□定到地址 sockaddr。成功，返回0，否□返回一个□□□□`errno`。建□采用的□□□□有：EACCES，用于□□不符的情况；EPERM，用于特□不足的情况。

#### 6.7.4.17. `mpo_check_socket_receive`

```
int
mpo_check_socket_receive (struct ucred,
                           struct socket,
                           struct label);
```

参数	□明	□定
cred	主体信任状	
so	套接字	
socketlabel	so的策略□□	

决定□主体是否有□□□套接字so的相□信息。

#### 6.7.4.18. `mpo_check_socket_send`

```
int
mpo_check_socket_send (struct ucred,
                      struct socket,
                      struct label);
```

参数	□明	□定
cred	主体信任状	
so	套接字	
socketlabel	so的策略□□	

决定□主体是否有□通□套接字so□送信息。

#### 6.7.4.19. `mpo_check_cred_visible`

```
int
mpo_check_cred_visible (struct ucred,
                        struct ucred);
```

参数	□明	□定
u1	主体信任状	
u2	□象信任状	

□定□主体信任状u1是否有□ "see" 具有信任状u2 的其他主体。 成功，返回0；否□，返回□□□□`errno`。建□采用的□□□□有： EACCES，用于□□不符的情况；EPERM，用于特□不足的情况；ESRCH， 用来提供不可□性。□函数可在□多□境下使用，包括命令ps所使用的□程□的状□ sysctl，以及通□procfs 的状□□□操作。

#### 6.7.4.20. `mpo_check_socket_visible`

```
int
mpo_check_socket_visible (struct ucred,
                          struct socket,
                          struct label);
```

参数	□明	□定
cred	主体信任状	
socket	客体；套接字	
socketlabel	socket的策略□□	

#### 6.7.4.21. `mpo_check_ifnet_relabel`

```
int
mpo_check_ifnet_relabel (struct ucred,
                         struct ifnet,
                         struct label,
                         struct label);
```

参数	□明	□定
cred	主体信任状	
ifnet	客体；网□接口	
ifnetlabel	ifnet□有的策略□□	
newlabel	将被□用到ifnet上的新的策略□□	

决定□主体信任状是否有□使用□入的□□更新参数□□定的网□接口的□□□行重新□置。

#### 6.7.4.22. `mpo_check_socket_relabel`

```
int
mpo_check_socket_relabel (struct ucred,
                          struct socket,
                          struct label,
                          struct label);
```

参数	□明	□定
cred	主体信任状	
socket	客体；套接字	
socketlabel	socket□有的策略□□	
newlabel	将被□用到socketlabel上的更新□□	

决定□主体信任状是否有□采用□入的□□□套接字参数的□□□行重新□置。

#### 6.7.4.23. `mpo_check_cred_relabel`

```
int
    mpo_check_cred_relabel (struct ucred,
                            struct label);
```

参数	□明	□定
cred	主体信任状	
newlabel	将被□用到cred上的更新□□	

决定□主体信任状是否有□将自己的□□重新□置□□定的更新□□。

#### 6.7.4.24. mpo\_check\_vnode\_relabel

```
int
    mpo_check_vnode_relabel (struct ucred,
                            struct vnode,
                            struct label,
                            struct label);
```

参数	□明	□定
cred	主体信任状	不可改□
vp	客体；vnode	已被□定
vnode_label	vp□有的策略□□	
newlabel	将被□用到vp上的策略□□	

决定□主体信任状是否有□将参数 vnode 的□□重新□置□□定□□。

#### 6.7.4.25. mpo\_check\_mount\_stat

```
int mpo_check_mount_stat (struct ucred,
                        struct mount,
                        struct label);
```

参数	□明	□定
cred	主体信任状	
mp	客体；文件系□挂□	
mountlabel	mp的策略□□	

□定相□主体信任状是否有□看在□定文件系□上□行 statfs 的□果。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。 □函数可能在下列情况下被□用： 在 statfs(2) 和其他相□□用期□， 或者当需要从文件系□列表中□□排除□个文件系□□， 比如， □用 getfsstat(2)□。

#### 6.7.4.26. mpo\_check\_proc\_debug

```
int  
mpo_check_proc_debug (struct ucred,  
                      struct proc);
```

参数	□明	□定
cred	主体信任状	不可改□
proc	客体；□程	

□定相□主体信任状是否有□ debug □定□程。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□； ESRCH， 用于□□目□的存在。 [ptrace\(2\)](#) 和 [ktrace\(2\)](#) API， 以及某些 procfs 操作将□用□函数。

#### 6.7.4.27. mpo\_check\_vnode\_access

```
int  
mpo_check_vnode_access (struct ucred,  
                        struct vnode,  
                        struct label,  
                        int flags);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
flags	<a href="#">access(2)</a> □志	

根据相□主体信任状决定其□□定 vnode 以□定□□□志□行的 [access\(2\)](#) 和其他相□□用的返回□。一般， □采用与 [mpo\\_check\\_vnode\\_open](#) 相同的□□来□□□函数。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.28. mpo\_check\_vnode\_chdir

```
int  
mpo_check_vnode_chdir (struct ucred,  
                       struct vnode,  
                       struct label);
```

参数	□明	□定
cred	主体信任状	
dvp	客体； <a href="#">chdir(2)</a> 的目的 vnode	

参数	□明	□定
dlabel	dvp的策略□□	

□定相□主体信任状是否有□将□程工作目□切□到□定 vnode。成功， □返回 0； 否□， 返回一个 errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.29. mpo\_check\_vnode\_chroot

```
int
mpo_check_vnode_chroot (struct ucred,
                        struct vnode,
                        struct label);
```

参数	□明	□定
cred	主体信任状	
dvp	目□ vnode	
dlabel	与dvp相□□的策略□□	

□定相□主体是否有□ chroot(2) 到由 (dvp)□定的目□。

#### 6.7.4.30. mpo\_check\_vnode\_create

```
int
mpo_check_vnode_create (struct ucred,
                        struct vnode,
                        struct label,
                        struct componentname,
                        struct vattr);
```

参数	□明	□定
cred	主体信任状	
dvp	客体；vnode	
dlabel	dvp的策略□□	
cnp	dvp中的成□名	
vap	vap的 vnode 属性	

□定相□主体信任状是否有□在□定父目□，以□定的名字和属性， 常□一个 vnode。成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES 来表示用于□□不匹配， 而用 EPERM， 用于□限不足。 以O\_CREAT□参数□用 open(2)， 或□ mknod(2), mknod(2) 等的□用将□致□函数被□用。

#### 6.7.4.31. mpo\_check\_vnode\_delete

```

int
    mpo_check_vnode_delete (struct ucred,
                           struct vnode,
                           struct label,
                           struct vnode,
                           void *label,
                           struct componentname);

```

参数	□明	□定
cred	主体信任状	
dvp	父目□ vnode	
dlabel	dvp的策略□□	
vp	客体；将被□除的 vnode	
label	vp的策略□□	
cnp	vp中的成□名	

□定相□主体信任状是否有□从□定的父目□中， □除□定名字的 vnode。 成功， □返回 0； 否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。 使用 unlink(2) 和 rmdir(2)， 将□致□函数被□用。 提供□入口函数的策略□必□□□一个 mpo\_check\_rename\_to， 用来授□由于重命名操作□致的目□文件的□除。

#### 6.7.4.32. mpo\_check\_vnode\_deleteacl

```

int
    mpo_check_vnode_deleteacl (struct ucred *cred,
                               struct vnode *vp,
                               struct label *label,
                               acl_type_t type);

```

参数	□明	□定
cred	主体信任状	不可改□
vp	客体；vnode	被□定
	label	vp的策略□□
	type	ACL □型

□定相□主体信任状是否有□除□定 vnode 的□定□型的 ACL。 成功， □返回 0； 否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.33. mpo\_check\_vnode\_exec

```

int
mpo_check_vnode_exec (struct ucred,
                      struct vnode,
                      struct label);

```

参数	□明	□定
cred	主体信任状	
vp	客体；将被□行的 vnode	
label	vp的策略□□	

□定相□主体信任状是否有□□行□定 vnode。 □于□行特□的决策与任何瞬□事件的决策是□格分□的。 成功，□返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.34. mpo\_check\_vnode\_getacl

```

int
mpo_check_vnode_getacl (struct ucred,
                        struct vnode,
                        struct label,
                        acl_type_t);

```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
type	ACL □型	

□定相□主体信任状是否有□□□定 vnode 上的□定□型的 ACL。 成功，□返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.35. mpo\_check\_vnode\_getextattr

```

int
mpo_check_vnode_getextattr (struct ucred,
                           struct vnode,
                           struct label,
                           int,
                           const char,
                           struct uio);

```

参数	□明	□定
cred	主体信任状	

参数	□明	□定
vp	客体；vnode	
label	vp的策略□□	
attrnamespace	□展属性名字空□	
name	□展属性名	
uios	I/O □□指□；参□ <a href="#">uios(9)</a>	

□定相□主体信任状是否有□□□□定 vnode 上□定名字空□和名字的□展属性。 使用□展属性□□□□存□的策略模□可能会需要□□些□展属性的操作□行特殊□理。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.36. [mpo\\_check\\_vnode\\_link](#)

```
int
mpo_check_vnode_link (struct ucred,
                      struct vnode,
                      struct label,
                      struct vnode,
                      struct label,
                      struct componentname);
```

参数	□明	□定
cred	主体信任状	
dvp	目□ vnode	
dlabell	与dvp相□□的策略□□	
vp	□接目的 vnode	
label	与vp相□□的策略□□	
cnp	将被□建的□接□□的成□名	

□定相□主体是否有□□参数vp□定的 vnode □建一个由参数cnp□定名字的□接。

#### 6.7.4.37. [mpo\\_check\\_vnode\\_mmap](#)

```
int
mpo_check_vnode_mmap (struct ucred,
                      struct vnode,
                      struct label,
                      int prot);
```

参数	□明	□定
cred	主体信任状	

参数	□明	□定
vp	将被映射的 vnode	
label	与 vp 相关的策略	
prot	mmap 保护 (参见 mmap(2))	

□定相□主体是否有□将□定 vnode vp 以 prot 指定的保护方式□行映射.

#### 6.7.4.38. mpo\_check\_vnode\_mmap\_downgrade

```
void
mpo_check_vnode_mmap_downgrade (struct ucred,
                                 struct vnode,
                                 struct label,
                                 int *prot);
```

参数	□明	□定
cred	See mpo_check_vnode_mmap.	
vp		label prot

根据主体和客体□□, 降低 mmap protections。

#### 6.7.4.39. mpo\_check\_vnode\_mprotect

```
int
mpo_check_vnode_mprotect (struct ucred,
                           struct vnode,
                           struct label,
                           int prot);
```

参数	□明	□定
cred	主体信任状	
vp	映射的 vnode	
prot	存□保□	

□定相□主体是否有□将□定 vnode vp 映射内存空□的存□保□参数□置□指定□。

#### 6.7.4.40. mpo\_check\_vnode\_poll

```

int
mpo_check_vnode_poll (struct ucred,
                      struct ucred,
                      struct vnode,
                      struct label);

```

参数	□明	□定
active_cred	主体信任状	
file_cred	与struct file相□的的信任状	
vp	将被□行 poll 操作的 vnode	
label	与vp相□的策略□	

□定相□主体是否有□□□定 vnode vp□行 poll 操作。

#### 6.7.4.41. mpo\_check\_vnode\_rename\_from

```

int
mpo_vnode_rename_from (struct ucred,
                       struct vnode,
                       struct label,
                       struct vnode,
                       struct label,
                       struct componentname);

```

参数	□明	□定
cred	主体信任状	
dvp	目□ vnode	
dlabel	与dvp相□的策略□	
vp	将被重命名的 vnode	
label	与vp相□的策略□	
cnp	vp中的成□名	

□定相□主体是否有□重命名□定vnode, vp。

#### 6.7.4.42. mpo\_check\_vnode\_rename\_to

```

int
mpo_check_vnode_rename_to (struct ucred,
                           struct vnode,
                           struct label,
                           struct vnode,
                           struct label,
                           int samedir,
                           struct componentname);

```

参数	□明	□定
cred	主体信任状	
dvp	目□ vnode	
dlabel	与dvp相□的策略□	
vp	被覆□的 vnode	
label	与vp相□的策略□	
samedir	布□型□量；如果源和目的目□是相 同的， □被置□1	
cnp	目□component名	

□定相□主体是否有□重命名□定 vnode vp, 至指定目□ dvp, 或更名□cnp。如果无需覆□已有文件， □vp 和 label 的□将□ NULL.

#### 6.7.4.43. mpo\_check\_socket\_listen

```

int
mpo_check_socket_listen (struct ucred,
                        struct socket,
                        struct label);

```

参数	□明	□定
cred	主体信任状	
socket	客体；套接字	
socketlabel	socket的策略□	

□定相□主体是否有□听□定套接字。 成功， □返回0；否□， 返回□□□□□errno。 建□使用的□□□□： EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.44. mpo\_check\_vnode\_lookup

```

int
mpo_check_vnode_lookup (struct ucred,
                        struct vnode,
                        struct label,
                        struct componentname);

```

参数	□明	□定
cred	主体信任状	
dvp	客体；vnode	
dlabel	dvp的策略□□	
cnp	被□□的成□名	

□定相□主体信任状是否有□在□定的目□ vnode 中□□□定名字□行lookup操作。 成功， □返回 0 ； 否□， 返回一个 errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.45. mpo\_check\_vnode\_open

```

int
mpo_check_vnode_open (struct ucred,
                      struct vnode,
                      struct label,
                      int);

```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
acc_mode	open(2) □□模式	

□定相□主体信任状是否有□在□定 vnode 上以□定的□□模式□行 open 操作。 如果成功， □返回 0 ； 否□， 返回一个□□□□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.46. mpo\_check\_vnode\_readdir

```

int
mpo_check_vnode_readdir (struct ucred,
                         struct vnode,
                         struct label);

```

参数	□明	□定
cred	主体信任状	

参数	□明	□定
dvp	客体；目□ vnode	
dlabel	dvp的策略□□	

□定相□主体信任状是否有□在□定的目□ vnode 上□行 `readdir` 操作。成功，□返回 0；否□，返回一个□□□□ `errno`。建□使用的□□□□：EACCES，用于□□不匹配；EPERM，用于□限不□。

#### 6.7.4.47. `mpo_check_vnode_readlink`

```
int
    mpo_check_vnode_readlink (struct ucred,
                              struct vnode,
                              struct label);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	

□定相□主体信任状是否有□在□定符号□接 vnode 上□行 `readlink` 操作。成功，□返回 0；否□，返回一个 `errno`。建□使用的□□□□：EACCES，用于□□不匹配；EPERM，用于□限不□。□函数可能在若干□境下被□用，包括由用□□程□式□行的 `readlink` □用，或者是在□程□行名字□□□□式□行的 `readlink` 。

#### 6.7.4.48. `mpo_check_vnode_revoke`

```
int
    mpo_check_vnode_revoke (struct ucred,
                              struct vnode,
                              struct label);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	

□定相□主体信任状是否有□撤□□定 vnode 的□□。成功，□返回 0；否□，返回一个 `errno`。建□使用的□□□□：EACCES，用于□□不匹配；EPERM，用于□限不□。

#### 6.7.4.49. `mpo_check_vnode_setacl`

```

int
mpo_check_vnode_setacl (struct ucred,
                        struct vnode,
                        struct label,
                        acl_type_t,
                        struct acl);

```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
type	ACL □型	
acl	ACL	

□定相□主体信任状是否有□□置□定 vnode 的□定□型的 ACL。 成功，□返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配；EPERM， 用于□限不□。

#### 6.7.4.50. mpo\_check\_vnode\_setextattr

```

int
mpo_check_vnode_setextattr (struct ucred,
                            struct vnode,
                            struct label,
                            int,
                            const char,
                            struct uio);

```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
attrnamespace	□展属性名字空□	
name	□展属性名	
uiο	I/O □□指□；参□ <a href="#">uio(9)</a>	

□定相□主体信任状是否有□□置□定 vnode 上□定名字空□中□定名字的□展属性的□。 使用□展属性□□安全□□的策略模□可能需要□其使用的属性□施□外的保□。 □外， 由于在□□和□□操作□□可能存在的□争， 策略模□□□避免根据来自uiο中的数据做出决策。 如果正在□行一个□除操作， □参数 uiο 的□也可能□ NULL。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配；EPERM， 用于□限不□。

#### 6.7.4.51. `mpo_check_vnode_setflags`

```
int  
    mpo_check_vnode_setflags (struct ucred,  
                                struct vnode,  
                                struct label,  
                                u_long flags);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
flags	文件□志；参□ <a href="#">chflags(2)</a>	

□定相□主体信任状是否有□□定的 vnode □置□定的□志。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.52. `mpo_check_vnode_setmode`

```
int  
    mpo_check_vnode_setmode (struct ucred,  
                                struct vnode,  
                                struct label,  
                                mode_t mode);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
mode	文件模式；参□ <a href="#">chmod(2)</a>	

□定相□主体信任状是否有□将□定 vnode 的模式□置□□定□。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.53. `mpo_check_vnode_setowner`

```
int  
    mpo_check_vnode_setowner (struct ucred,  
                                struct vnode,  
                                struct label,  
                                uid_t uid,  
                                gid_t gid);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	
uid	用□ID	
gid	□ID	

□定相□主体信任状是否有□将□定 vnode 的文件 uid 和文件 gid □置□□定□。如果无需更新， 相□参数□可能被□置□(-1)。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.54. mpo\_check\_vnode\_setutimes

```
int
mpo_check_vnode_setutimes (struct ucred,
                           struct vnode,
                           struct label,
                           struct timespec,
                           struct timespec);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vp	
label	vp的策略□□	
atime	□□□□；参□ utimes(2)	
mtime	修改□□；参□ utimes(2)	

□定相□主体信任状是否有□将□定 vnode 的□□□□□□□置□□定□。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.55. mpo\_check\_proc\_sched

```
int
mpo_check_proc_sched (struct ucred,
                      struct proc);
```

参数	□明	□定
cred	主体信任状	
proc	客体；□程	

□定相□主体信任状是否有□改□□定□程的□度参数。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□□：EACCES， 用于□□不匹配； EPERM， 用于□限不□； ESRCH， 用于提供不可□性□。

See [setpriority\(2\)](#) for more information.

#### 6.7.4.56. `mpo_check_proc_signal`

```
int
    mpo_check_proc_signal (struct ucred,
                           struct proc,
                           int signal);
```

参数	□明	□定
cred	主体信任状	
proc	客体；□程	
signal	信号；参□ <a href="#">kill(2)</a>	

□定相□主体信任状是否有□向□定□程□送□定信号。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□：EACCES， 用于□□不匹配；EPERM， 用于□限不□；ESRCH， 用于提供不可□性□。

#### 6.7.4.57. `mpo_check_vnode_stat`

```
int
    mpo_check_vnode_stat (struct ucred,
                           struct vnode,
                           struct label);
```

参数	□明	□定
cred	主体信任状	
vp	客体；vnode	
label	vp的策略□□	

□定相□主体信任状是否有□在□定 vnode 上□行 stat 操作。 成功， □返回 0；否□， 返回一个errno□。 建□使用的□□□：EACCES， 用于□□不匹配；EPERM， 用于□限不□。

See [stat\(2\)](#) for more information.

#### 6.7.4.58. `mpo_check_ifnet_transmit`

```
int
    mpo_check_ifnet_transmit (struct ucred,
                               struct ifnet,
                               struct label,
                               struct mbuf,
                               struct label);
```

参数	□明	□定
cred	主体信任状	
ifnet	网口接口	
ifnetlabel	ifnet的策略□□	
mbuf	客体；将被□送的 mbuf	
mbuflabel	mbuf的策略□□	

□定相□网□接口是否有□□送□定的 mbuf。成功， □返回 0； 否□， 返回一个errno□。 建□使用的□□□□： EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.59. `mpo_check_socket_deliver`

```
int
mpo_check_socket_deliver (struct ucred,
                           struct ifnet,
                           struct label,
                           struct mbuf,
                           struct label);
```

参数	□明	□定
cred	主体信任状	
ifnet	网口接口	
ifnetlabel	ifnet的策略□□	
mbuf	客体；将被□送的 mbuf	
mbuflabel	mbuf的策略□□	

□定相□套接字是否有□从□定的 mbuf 中接收数据□。 成功， □返回 0； 否□， 返回一个errno□。 建□使用的□□□□： EACCES， 用于□□不匹配； EPERM， 用于□限不□。

#### 6.7.4.60. `mpo_check_socket_visible`

```
int
mpo_check_socket_visible (struct ucred,
                           struct socket,
                           struct label);
```

参数	□明	□定
cred	主体信任状	不可改□
so	客体；套接字	
socketlabel	so的策略□□	

□定相□主体信任状cred 是否有□使用系□控函数，比如，由[netstat\(8\)](#) 和 [sockstat\(1\)](#)使用的程序来□察□定的套接字(socket)。成功，□返回 0；否□，返回一个errno□。建□使用的□□□：EACCES，用于□□不匹配；EPERM，用于□限不□；ESRCH，用于提供不可□性□。

#### 6.7.4.61. [mpo\\_check\\_system\\_acct](#)

```
int
    mpo_check_system_acct (struct ucred,
                           struct vnode,
                           struct label);
```

参数	□明	□定
ucred	主体信任状	
vp	□□文件； <a href="#">acct(5)</a>	
vlabel	与vp相□□的□□	

根据主体□□和□□日志文件的□□，□定□主体是否有□□□□□。

#### 6.7.4.62. [mpo\\_check\\_system\\_nfsd](#)

```
int
    mpo_check_system_nfsd (struct ucred);
```

参数	□明	□定
cred	主体信任状	

□定相□主体是否有□□用 [nfssvc\(2\)](#)。

#### 6.7.4.63. [mpo\\_check\\_system\\_reboot](#)

```
int
    mpo_check_system_reboot (struct ucred,
                             int howto);
```

参数	□明	□定
cred	主体信任状	
howto	来自 <a href="#">reboot(2)</a> 的howto 参数	

□定相□主体是否有□以指定方式重□系□。

#### 6.7.4.64. [mpo\\_check\\_system\\_settime](#)

```
int
mpo_check_system_settime (struct ucred);
```

参数	□明	□定
cred	主体信任状	

□定相□用□是否有□置系□□。

#### 6.7.4.65. mpo\_check\_system\_swapon

```
int
mpo_check_system_swapon (struct ucred,
                         struct vnode,
                         struct label);
```

参数	□明	□定
cred	主体信任状	
vp	swap□□	
vlabel	与vp相□□的□□	

□定相□主体是否有□加一个作□swap□□的 vp 。

#### 6.7.4.66. mpo\_check\_system\_sysctl

```
int
mpo_check_system_sysctl (struct ucred,
                        int *name,
                        u_int *namelen,
                        void *old,
                        size_t,
                        int inkernel,
                        void *new,
                        size_t newlen);
```

参数	□明	□定
cred	主体信任状	
name	参□ <a href="#">sysctl(3)</a>	
namelen		
old		
oldlenp		

参数	说明	约定
inkernel	布尔量；如果从内核被引用，其值被置为1	
new	参见 <a href="#">sysctl(3)</a>	

约定相主体是否被允许进行指定的 [sysctl\(3\)](#) 事。

### 6.7.5. 管理使用

当用程序请求某个对象的修改时，将引起重新事件。更新操作分两步进行：首先，执行控制命令，此次更新操作是有效且被允的；然后，用一个独立的入口函数进行修改。重新入口函数通常接收由程序提交的对象、对象指针和更新对象，作为输入参数。对象重新操作的失败将由先期的警告，所以，不允许在接下来的修改过程中失败，故而不提倡在此进程中新分配内存。

## 6.8. 使用体系

TrustedBSD MAC 框架包含了一策略无关的组成元素，包括管理抽象的 MAC 接口，系统信任管理体系的修改，用分配 MAC 提供支持的 login 函数，以及若干管理和更新内核对象（程序、文件和网接口等）安全的工具。不久，将有更多关于使用体系的更多信息被包含进来。

### 6.8.1. 策略无关的管理 API

TrustedBSD MAC 提供的大量函数和系公用，允许应用程序使用一一一的、策略无关的接口来管理对象的 MAC。如此，应用程序可以轻松管理各策略的，无需附加某个特定策略的支持而重新。许多通用工具，比如 [ifconfig\(8\)](#), [ls\(1\)](#) 和 [ps\(1\)](#)，使用一些策略无关的接口网、文件和程序的。这些 API 也被用于支持 MAC 管理工具，比如，[getfmac\(8\)](#), [getpmac\(8\)](#), [setfmac\(8\)](#), [setfsmac\(8\)](#), 和 [setpmac\(8\)](#)。MAC API 的可参考 [mac\(3\)](#).

应用程序管理的 MAC 有存在形式：内部形式，用来返回和设置程序和对象的（[mac\\_t](#)）；基于 C 字符串的外部形式，作在配置文件中的存放形式，用于向用示或者由用入。一个 MAC 由一个元素组成，其中个元素是一个形如（名字，）的二元。内核中的个策略模块被指定一个特定的名字，由它中与名字的采用其策略特有的方式行解析。采用外部形式表示的，其元素表示名字 / ，元素之间以逗号分隔。应用程序可以使用 MAC 框架提供的 API 将一个安全在内部形式和文本形式之行。当向内核某个对象的安全，内部形式的必须所需的元素集合作好内部存准备。此，通常采用下面方式之一：使用 [mac\\_prepare\(3\)](#) 和一个包含所需元素的任意列表；或者，使用从 [mac.conf\(5\)](#) 配置文件中加缺省元素集合的某个系用。在对象置缺省，将允许应用程序在不定系是否采用相策略的情况下，也能向用返回与对象相的有意的安全。



目前的 MAC 不支持直接修改内部形式的元素，所有的修改必须按照下列的行：将内部形式的成文本字符串，字符串行，最后将其成内部形式。如果应用程序的作者明有需要，可以在将来的版本中加入内部形式行直接修改的接口。

### 6.8.2. 使用指定

用上下文管理的接口，[setusercontext\(3\)](#)，的行已被修改，从 [login.conf\(5\)](#) 中与某个用户登录相的 MAC 安全。当 [LOGIN\\_SETALL](#) 被置，或者当 [LOGIN\\_SETMAC](#) 被明指定，一些安全将和其他用

与上下文参数一起被置。



可以期，在今后的某个版本中，FreeBSD 将把 MAC 从 login.conf 的用数据中抽出，其一个独立的数据。不在此前后，[setusercontext\(3\)](#) API 保持不。

## 6.9. 小

TrustedBSD MAC 框架使得内核模能以一集中的方式，完善系的安全策略。它既可利用有的内核象属性，又能使用由 MAC 框架助的安全数据，来施控制。框架提供的活性使得人可以在其上各策略，如利用 BSD 有的信任状（credential）与文件保护机制的策略，以及信息流安全策略（如 MLS 和 Biba）。革新安全服的策略程人，可以参考本文，以了解有安全模的信息。

# Chapter 7. 虚内存系

## 7.1. 物理内存的管理-`vm_page_t`

物理内存通常体`vm_page_t`以基行管理。物理内存的由它各自体`vm_page_t`所代表，些体存放在若干个管理列中的一个里面。

一可以于在(wired)、活(active)、去活(inactive)、存(cache)、自由(free)状态。除了在状态，一般被放置在一个双向表列里，代表了它所的状。在不放置在任何列里。

FreeBSD内存和自由了一个更的列机制，以的分管理。一状态都着多个列，列的安排着处理器的一、二存。当需要分配一个新，FreeBSD会把一个按一、二存的面分配虚内存象。

此外，一个可以有一个引用数，可以被一个忙数定。虚内存系也了"锁定"(ultimate locked)状态，一个可以用标志`PG_BUSY`表示一状态。

之，个列都按照LRU(Least-Recently Used)的原则工作。

者注



短Least-Recently Used有理解方式：1.将"least-recently"理解反向比，意"最早"，整个短理解“最近的使用最早”；2.将"least"和"recently"理解副，都修"used"，整个短理解"最近最少使用"。理解方式的意基本相同。

一个常常最初于在或活状态。在，常常于某的表。虚内存系通常描在一个活的列(LRU)定的年，以便将他移到一个不活的列中。移到存中的依然与一个VM象，但被作立即再用的候。在自由列中的是真正未被使用的。FreeBSD尽量不将放在自由列中，但是必保持一定数量的自由，以便中断分配。

如果一个程不在表中而在某一列中的（例如去活列或存列），一个相耗源少的产生，致被重激活。如果根本不存在于系内存之中，程必被阻塞，此被从磁中入。

者注



Intel等厂商的CPU工作在保模式，可用来虚内存。当址的地址空着真内存，正常写；当址的地址空没有真的真内存，CPU会生一个"”，通知操作系统与磁等行交，址写入存内容，写址写出存内容。一个"”并非操作系统或应用程序人犯下的，尽管在CPU硬件中与应用程序或操作系统内核崩的产生机制相同。参Intel的CPU保模式手册。

FreeBSD的整列，将各个列中的数在一个适当的比例上，同管理程序崩的已清理和未清理。重新平衡的比例数决定于系内存的负担。重新平衡由pageout守程，包括清理未清理（与他的后存同）、被引用的活程度（重置它在LRU列中的位置或在不同活程度的列移）、当比例不平衡在列移，如此等等。FreeBSD的VM系会将重激活而生的率降低到一个合理的数，由此决定某一活/置的度。可以更好的决定何清理/分配一个做出决策。

## 7.2. 一的存信息体-`vm_object_t`

FreeBSD一的"虚内存象"(VM象)的思想。

直接使用(unbacked)、 交(swap)、 物理、 文件。 由于文件系使用相同的VM象管理核内数据-文件的存， 所以些存的也是唯一的。

VM象可以被影制(shadowed)。

它可以被堆放到其它VM象堆的端。例如，可以有一个交VM象， 放置在文件VM象堆的端，以MAP\_PRIVATE的mmap()操作。 的入操作也可以用来各各的共享特性，包括写入制(copy-on-write，用于日志文件系)，以派生出地址空。

当注意，一个`vm_page_t`

体在任一个时刻只能与一个VM象相。

VM象影本可以跨

例的共享相同的。

## 7.3. 文件系入/出-`buf`体

vnode VM象，比如文件VM象，一般需要它自己的清理(clean)/未清理(dirty)信息，而不依赖于文件系的清理/未清理。 例如，当VM系要同一个物理和其存器，VM系就需要在写入到存器前将已清理。 另外，文件系要能将文件或文件元数据的各部分映射到内核虚内存(KVM)中以便操作。

用行一些管理的体就是所周知的文件系存，`struct buf`或`bp`。 当文件系需要一个VM象的一部分操作，它常会将象的部分映射到`struct buf`，并且将`struct buf`中映射到内核虚内存(KVM)中。 同的，磁入/出通常要先将VM象的各部分映射到buf体中，然后buf体行/入/出操作。 下的`vm_page_t`在入/出期通常被"忙"。 文件系存也会"忙"，于文件系程序非常有用，文件系存操作比VM真(hard)操作更好。

FreeBSD保留一定数量的内核虚内存来存放`struct` 体， 些内核虚内存用来存放映射，

`vm_page_t`的一个功能，不是文件系存的功能。

固有的限制了同行入/出可能的数量。

并不会造成。

buf的映射，但是些buf体是被清理的。

并不限制存数据的能力。格的，物理数据存是然而，由于文件系存被用来理入/出，他但是，由于通常有数千文件系存可供使用，所以并不会造成。

## 7.4. 映射表-`vm_map_t`, `vm_entry_t`

FreeBSD将物理表从VM系中分了出来。各程的所有表可以脱程(on the fly)重建，并且通常被是一次性的。特殊的表，如内核虚内存(KVM)，常常是被永久性分配的；些表不是一次性的。

FreeBSD通`vm_map_t`和`vm_entry_t` 将虚内存中`vm_objects`的各地址部分起来。 表被直接的从`vm_map_t/vm_entry_t/vm_object_t` 中有次的合成出来。里需要重申一下，我曾提到的"物理直接与`vm_object`相并不很正。 `vm_page_t` 也被会被接到正在与之相的表中。当表被用，一个`vm_page_t`体可以被接到几个`pmaps`。 然而，由于有了次的，因此在象中所有同一的引用会引用同一`vm_page_t`体，就跨区域(board)的存的唯一。

## 7.5. KVM存映射

FreeBSD使用KVM存放各各的内核体。在KVM中最大的个体是文件系存。那是与`struct buf`体有映射。

不像Linux, FreeBSD不将所有的物理内存映射到KVM中。这意味着FreeBSD可以在32位平台上管理超过4GB的内存配置。事实上，如果mmu(译者注：可能是指“**内存管理单元**”，“Memory Management Unit”)有足够的能力，FreeBSD理论上可以在32位平台上管理最多8TB的内存配置。然而，大多数32平台只能映射4GB内存，也只能是一个争执点。

有几项机制可以管理KVM。管理KVM的主要机制是区域分配器(zone allocator)。区域分配器管理着KVM的大块，再将大块切分成恒定大小的小块，以便按照某一类型的对象分配。可以使用命令`vmstat -m` —当前KVM分区使用情况。

## 7.6. 整理FreeBSD的虚拟内存系统

作者的共同努力使得FreeBSD可以自行整理内核。一般来说，不需要做任何混乱的事情。一些内核配置项(一般)被指定在`/usr/src/sys/i386/conf/CONFIG_FILE`之中。所有可用内核配置项的描述可在`/usr/src/sys/i386/conf/LINT`中找到。

在一个大的配置中，可能需要增加`maxusers`的值。数值通常在10到128。注意，过度增加`maxusers`的值可能导致系统从可用的KVM中溢出，从而引起无法感知的操作。最好将`maxusers`设置一个合理的数值，并添加其它项，如`NMBCLUSTERS`，来增加特定的资源。

如果系统的硬盘要被重荷的使用网卡，需要增加`NMBCLUSTERS`的值。数值通常在1024到4096。

`NBUF`也是系统的参数。这个参数决定系统可用来映射文件系统入/出存储的KVM的数量。注意：这个参数与单一的内存没有任何关系。这个参数可在3.0-CURRENT 和以后的内核中被修改，通常不当被手工修改。我推荐不要指定`NBUF`。系统自行决定它。太小的值会致非常低效的文件系统操作；太大的值会使用内存中缺少页面，而大量的位于在状态。

缺省情况下，FreeBSD内核通常是不被优化的。可以在内核配置文件中用`makeoptions`指定排列(debugging)和优化标志。注意，一般不使用`-g`，除非能对付由此产生的大内核(典型的是7MB或更多)。

```
makeoptions      DEBUG="-g"  
makeoptions      COPTFLAGS="-O -pipe"
```

Sysctl提供了在自行整理内核的方式。通常不需要指定任何sysctl变量，尤其是与VM相关的那些变量。

自行整理VM和系统的影子相对直接一些。首先，当尽可能在UFS/FFS文件系统上使用Soft Updates。在`/usr/src/sys/ufs/ffs/README.softupdates`里有关于如何配置的指示。

其次，当配置足够的交换空间。当在每个物理磁盘上配置一个交换分区，最多4个，甚至在大的“工作”磁盘上。当有至少2倍于主内存的交换空间；假如没有足够的内存，交换分区会更多。你也应当按照期望中的最大内存配置决定交换分区的大小，以后就不再需要重新磁盘分区了。如果物理系统崩溃后的内存倒(crash dump)，第一个交换分区必须至少与主内存一样大，/var/crash必须有足够的空间来承载倒。

NFS上的交换分区可以很好的被4.X或后来的系统使用，但是必须明白NFS服务器将要遭受操作很慢的冲击。

# Chapter 8. SMPng 文档

## 8.1. 前言

本文目前 SMPng 架构的介绍与本章行了介绍。它首先介绍了基本的原语和相关工具，其后是关于 FreeBSD 内核的同步与并行模型，接下来介绍了具体体系中的策略，并描述了在各个子系统中引入粒度的同步和并行化的尝试，最后是结论性的说明，用以解释最初做出某些决策的动机，并使读者了解使用特定的原语所可能产生的重大影响。

本文仍在撰写当中，并将不断更新以反映与 SMPng 目前的最新同步与并行的情况。其中有多少目前只是提到了，但我将会逐步充实内容。至于本文的更新和建设，本文将更新。

SMPng 的目标是使内核能并行。基本上，内核是一个很大而复杂的程序。要让内核能多线程地运行，我需要使用某些其它多线程程序在其中所用到的工具，包括互斥体(mutex)、共享/排他锁(shared/exclusive lock)、信号量(semaphores) 和条件变量(condition variable)。如果希望了解它们以及其他 SMP 相关，参见本文的 [附录一](#)。

## 8.2. 基本工具与上层的基础知识

### 8.2.1. 原子操作指令和内存屏障

对于内存屏障和原子操作指令已有很多介绍材料，因此这一节并不打算对其行尽的介绍。总而言之，如果有某一处写入，就不能在不得相容的情况下行取操作。也就是说，内存屏障的作用在于保证内存操作的相容性，但并不保证内存操作的严格顺序。而言之，内存屏障并不保证 CPU 将本地快取内存或存命中内容刷写回内存，而是在存放其所保证的数据，于能看到存放的那个 CPU 或可。持有内存屏障的 CPU 可以在其快取内存或存命中将数据保持其所希望的、任意的值，但如果其它 CPU 在同一数据元上行原子操作，第一个 CPU 必须保证，其所更新的数据，以及内存屏障所要求的任何其它操作，第二个 CPU 可以。

例如，假在一模型中，要在主存(或某一全局快取内存)中的数据是可见的，当某一 CPU 上触原子操作时，其它 CPU 的存命中和快取内存就必须同一快取内存上的全部写操作，以及内存屏障之后的全部未完成操作行刷写。

一来，在使用由原子操作保证的内存单元就需要特别小心。例如，在 sleep mutex 时，我就必须使用 `atomic_cmpset` 而不是 `atomic_set` 来打 MTX\_CONTESTED 位。做的原因是，我需要把 `mtx_lock` 的值到某个量，并据此行决策。然而，我得到的可能是 0 的，也可能在我行决策的进程中产生变化。因此，当行 `atomic_set` 时，最可能会同一行置位，而不是我行决策的那个。就必须通过 `atomic_cmpset` 来保证只有在我行的决策依据是最新的，才相的量行置位。

最后，原子操作只允一次更新或一个内存单元。需要原子地更新多个单元，就必须使用来代替它了。例如，如果需要更新多个相互的计数器，就必须使用，而不是次独的原子操作了。

### 8.2.2. 读与写

并不需要像写那般。不同的，都需要保证它所的不是的数据。然而，只有写操作必须是排他的，而多个线程可以安全地同一量的。使用不同类型的用于读和写操作有多各自不同的方式。

第一种方法是用 sx，它可用于读使用的排他，而作共享。这种方法十分明了。

第二种方法略晦涩。可以用多个锁来保护同一数据元。但是，只需其中一个锁即可。然而，如果要写数据的，则需要首先上所有的写锁。这会大大提高写操作的代价，但当可能以多线程方式访问数据时可能非常有用。例如，父线程指针是同步受 `proclock` 锁和线程 mutex 保护的。在只希望已知线程的父线程时，用 proc 更方便。但是，其它一些地方，例如 `inferior` 需要通过父指针在进程中进行搜索，并且一个进程中不能做了，否则，将无法保证在我所得的结果上操作，之前所有的状况依旧有效。

### 8.2.3. 上锁和后果

如果需要使用锁来保持所共享量的状态，并据此执行某些操作，则是不能在共享量之前上其锁，并在线程操作之前解锁的。过早解锁将使共享量再次可见，可能会导致之前所做的决策失效。因此，在所做修改的操作结束之前，必须保持上锁状态。

## 8.3. 架构与概览

### 8.3.1. 中断的管理

与许多其它多线程 UNIX® 内核所采取的模式类似，FreeBSD 会授予中断处理程序独立的线程上下文，这样做能保证中断线程在遇到阻塞时被唤醒。但为了避免不必要的延时，中断线程在内核中，是以线程的优先级运行的。因此，中断处理程序不运行太久，以免饿死其它内核线程。此外，由于多个处理程序可以分享同一中断线程，中断处理程序不休眠，或使用可能导致休眠的锁，以避免将其它中断处理程序饿死。

目前在 FreeBSD 中的中断线程是指重量级中断线程。之所以称呼它是因为在于，回到中断线程需要执行一次完整的上下文切换操作。在最初的版本中，内核不允许占有，因此中断在打断内核线程之前，必须等待内核线程阻塞或返回可用之后才能运行。

为了解决这个问题，FreeBSD 内核在采用了抢占式调度策略。目前，只有释放休眠 mutex 或发生中断才能打断内核线程，但最明显是在 FreeBSD 上面所描述的全抢占式调度策略。

并非所有的中断处理程序都在独立的线程上下文中运行。相反，某些处理程序会直接在主中断上下文中运行。有些中断处理程序，如在被本地命名的“快速”中断处理程序，因早期版本的内核中使用了 INTR\_FAST 标志来某些处理程序。目前只有中断和串口 I/O 中断采用合一型。由于这些处理程序没有独立的上下文，因而它们都不能获得阻塞性锁，因此也就只能使用自旋 mutex。

最后，有一种称重量级上下文切换的优化，可以在 MD 代码中使用。因为中断线程都是在内核上下文中运行的，所以它可以借用任意线程的 vmspace (虚内存地址空间)。因此，在重量级上下文切换中，切换到中断线程并不切换它的 vmspace，而是借用被中断线程的 vmspace。确保被中断线程的 vmspace 不在中断线程中消失，被中断线程在中断线程不再借用其 vmspace 之前是不允许运行的。刚才提到的情况可能在中断线程阻塞或完成再生。如果中断线程发生阻塞，当它再次进入可运行状态将使用自己的上下文，这样一来，就可以释放被中断的线程了。

这种优化的坏处在于它和硬件紧密相关，而且相比而言，因此只有在可能做能来大幅性能改善才采用。目前可能太早，而且事实上可能会反而导致性能下降，因为几乎所有的中断处理程序都会立即被全局 (Giant) 阻塞，而阻塞将需要线程修正。另外，Mike Smith 提出采用合一方式来管理中断线程：

1. 每个中断处理程序分两部分，一个在主中断上下文中运行的主体 (predicate) 和一个在自己的线程上下文中运行的处理程序 (handler)。
2. 如果中断处理程序有主体，当触发中断时，运行主体。如果主体返回真，则中断被处理完。

内核从中断返回。如果主体返回假，或者中断没有主体，调度执行程式管理程序。

在一模式中当地采用量上下文切换是非常的。因此我可能会希望在未来改一模式，因此在最好的方案，我是会推量上下文切换之上的工作，以便一完善中断处理架构，随后再考察量上下文切换是否可用。

### 8.3.2. 内核抢占与界区

#### 8.3.2.1. 内核抢占简介

内核抢占的概念很，其基本思想是 CPU 调度将要执行的线程并放入待执行队列，看它的优先级是否高于目前正在执行的线程。当然，至少在理想情况下是的。有些时候，造成一理想的代价会十分高昂，以至于在一些情况下抢占会得不偿失。

完全的内核抢占十分：在调度将要执行的线程并放入待执行队列，看它的优先级是否高于目前正在执行的线程。如果是的，执行一次上下文切换并立即开始执行线程。

尽管能抢占多数数据，但内核并不是可以安全地完全抢占的。例如，如果持有自旋 mutex 的线程被抢占，而新线程也需要同一自旋 mutex，新线程就可能一直自旋下去，因为被中断的线程可能永远没有机会执行了。此外，某些代码，例如在 Alpha 上的 exec 线程地址空间号执行的代码也不能被中断，因为它被用来支持的上下文切换操作。在这些代码段中，会通过使用界区来禁用抢占。

#### 8.3.2.2. 界区

界区 API 的任务是避免在界区内发生上下文切换。对于完全抢占式内核而言，除了当前线程之外的其它线程的每个 setrunqueue 都是中断点。`critical_enter` 的一种方式是设置一线程私有标志，并由其所有方清除。如果用 setrunqueue 置了多个标志，无论新线程和当前线程相比其优先级高低，都不会发生抢占。然而，由于界区会在自旋 mutex 中用于避免上下文切换，而且能同时获得多个自旋 mutex，因此界区 API 必须支持嵌套。由于这个原因，目前的 API 中采用了嵌套函数，而不是多个的线程标志。

为了尽可能缩短时间，在界区中的抢占被推到，而不是直接。如果线程被中断，并被置为可执行，而当前线程位于界区，会设置一线程私有标志，表示有一个尚未执行的中断操作。当最外层界区退出时，会重置一个标志，如果它被置位，当前线程会被中断，以允更高优先级的线程开始执行。

中断会引入一个和自旋 mutex 有关的。如果低中断处理程序需要，它就不能中断任何需要的代码，以避免可能产生的坏数据的情况。目前，这一机制是通过界区 API 以 `cpu_critical_enter` 和 `cpu_critical_exit` 函数的形式实现的。目前唯一 API 会在所有 FreeBSD 所支持的平台上禁用和重新启用中断。这种方法并不是最的，但它更易理解，也更容易正确地实现。理论上，这一帮助 API 只需要配合在主中断上下文中的自旋 mutex 使用。然而，为了代码更清晰，它被用在了全部自旋 mutex，甚至包括所有界区上。将其从 MI API 中剥出来放入 MD API，并只在需要使用它的 MI API 的自旋 mutex 中使用可能会有更好的效果。如果我最初采用了这种方式，MD API 可能需要改名，以彰显其独一无二 API 一事。

#### 8.3.2.3. 折衷

如前面提到的，当完全抢占并不能提供最佳性能时，采取了一些折衷的措施。

第一折衷是，抢占并不考虑其它 CPU 的存在。假设我有两个 CPU，A 和 B，其中 A 上线程的优先级是 4，而 B 上线程的优先级是 2。如果 CPU B 令一线程进入可执行状态，理论上，我希望 CPU A 切换至另一新线程，就有个优先级最高的线程在执行了。然而，注定一个 CPU 在抢占更合适，并通过 IPI 向那个 CPU 发出信号，并完成相工作的代价十分高昂。因此，目前的代码会强制 CPU B 切换至更高优先级的线程。

注意做仍会系入更好的状，因 CPU B 会去行先 1 而不是 2 的那个程。

第二折衷是限制于先的内核程的立即占。在前面所定的占操作的的情形中，低先会被立即断（或在其退出界区后被断）。然而，多在内核中行的程，有很多只会行很短的就会阻塞或返回用。因此，如果内核断些程并行其它非的内核程，内核可能会在些程上要休眠或行完之前切出去。一来，CPU 就必整快取存以配合新程的行。当内核返回到被断的程，它又需要重新填充之前失的快取存信息。此外，如果内核能将将阻塞或返回用的那个程的断延到之后的，能免去次外的上下文切。因此，默情况下，只有在先高的程是程，占代才会立即行断操作。

用所有内核程的完全占于非常有帮助，因它会暴露出更多的条件（race conditions）。在以模些条件的处理器系中，得尤其有用。因此，我提供了内核 **FULL\_PREEMPTION** 来用所有内核程的占，一主要用于目的。

### 8.3.3. 程移

地，程从一个 CPU 移到一个上的程称作移。在非占式内核中，只会在明定的点，例如用 **msleep** 或返回至用才会生。但是，在占式内核中，中断可能会在任何时候制断，并致移。于 CPU 私有的数据而言可能会来一些面影，因除 **curthread** 和 **curpcb** 以外的数据都可能在移程中生化。由于存在潜在的程移，使得未受保的 CPU 私有数据显得无用。就需要在某些代码段禁止移，以得一定的 CPU 私有数据。

目前我采用界区来避免移，因它能阻止上下文切。但是，有可能是一于的限制，因界区上会阻止当前处理器上的中断程。因而，提供了另一个 API，用以指示当前程在被断，不移到一 CPU。

API 也叫程制，它由度器提供。API 包括个函数：**sched\_pin** 和 **sched\_unpin**。个函数用于管理程私有的数 **td\_pinned**。如果嵌套数大于零，程将被住，而程始行其嵌套数零，表示于未制状。所有的度器中，都要求保制程只在它首次用 **sched\_pin** 所在的 CPU 上行。由于只有程自己会写嵌套数，而只有其它程在受制程没有行，且持有 **sched\_lock** 才会嵌套数，因此 **td\_pinned** 不必上。**sched\_pin** 函数会使嵌套数，而 **sched\_unpin** 使其。注意，些函数只操作当前程，并将其定到其行它所的 CPU 上。要将任意程定到指定的 CPU 上，使用 **sched\_bind** 和 **sched\_unbind**。

### 8.3.4. 叫出 (Callout)

内核机制 **timeout** 允内核服注册函数，以作 **softclock** 件中断的一部分来行。事件将基于所希望的数目行，并在大指定的回用提供的函数。

未决 **timeout** (超) 事件的全局表是由一全局 mutex，**callout\_lock** 保的；所有 timeout 表的，都必首先拿到个 mutex。当 **softclock** 醒，它会描未决超表，并出的那些。为了避免逆序，**softclock** 程会在用所提供的 **timeout** 回函数首先放 **callout\_lock** mutex。如果在注册没有置 **CALLOUT\_MPSAFE** 志，在用函数之前，会取全局，而在之后放。其后，**callout\_lock** mutex 会在理前再次得。**softclock** 代在放个 mutex 会非常小心地保持表的一致状。如果用了 **DIAGNOSTIC**，个函数的行会被，如果超了某一，会生警告。

## 8.4. 特定数据的策略

### 8.4.1. 凭据

`struct ucred` 是内核内部的凭据对象，它通常作为内核中以编程方式的控制的依据。BSD派生的系统采用一种“写时复制”的模型来管理凭据数据：同一凭据对象可能存在多个引用，如果需要对其进行修改，整个对象将被复制、修改，然后替换引用。由于在打算用于共享控制的凭据快取内存广泛存在，这样做会极大地节省内存。在移到粒度的 SMP 上，这一模型也省去了大量的操作，因为只有未共享的凭据才能施修改，因而避免了在使用共享凭据外的同步操作。

凭据对象只有一个引用，被修改是可行的；不允许改用共享的凭据对象，否则将可能导致生成条件。`cr_mtx` 用于保持 `struct ucred` 的引用计数，以确保其一致性。使用凭据对象时，必须在使用过程中保持有效的引用，否则它就可能在某个不合理的消费者使用过程中被释放。

`struct ucred mutex` 是一个叶 mutex，出于性能考虑，它通过 mutex 池。

由于多用于控制决策，凭据通常情况下是以只读方式访问的，此一般使用 `td_ucred`，因为它不需要上锁。当更新进程凭据时，必须和更新进程中持有 `proc`。读和更新操作必须使用 `p_ucred`，以避免死锁和使用中的条件。

如果所系线程将在更新进程凭据之后执行控制操作，则 `td_ucred` 也必须刷新当前进程的线程。这样做能避免修改后使用旧的凭据。内核会自己在进程中入内核，将进程对象的 `td_ucred` 指向刷新进程的 `p_ucred`，以保证内核线程能用到新的凭据。

### 8.4.2. 文件描述符和文件描述符表

所有内容将在后面加。

### 8.4.3. Jail 对象

`struct prison` 保存了用于那些通过 `jail(2)` API 建立的 jail 所用到的管理信息。包括 jail 的主机名、IP 地址，以及一些相关的位置。一个对象包含引用计数，因为指向同一对象实例的指针会在多个凭据对象之间共享。用了一个 mutex，`pr_mtx` 来保持引用计数以及所有 jail 对象中可共享的写。有一些量只会在创建 jail 的时刻产生变化，只需持有有效的 `struct prison` 就可以开始某些操作了。对于个别具体的上层操作的文档，可以在 `sys/jail.h` 的注释中找到。

### 8.4.4. MAC 框架

TrustedBSD MAC 框架会以 `struct label` 的形式提供一系列内核对象的数据。一般来说，内核中的 label (对象) 是由与其相关的内核对象同保有的。例如，`struct vnode` 上的 `v_label` 对象是由其所在 vnode 上的 vnode 保有的。

除了嵌入到标准内核对象中的对象之外，MAC 框架也需要一个包含已注册的和激活策略的列表。策略表和忙对象由一个全局 mutex (`mac_policy_list_lock`) 保护。由于能同时并行地执行多线程控制操作，策略表的只读，忙对象数，在入口处需要首先持有这个 mutex。MAC 入口操作的线程中并不需要一直持有此 mutex—有些操作，例如文件系统上的操作—是持久的。要修改策略表，例如在注册和解除注册策略时，需要持有此 mutex，而且要求引用计数为零，以避免在用表时对其进行修改。

由于需要等待表插入置状态的线程，提供了一个条件变量 `mac_policy_list_not_busy`，但该条件变量只能在使用者没有持有其它对象才能使用，否则可能会引起逆序问题。忙对象在整个框架中事实上扮演了某种形式的共享/排他对象：与 sx 不同的地方在于，等待列表插入置状态的线程可以死，而不是允忙对象和其他在 MAC 框架入口 (或内部) 的之中的逆序情况。

## 8.4.5. 模块

于模块子系统，用于保证共享数据使用了一个唯一的，它是一个 共享/排他 (SX)，很多时候需要获得它（以共享或排他的方式），因此我提供了几个方便使用的宏来简化这个过程，这些宏可以在 sys/module.h 中找到，其用法都非常明确了。一个模块的主要部分是 `module_t` (当以共享方式上) 和全局的 `modulelist_t` 两个实体，以及模块。要更进一步理解这些策略，需要仔细看 kern/kern\_module.c 的源代码。

## 8.4.6. Newbus 系统

newbus 系统使用了一个 sx 锁。写的一方持有共享 (S) 锁 (`sx_slock(9)`) 而写的一方持有排他 (写) 锁 (`sx_xlock(9)`)。内部函数一般不需要自行上锁，而外部可调的可以根据需要上锁。有些项目不需上锁，因为这些项目在全程是只读的，(例如 `device_get_softc(9)`)，因而并不会产生冲突条件。对于 newbus 数据结构的修改而言非常少，因此这个锁已足够使用，而不致造成性能折损。

## 8.4.7. 管道

...

## 8.4.8. 进程和线程

- 线程层次
- proc 线程及其参考
- 在系统中用进程中私有的 proc 副本，包括 td\_ucred
- 线程操作
- 线程和会话

## 8.4.9. 计时器

本文在其它地方已经提供了很多关于 `sched_lock` 的参考和注释。

## 8.4.10. Select 和 Poll

`select` 和 `poll` 两个函数允许多线程阻塞并等待文件描述符上的事件 — 最常见的情况是文件描述符是否可读或可写。

..

## 8.4.11. SIGIO

SIGIO 服务允许多线程请求在特定文件描述符的读/写状态变化，将 SIGIO 信号群映射到其线程。任意线程内核对象上，只允许多线程或线程对象注册 SIGIO，一个线程或线程对象称为主 (owner)。支持 SIGIO 注册的对象，都包含一个指向 `sigio` 结构的指针，如果对象未注册，则为 NULL，否则是指向描述符注册的 `struct sigio` 的指针。该字段由全局 mutex, `sigio_lock` 保护。使用 SIGIO 函数时，必须以“引用”方式设置该字段，以便保本地注册副本的中该字段不脱销的保护。

一个线程或线程对象的注册对象，都会分配一个 `struct sigio` 对象，并包括指向对象的指针、属主、信号信息、凭据，以及关于注册的一般信息。一个线程或线程对象都包含一个已注册 `struct sigio` 对象的列表，线程来

是 `p_sigiolst`, 而进程是 `pg_sigiolst`。有些表由相同的进程或线程共享。除了用以将 `struct sigio` 接到线程上的 `sio_pgsigio` 字段之外，在 `struct sigio` 中的多数字段在注册进程中都是不共享的。一般而言，人们在新的支持 SIGIO 的内核对象，会希望避免在公用 SIGIO 支持函数，例如 `fsetown` 或 `funsetown` 持有线程体，以免去需要在线程体和全局 SIGIO 之间锁定顺序。通常可以通提高线程体上的引用计数来达到目的，例如，在进行管道操作时，使用引用某个管道的文件描述符的操作，就可以照此办理。

## 8.4.12. Sysctl

`sysctl` MIB 服务会从内核内部，以及用它的公用程序以系统调用的方式触。它会引入至少两个和两个的接口：其一是持久命名空间的数据的保护，其二是与那些通过 `sysctl` 接口的内核变量和函数之的交互。由于 `sysctl` 允许直接读出（甚至修改）内核数据以及配置参数，`sysctl` 机制必须知道哪些变量相关的上层。目前，`sysctl` 使用一个全局 `sx` 来保证 `sysctl` 操作的串行化；然而，这些是假定用全局保护的，并且没有提供其它保护机制。唯一其余部分将介入上层和 `sysctl` 相关的。

- 需要将 `sysctl` 更新所执行的操作的顺序，从原先的旧、`copyin` 和 `copyout`、写新，改用 `copyin`、上、旧、写新、解、`copyout`。一般的 `sysctl` 只是 `copyout` 旧并置它用 `copyin` 所得到的新，仍然可以采用旧式的模型。然而，所有 `sysctl` 处理程序采用第二模型并避免操作方面，第二方式可能更直接一些。
- 由于通常的情况，`sysctl` 可以内嵌一个 `mutex` 指向到 `SYSCTL_FOO` 宏和线程体中。大多数 `sysctl` 都是有效的。由于使用 `sx`、自旋 `mutex`，或其它除休眠 `mutex` 之外的策略，可以用 `SYSCTL_PROC` 点来完成正确的上层。

## 8.4.13. 任务队列 (Taskqueue)

任务队列 (taskqueue) 的接口包括两个与之相关的用于保护共享数据的。`taskqueue_queues_mutex` 是用于保护 `taskqueue_queues` TAILQ 的。与两个系统相关的每一个 mutex 是位于 `struct taskqueue` 线程体上。在此使用同一原的目的在于保护 `struct taskqueue` 中数据的完整性。注意的是，并没有单独的、帮助用其自身的工作执行的宏，因为这些基本上不会在 `kern/subr_taskqueue.c` 以外的地方用到。

# 8.5. 休眠队列

## 8.5.1. 休眠队列

休眠队列是一个用于保存同一个等待通道 (wait channel) 上休眠线程列表的数据结构。在等待通道上，所有非睡眠状态的线程都会携一个休眠队列。当线程在等待通道上发生阻塞时，它会将休眠队列线程送那个等待通道。与等待通道的休眠队列保存在一个散列表中。

休眠队列散列表中保存了包含至少一个阻塞线程的等待通道上的休眠队列。一个散列表上的称作 `sleepqueue` (休眠队列)。它包含了一个休眠队列的表，以及一个自旋 `mutex`。此的自旋 `mutex` 用于保护休眠队列，以及其上休眠队列的内容。一个等待通道上只会有一个休眠队列。如果有多个线程在同一等待通道上阻塞，休眠队列中将除第一个线程之外的全部线程。当从休眠队列中除线程时，如果它不是唯一的阻塞的休眠线程，会得主休眠队列的空表上的休眠队列。最后一个线程会在恢复正常得主休眠队列。由于线程有可能以和加入休眠队列不同的次序从其中除，因此，线程队列可能会携与其插入不同的休眠队列。

`sleepq_lock` 函数会住指定等待通道上休眠队列的自旋 `mutex`。`sleepq_lookup` 函数会在主休眠队列散列表中确定的等待通道。如果没有到主休眠队列，它会返回 `NULL`。`sleepq_release` 函数会锁定等待通道所的自旋 `mutex` 行解。

将`程`加入休眠`列`是通过 `sleepq_add` 来完成的。该函数的参数包括等待通道、指向保`持`等待通道的 `mutex` 的指`针`、等待消息描述串，以及一个`标志掩码`。在用此函数之前，必须通过 `sleepq_lock` 休眠`列`上`锁`。如果等待通道不是通过 `mutex` 保`持`的（或者它由全局`保`），必须将 `mutex` 指`址`置`为 NULL`。而 `flags`（`标志`）参数`包括`了一个`类型`字段，用以表示`程`即将加入到的休眠`列`的`类型`，以及休眠是否是可中断的（SLEEPQ\_INTERRUPTIBLE）。目前只有`两种`类型的休眠`列`：通过 `msleep` 和 `wakeup` 函数管理的`休眠``列`（SLEEPQ\_MSLEEP），以及基于条件`量`的休眠`列`（SLEEPQ\_CONDVAR）。休眠`列``类型`和`指`定`参数`完全是用于内部的断言`。用` `sleepq_add` 的代`，明示地在`的 `sleepqueue` 透`过` `sleepq_lock` 行上`锁`之后，并使用等待函数在休眠`列`上阻塞之前解`所有`用于保`持`等待通道的 `interlock`。

通常使用 `sleepq_set_timeout` 可以`休眠`置超`。该`函数的参数包括等待通道，以及以相`应`数`数`位的超`。如果`休眠`被某个到来的信号打断，用` `sleepq_catch_signals` 函数，该函数唯一的参数就是等待通道。如果此`程`已有未决信号，`sleepq_catch_signals` 将返回信号`号`；其它情况下，其返回`0`。

一旦将`程`加入到休眠`列`中，就可以使用 `sleepq_wait` 函数族之一将其阻塞了。目前共提供了四个等待函数，使用`个`取决于`用`是否希望允`使用超`、收到信号，或用`程`度器打断休眠`。其中，``sleepq_wait` 函数`直接`地等待，直到当前`程`通过某个`醒`（wakeup）函数式地恢`复`；`sleepq_timedwait` 函数`间接`地等待，直到当前`程`被式地`醒`，或者`到早前使用` `sleepq_set_timeout` 置的超`；` `sleepq_wait_sig` 函数会等待式地`醒`，或者其休眠被中断；而 `sleepq_timedwait_sig` 函数`间接`地`醒`、`到用` `sleepq_set_timeout` 置的超`，或`程的休眠被中断三`条件`之一。`所有`些等待函数的第一个参数都是等待通道。除此之外，`sleepq_timedwait_sig` 的第二个参数是一个布`，表示之前用` `sleepq_catch_signals` 是否有`未决信号`。

如果`程`被式地恢`复`，或其休眠被信号`止`，等待函数会返回零，表示休眠成功。如果`程`的休眠被超`或用`程`度器`打断，会返回相`的` `errno` 数`。需要注意的是，因` `sleepq_wait` 只能返回 0，因此`用者不能指望它返回什`有用信息，而`假定它完成了一次成功的休眠。同`时，如果`程`的休眠`超`，并同`被`止，`sleepq_timedwait_sig` 将返回一个表示`生超`的`代`。如果返回`代`是 0 而且使用 `sleepq_wait_sig` 或 `sleepq_timedwait_sig` 来`行阻塞`，`用` `sleepq_calc_signal_retval` 来`是否`有`未决信号`，并据此`合`的返回`。早前用` `sleepq_catch_signals` 得到的信号`号`，`作`参数`sleepq_calc_signal_retval`。

在同一休眠通道上休眠的`程`，可以由 `sleepq_broadcast` 或 `sleepq_signal` 函数来式地`醒`。该函数的参数均包括希望`醒`的等待通道、将`醒`的`优先级`（priority）提高到多少，以及一个`标志`（flags）参数表示将要恢`复`的休眠`列``型`。`优先级`参数将作`最低`先`，如果将恢`的`程`的`优先级`比此参数更高（数`更低`）`其`先`不会`整。`标志`参数主要用于函数内部的断言，用以`休眠``列`没有被当做`的`型`待`。例如，条件`量`函数不`恢`休眠`列`的`行`。`sleepq_broadcast` 函数将恢`所有`指定休眠通道上的阻塞`程`，而 `sleepq_signal` 只恢`在`等待通道上`先`最高的阻塞`程`。在用`些`函数之前，首先使用 `sleepq_lock` 休眠`列`上`。`

休眠`程`也可以通过用 `sleepq_abort` 函数来中断其休眠`。该`函数只有在持有  `sched_lock` 才能`用`，而且`程`必须`于`休眠`列`之上。`程`也可以通过使用 `sleepq_remove` 函数从指定的休眠`列`中`除`。该函数包括`个`参数，即休眠通道和`程`，它只在`程`于指定休眠通道的休眠`列`之上`才`将其`醒`。如果`程`不在那个休眠`列`之上，或同`于`一等待通道的休眠`列`上，`个`函数将什`都`做而直接返回。

### 8.5.2. 十字口（turnstile）

- 与休眠`列`的比`和不同。`
- `/等待/放`（lookup/wait/release）- 介`TDF_TSNOBLOCK` 条件。

- 首先广播。

### 8.5.3. 用于 mutex 的一些...

- 我是否要求 mtx\_destroy() 持有 mutex， 因无法安全地断言它没有被其它对象持有？

#### 8.5.3.1. 自旋 mutex

- 使用一界区...

#### 8.5.3.2. 休眠 mutex

- 描述 mutex 冲突的条件
- 何在持有十字路口， 可以安全地冲突 mutex 的 mtx\_lock。

### 8.5.4. Witness

- 它能做什么
- 它如何工作

## 8.6. 其它...

### 8.6.1. 中断源和 ICU 抽象

- struct isrc
- pic ...

### 8.6.2. 其它.../...

- 是否将 interlock 为 sema\_wait？
- 是否提供非休眠式 sx ？
- 加一些于正使用引用数的介。

## ...

### 原子

当遵循适当的规则时，如果一操作的效果于其它所有 CPU 均可，称其原子操作。狭义的原子操作是机器直接提供的。就更高的抽象层次而言，如果个体的多个成员由一个保证，如果它们的操作都是在上后、解前进行的，也可以称其原子操作。

### 阻塞

进程等待、源或条件被阻塞。单一也因此被予了太多的意涵。

### 界区

不允许生占的代码段。使用 [critical\\_enter\(9\)](#) API 来表示进入和退出界区。

## **MD**

表示与机器/平台有口。

### 内存操作

内存操作包括口或写内存中的指定位置。

## **MI**

表示与机器/平台无口。

### 操作

#### 主中断上下文

主中断上下文表示当口生中断口所口行的那段代口。一些代口可以直接口行某个中断口理程序，或口度一口端口程，以便口定的中断源口行中断口理程序。

#### 内核口程

一口高口先口的内核口程。目前，只有中断口程属于口先口的内核口程。

### 休眠

当口程由条件口量或通口 msleep 或 tsleep 阻塞并口入休眠口列口，称其口入休眠状口。

### 可休眠口

可休眠口是一口在口程休眠口仍可持有的口。管理器 (lockmgr) 口和 sx 口是目前 FreeBSD 中口有的可休眠口。最口，某些 sx 口，例如 allproc (全部口程) 和 proctree (口程口) 口将成口不可休眠口。

### 口程

由 struct thread 所表口的内核口程。口程可以持有口，并口有独立的口行上下文。

### 等待通道

口程可以在其上休眠的内核虚口地址。

## 部分 II: 算法程序

# 写 FreeBSD 驱动程序

## 1. 介绍

本章要介绍了如何写 FreeBSD 驱动程序。这儿的上下文中多用于指代系统硬件相关的東西，如磁盘，打印机，图形显示器及其驱动。驱动程序是操作系统中用于控制特定硬件的软件。也有所谓的模块，即驱动程序用文件模版的形式存在，而没有特定的底层硬件。驱动程序可以被静态地链接到系统，或者通过内核直接工具‘kld’在需要时加载。

UNIX®操作系统的大多数都是通过点来访问的，有些也被称作特殊文件。这些文件在文件系统的层次中通常位于 /dev 目录下。在 FreeBSD 5.0-RELEASE 以前的发行版中，devfs(5) 的支持没有被集成到 FreeBSD 中，所以所有点必须要静态地建立，并且独立于相关驱动程序的存在。系统中大多数点是通过 MAKEDEV 建立的。

驱动程序可以粗略地分三类，字符和网络驱动程序。

## 2. 内核直接工具-KLD

kld 接口允许系统管理从运行的系统中动态地添加和删除功能。在运行的内核中，而不用为了新改动而繁地重编。

允许程序的作者将他的新改动加到

kld 接口通过下面的特命令使用：

- **kldload** - 加载新内核模块
- **kldunload** - 卸载内核模块
- **kldstat** - 列出当前加载的模块

内核模块的程序框架

```

/*
 * KLD程序框架
 * 受Andrew Reiter在Daemonnews上的文章所啓
 */

#include sys/types.h
#include sys/module.h
#include sys/sysctl.h /* uprintf */
#include sys/errno.h
#include sys/param.h /* kernel.h中用到的定义 */
#include sys/kernel.h /* 模块初始化中使用的类型 */

/*
 * 加载函数，管理KLD的加载和卸载。
 */

static int
skel_loader(struct module *m, int what, void *arg)
{
    int err = 0;

    switch (what) {
    case MOD_LOAD:           /* kldload */
        uprintf("Skeleton KLD loaded.\n");
        break;
    case MOD_UNLOAD:
        uprintf("Skeleton KLD unloaded.\n");
        break;
    default:
        err = EOPNOTSUPP;
        break;
    }
    return(err);
}

/* 向内核其余部分声明此模块 */

static moduledata_t skel_mod = {
    "skel",
    skel_loader,
    NULL
};

DECLARE_MODULE(skeleton, skel_mod, SI_SUB_KLD, SI_ORDER_ANY);

```

## 2.1. Makefile

FreeBSD提供了一个makefile包含文件，利用它可以快速地附加到内核的东西。

```
SRCS=skeleton.c  
KMOD=skeleton  
  
.include bsd.kmod.mk
```

本地用一个makefile行`make`就能创建文件 `skeleton.ko`, 运入如下命令可以把它加载到内核：

```
# kldload -v ./skeleton.ko
```

## 3. 系统程序

UNIX® 提供了一套公共的系统用供公用的应用程序使用。当用命令 `MAKEDEV`，内核的上层将一些公用分派到相应的程序。脚本 `/dev/MAKEDEV` 的系统生成了大多数的设备点，但如果正在编写自己的程序，可能需要用 `mknod` 创建自己的设备点。

### 3.1. 创建静态设备点

`mknod`命令需要四个参数来创建设备点。必须指定设备点的名字，它的类型，它的主号和它的从号。

### 3.2. 动态设备点

文件系统，或者`devfs`，在全局文件系统名字空间中提供内核名字空间的接口。消除了由于有程序而没有静态设备点，或者有设备点而没有安装程序而来的潜在问题。`Devfs`仍在发展中，但已能工作得相当好了。

## 4. 字符IO

字符设备程序直接从用户程序读数据，或把数据写到用户程序。

是最普通的一类程序，源码中有大量的

例子。一个例子会读写它的任何，而且当取它的时候

会将一些返回。下面展示了两个版本，一个

用于FreeBSD 4.X，一个用于FreeBSD 5.X。

#### 例4. 用于FreeBSD 4.X的回显驱动程序例

```
/*
 * echo'KLD
 *
 * Murray Stokely
 */

#define MIN(a,b) (((a) < (b)) ? (a) : (b))

#include sys/types.h
#include sys/module.h
#include sys/sysctl.h /* uprintf */
#include sys/errno.h
#include sys/param.h /* kernel.h中用到的定义 */
#include sys/kernel.h /* 模块初始化中使用的类型 */
#include sys/conf.h /* cdevsw */
#include sys/uio.h /* uio */
#include sys/malloc.h

#define BUFFERSIZE 256

/* 函数原型 */
d_open_t echo_open;
d_close_t echo_close;
d_read_t echo_read;
d_write_t echo_write;

/* 字符入口点 */
static struct cdevsw echo_cdevsw = {
    echo_open,
    echo_close,
    echo_read,
    echo_write,
    noioctl,
    nopoll,
    nommap,
    nosstrategy,
    "echo",
    33,           /* lkms保留 - /usr/src/sys/conf/majors */
    nodump,
    nopsize,
    D_TTY,
    -1
};

typedef struct s_echo {
    char msg[BUFFERSIZE];
    int len;
} t_echo;
```

```

/* 定义 */
static dev_t sdev;
static int count;
static t_echo *echomsg;

MALLOC_DECLARE(M_ECHOBUF);
MALLOC_DEFINE(M_ECHOBUF, "echobuffer", "buffer for echo module");

/*
 * 该函数被kld[un]load(2)系统调用使用,
 * 以决定加载和卸载模块需要采取的操作。
 */

static int
echo_loader(struct module *m, int what, void *arg)
{
    int err = 0;

    switch (what) {
    case MOD_LOAD:                      /* kldload */
        sdev = make_dev(echo_cdevsw,
                        0,
                        UID_ROOT,
                        GID_WHEEL,
                        0600,
                        "echo");
        /* kmalloc分配供程序使用的内存 */
        MALLOC(echomsg, t_echo *, sizeof(t_echo), M_ECHOBUF, M_WAITOK);
        printf("Echo device loaded.\n");
        break;
    case MOD_UNLOAD:
        destroy_dev(sdev);
        FREE(echomsg, M_ECHOBUF);
        printf("Echo device unloaded.\n");
        break;
    default:
        err = EOPNOTSUPP;
        break;
    }
    return(err);
}

int
echo_open(dev_t dev, int oflags, int devtype, struct proc *p)
{
    int err = 0;

    uprintf("Opened device \"echo\" successfully.\n");
    return(err);
}

```

```

int
echo_close(dev_t dev, int fflag, int devtype, struct proc *p)
{
    uprintf("Closing device \"echo.\\"\n\"");
    return(0);
}

/*
* read函数接受由echo_write()存入的buf，并将其返回到用户空间，  

* 以供其他函数使用。  

* uio(9)  

*/
int
echo_read(dev_t dev, struct uio *uio, int ioflag)
{
    int err = 0;
    int amt;

    /*
     * 读操作有多大？
     * 与用户要求的大小一致，或者等于剩余数据的大小。
     */
    amt = MIN(uio-uio_resid, (echomsg-len - uio-uio_offset) ?  

              echomsg-len - uio-uio_offset : 0);
    if ((err = uiomove(echomsg-msg + uio-uio_offset,amt,uio)) != 0) {
        uprintf("uiomove failed!\n");
    }
    return(err);
}

/*
* echo_write接受一个字符串并将它保存到缓冲区，用于以后的使用。
*/
int
echo_write(dev_t dev, struct uio *uio, int ioflag)
{
    int err = 0;

    /* 将字符串从用户空间的内存复制到内核空间 */
    err = copyin(uio-uio iov-iov_base, echomsg-msg,
                 MIN(uio-uio iov-iov_len, BUFSIZE - 1));

    /* 在需要以null结束字符串，并且长度 */
    *(echomsg-msg + MIN(uio-uio iov-iov_len, BUFSIZE - 1)) = 0;
    echomsg-len = MIN(uio-uio iov-iov_len, BUFSIZE);

    if (err != 0) {
        uprintf("Write failed: bad address!\n");
}

```

```
    }
    count++;
    return(err);
}

DEV_MODULE(echo,echo_loader,NULL);
```

## 例 5. 用于 FreeBSD 5.X 的简单程序示例

```
/*
 * 'echo' KLD
 *
 * Murray Stokely
 *
 * 此代码由 Søren (Xride) Straarup 贡献到 5.X
 */

#include sys/types.h
#include sys/module.h
#include sys/sysctl.h /* uprintf */
#include sys/errno.h
#include sys/param.h /* kernel.h 中用到的定义 */
#include sys/kernel.h /* 模块初始化中使用的类型 */
#include sys/conf.h /* cdevsw */
#include sys/uio.h /* uio */
#include sys/malloc.h

#define BUFFERSIZE 256

/* 函数原型 */
static d_open_t echo_open;
static d_close_t echo_close;
static d_read_t echo_read;
static d_write_t echo_write;

/* 字符串入口点 */
static struct cdevsw echo_cdevsw = {
    .d_version = D_VERSION,
    .d_open = echo_open,
    .d_close = echo_close,
    .d_read = echo_read,
    .d_write = echo_write,
    .d_name = "echo",
};

typedef struct s_echo {
    char msg[BUFFERSIZE];
    int len;
} t_echo;

/* 宏 */
static struct cdev *echo_dev;
static int count;
static t_echo *echomsg;

MALLOC_DECLARE(M_ECHOBUF);
MALLOC_DEFINE(M_ECHOBUF, "echobuffer", "buffer for echo module");
```

```

/*
 * 一个函数被kld[un]load(2)系统用来使用,
 * 以决定加载和卸载模块需要采取的工作.
 */

static int
echo_loader(struct module *m, int what, void *arg)
{
    int err = 0;

    switch (what) {
    case MOD_LOAD:           /* kldload */
        echo_dev = make_dev(echo_cdevsw,
            0,
            UID_ROOT,
            GID_WHEEL,
            0600,
            "echo");
        /* kmalloc分配供程序使用的内存 */
        echomsg = malloc(sizeof(t_echo), M_ECHOBUF, M_WAITOK);
        printf("Echo device loaded.\n");
        break;
    case MOD_UNLOAD:
        destroy_dev(echo_dev);
        free(echomsg, M_ECHOBUF);
        printf("Echo device unloaded.\n");
        break;
    default:
        err = EOPNOTSUPP;
        break;
    }
    return(err);
}

static int
echo_open(struct cdev *dev, int oflags, int devtype, struct thread *p)
{
    int err = 0;

    uprintf("Opened device \"echo\" successfully.\n");
    return(err);
}

static int
echo_close(struct cdev *dev, int fflag, int devtype, struct thread *p)
{
    uprintf("Closing device \"echo.\n\"");
    return(0);
}

```

```

/*
 * read函数接受由echo_write()存入的buf，并将其返回到用户空间，以供其他函数使用。
 * uio(9)
 */

static int
echo_read(struct cdev *dev, struct uio *uio, int ioflag)
{
    int err = 0;
    int amt;

    /*
     * 一个操作有多大？
     * 等于用uio求的大小，或者等于剩余数据的大小。
     */
    amt = MIN(uio->uio_resid, (echomsg->len - uio->uio_offset) ?
              echomsg->len - uio->uio_offset : 0);
    if ((err = uiomove(echomsg->msg + uio->uio_offset, amt, uio)) != 0) {
        printf("uiomove failed!\n");
    }
    return(err);
}

/*
 * echo_write接受一个字符串并将它保存到缓冲区，用于以后的读取。
 */
static int
echo_write(struct cdev *dev, struct uio *uio, int ioflag)
{
    int err = 0;

    /* 将字符串从用户空间的内存复制到内核空间 */
    err = copyin(uio->uio_iov->iov_base, echomsg->msg,
                 MIN(uio->uio_iov->iov_len, BUFSIZE - 1));

    /* 在需要以null结束字符串，并且长度 */
    *(echomsg->msg + MIN(uio->uio_iov->iov_len, BUFSIZE - 1)) = 0;
    echomsg->len = MIN(uio->uio_iov->iov_len, BUFSIZE);

    if (err != 0) {
        printf("Write failed: bad address!\n");
    }
    count++;
    return(err);
}

DEV_MODULE(echo, echo_loader, NULL);

```

在FreeBSD 4.X上安装此程序，将首先需要用如下命令在文件系上建一个点：

```
# mknod /dev/echo c 33 0
```

程序被加后，能输入一些东西，如：

```
# echo -n "Test Data" > /dev/echo
# cat /dev/echo
Test Data
```

真正的硬件在下一章描述。

充源

- [Dynamic Kernel Linker \(KLD\) Facility Programming Tutorial - Daemonnews](#) October 2000
- [How to Write Kernel Drivers with NEWBUS - Daemonnews](#) July 2000

## 5. (消亡中)

其他UNIX®系支持单一型的磁盘。称。是内核  
几乎没有用，或者非常不可用。冲会重新安排写操作的次序，使得用程序失了在任何时刻及知道准  
的磁内容的能力。致磁数据（文件系，数据等）的  
不可能。由于写操作被延，内核无法向用  
程序告个特定的写操作遇到了写，又一加了一致性  
。由于个原因，真正的用程序从不依于，事上，几乎所有  
磁的用程序都尽力指定  
是使用字符（或"raw"）  
。由于将个磁（分区）同具有不同类型的个混一，从而致使  
相  
内核代大大地化，作推磁I/O基化的一部分，FreeBSD 抛了冲的磁的支持。

## 6. 网程序

网的程序不需要使用点。个程序是  
的使用通常由系用socket(2)引入。

基于内核内部的其他决定而不是用open()

更多，参ifnet(9)机手册、回的源代码，以及 Bill Paul 撰写的网程序。

# Chapter 9. ISA程序

## 9.1. 概述

本章介绍了写ISA程序的一些。这儿展示的代码相当，很容易人想到真正的代码，不过依然不是代码。它避免了与所的主无的。真的例子可以在程序的源代码中到。ep和aha更是信息的好来源。

## 9.2. 基本信息

典型的ISA程序需要以下包含文件：

```
#include sys/module.h
#include sys/bus.h
#include machine/bus.h
#include machine/resource.h
#include sys/rman.h

#include isa/isavar.h
#include isa/pnpvar.h
```

它描述了ISA和通用子系的东西。

子系是以面向象的方式的，其主要通过相的方法函数来。

ISA程序的方法的列表与任何其他真的很相似。于名字"xxx"的假想程序，它将是：

- `static void xxx_isa_identify (driver_t *, device_t);` 通常用于程序而不是程序。但由于ISA，这个方法有特殊用途：如果提供某些特定的（非PnP）方法自，这个例程可以它。
- `static int xxx_isa_probe (device_t dev);` 在已知（或PnP）位置探。于已部分配置的，这个例程也能提供特定的某些参数的自。
- `static int xxx_isa_attach (device_t dev);` 挂接和初始化。
- `static int xxx_isa_detach (device_t dev);` 卸模块前解挂。
- `static int xxx_isa_shutdown (device_t dev);` 系统前执行的。
- `static int xxx_isa_suspend (device_t dev);` 系统入能状态前挂起。也可以中止切换到能状态。
- `static int xxx_isa_resume (device_t dev);` 从能状态返回后恢的活动状态。

`xxx_isa_probe()`和`xxx_isa_attach()`是必须提供的，其余例程根据的需要可以有地。

使用下面一描述符将驱动接到系统。

```

/* 支持的方法表 */
static device_method_t xxx_isa_methods[] = {
    /* 列出程序支持的所有方法函数 */
    /* 略去不支持的函数 */
    DEVMETHOD(device_identify, xxx_isa_identify),
    DEVMETHOD(device_probe,     xxx_isa_probe),
    DEVMETHOD(device_attach,   xxx_isa_attach),
    DEVMETHOD(device_detach,   xxx_isa_detach),
    DEVMETHOD(device_shutdown, xxx_isa_shutdown),
    DEVMETHOD(device_suspend,  xxx_isa_suspend),
    DEVMETHOD(device_resume,   xxx_isa_resume),

{ 0, 0 }
};

static driver_t xxx_isa_driver = {
    "xxx",
    xxx_isa_methods,
    sizeof(struct xxx_softc),
};

static devclass_t xxx_devclass;

DRIVER_MODULE(xxx, isa, xxx_isa_driver, xxx_devclass,
    load_function, load_argument);

```

此的`xxx_softc`是一个特定的，它包含私有的程序数据和程序源的描述符。代会自动按需要每个分配一个softc描述符。

如果程序作可加载模块，当程序被加载或卸载，会用`load_function()`函数执行程序特定的初始化或清理工作，并将`load_argument`作为函数的一个参数。如果程序不支持加载（即没有，它必须被接到内核中），有些当被置0，最后的定义将看起来如下所示：

```
DRIVER_MODULE(xxx, isa, xxx_isa_driver,
    xxx_devclass, 0, 0);
```

如果程序是支持PnP的而写的，那就必须定义一个包含所有支持的PnP ID的表。这个表由此程序所支持的PnP ID的列表和以人可读的形式列出的、与这些ID对应的硬件型号的描述组成。看起来如下：

```

static struct isa_pnp_id xxx_pnp_ids[] = {
    /* 每个所支持的PnP ID占一行 */
    { 0x12345678, "Our device model 1234A" },
    { 0x12345679, "Our device model 1234B" },
    { 0,         NULL }, /* 表结束 */
};

```

如果程序不支持PnP，它仍然需要一个空的PnP ID表，如下所示：

```
static struct isa_pnp_id xxx_pnp_ids[] = {  
    { 0,           NULL }, /* 表结束 */  
};
```

## 9.3. Device\_t指针

Device\_t是面向而定的指针型，这里我只关心从程序写者的角度看感兴趣的方法。下面的方法用来操作中的：

- device\_t device\_get\_parent(dev) 取父。
- driver\_t device\_get\_driver(dev) 取指向其程序的指。
- char \*device\_get\_name(dev) 取程序的名字，在我的例子中"xxx"。
- int device\_get\_unit(dev) 取元号（与个程序的从0始）。
- char \*device\_get\_nameunit(dev) 取名，包括元号。例如"xxx0", "xxx1" 等。
- char \*device\_get\_desc(dev) 取描述。通常它以人可的形式描述的切型号。
- device\_set\_desc(dev, desc) 设置描述信息。使得描述指向desc字符串，此后该字符串就不能被解除分配。
- device\_set\_desc\_copy(dev, desc) 设置描述信息。描述被拷到内部分配的缓冲区，desc字符串在以后可以被改而不会有害的结果。
- void \*device\_get\_softc(dev) 取指向与的描述符 (xxx\_softc) 的指。
- u\_int32\_t device\_get\_flags(dev) 取配置文件中特定于的标志。

可以使用一个很方便的函数device\_printf(dev, fmt, ...)从程序中打印信息。它自动在信息前添加元名和冒号。

device\_t的一些方法在文件kern/bus\_subr.c 中。

## 9.4. 配置文件与自配置周期和探的序

ISA在内核配置文件中的描述如下：

```
device xxx0 at isa? port 0x300 irq 10 drq 5  
iomem 0xd0000 flags 0x1 sensitive
```

端口、IRQ和其他被当成与的源。根据 自配置需要和支持程度的不同，有些是可的。例如，某些根本不需要DRQ，而有些允许从配置端口取 IRQ置。如果机器有多个ISA，可以在配置文件中明指定条，如isa0或isa1，否则将在所有ISA上搜索。

敏感(sensitive)是一源，它指示 必在所有非敏感之前探。此特性未被支持，但似乎从未在目前的任何程序中使用。

对于老的ISA设备，很多情况下程序仍然能配置参数。

但是系统中配置的多个设备

具有一个配置行。如果系统中装有同一类型的多个设备，但程序却只有一个配置行，例如：

```
device xxx0 at isa?
```

那么只有一个设备会被配置。

但对于支持通用PnP或只有自行自配置的设备，一个配置行

就足够配置系统中的所有设备，如上面的配置行，或者

等地：

```
device xxx at isa?
```

如果通用程序既支持能自行自配置的又支持老设备，并且

同时安装在一台机器上，那

只要在配置文件中描述老设备就足够了。自行自配置的设备将被自动添加。

如果ISA设备是自行自配置的，发生的事件如下：

所有通用程序的探测例程（包括所有PnP设备的PnP例程）

以随机顺序被调用。他们调用后就把设备

添加到ISA设备上的列表中。通常通用程序的探测例程将新设备与它的原有程序结合起来。而PnP例程并不知道其他通用程序，因此不能将新程序与它所添加的新设备结合起来。

使用PnP进入睡眠，以防止它被探测老设备。

被探测敏感(sensitive)的非PnP设备的探测例程被调用。如果探测成功，那就调用挂接(attach)例程。

所有非PnP设备的探测和挂接例程以同样的方式被调用。

PnP从睡眠中恢复来，并分配所要求的资源：I/O、

内存地址、IRQ和DRQ，所有这些与已连接的老设备不会冲突。

对于每个PnP设备，所有ISA设备的探测例程都会被调用。

第一个要求此设备的程序将被挂接。多个

设备以不同的优先级要求一个设备的情况是可能的，这种情况下，具有最高优先级的设备将被挂接。探测例程必须调用 **ISA\_PNP\_PROBE()** 将真正的PnP ID和该设备支持的ID列表作比较，如果ID不在表中返回失败。这意味着一个设备，包括不支持任何PnP的设备，都必须未知的PnP设备无条件调用 **ISA\_PNP\_PROBE()**，对于未知设备，至少要用一个空的PnP ID表调用并返回失败。

探测例程遇到失败会返回一个正数（非零），成功返回零或负数。

负数的返回用于PnP支持多个接口的情况。例如，老的兼容接口

和新的高接口通过不同的

设备来提供支持。多个设备都存在。在探测例程中返回高优先级的设备首先（句号，然后返回0的设备具有最高的优先级，返回-1的其次，返回-2的更在其后，如此下去）。如果多个设备返回相同的负数，那么最先调用的设备。因此，如果设备返回0，就基本能相信它先于仲裁。

特定的探测例程也能将一个而不是多个设备指派给PnP。就象使用PnP的情况一样，对于某一设备，会探测系统中所有的设备。由于特性在任何现存的设备中均未实现，故本文中不再予以考虑。

由于探测老设备的时候PnP被禁用，它将不会被探测两次（一次作为老设备，一次作为PnP）。但如果探测例程相对的，情况下两个设备有任保同一设备不会被探测两次：一次作为老的由用配置的，一次作为自行自配置的。

对于自定义的设备（包括PnP和特定的）的一个实践是，不能从内核配置文件中向它添加旗子。因此它必须根本不使用旗子，要将所有自定义的设备使用元号的旗子，或者使用sysctl接口而不是旗子。

通常使用函数族`resource_query_*`()和`resource_*_value()`直接配置资源，从而可以提供其他不常用的配置。它们位于`kern/subr_bus.c`。老的IDE磁盘驱动器`i386/isa/wd.c`包含使用的例子。但必须先使用配置的准确方法。将解析配置资源的事情留给配置代码。

## 9.5. 资源

写入到内核配置文件中的信息被作配置资源处理，并加载到内核。

配置代码解析

部分信息并将其映射到`device_t`的和与之相关的资源。对于某些情况下的配置，程序可以直接使用`resource_*`函数配置资源。然而，通常既不需要也不推荐这样做，因此这儿不再一一列举。

资源与个数相关。通过型和型中的数字指定它们。对于ISA，定义了下面的型：

- `SYS_RES_IRQ` - 中断号
- `SYS_RES_DRQ` - ISA DMA通道号
- `SYS_RES_MEMORY` - 映射到系统内存空间的内存的地址
- `SYS_RES_IOPORT` - ISA I/O寄存器的地址

型内的枚举从0开始，因此如果有个内存区域，它的`SYS_RES_MEMORY`型的资源号为0和1。资源型与C语言的型无关，所有资源具有C语言`unsigned long`型，并且必须进行强制转换(cast)。资源号不必为0，尽管对于ISA它一般是0的。ISA允许的资源号为：

```
IRQ: 0-1  
DRQ: 0-1  
MEMORY: 0-3  
IOPORT: 0-7
```

所有资源被表示为有起始和结束的范围。对于IRQ和DRQ资源，结束一般等于1。内存的引用物理地址。

资源能执行三种类型的操作：

- set/get
- allocate/release
- activate/deactivate

Set置资源使用的。Allocation保留出请求的，使得其它不能再占用（并且在此之前没有被其它占用）。Activation执行必要的操作使得程序可以访问资源（例如，对于内存，它将被映射到内核的虚拟地址空间）。

操作资源的函数有：

- `int bus_set_resource(device_t dev, int type, int rid, u_long start, u_long count)`

资源置。成功返回0，否则返回非0。一般此函数只有在`type`, `rid`, `start`或`count`之一的值超出了允许的范围才会返回非0。

- dev - 程序的设备
- type - 资源类型， SYS\_RES\_\*
- rid - 型内部的资源号 (ID)
- start, count - 资源大小
- `int bus_get_resource(device_t dev, int type, int rid, u_long *startp, u_long *countp)`

取得资源。成功返回0，如果资源尚未定位返回ENOENT。

- `u_long bus_get_resource_start(device_t dev, int type, int rid) u_long bus_get_resource_count (device_t dev, int type, int rid)`

便捷函数，只用来取start或count。出错的情况下返回0，因此如果0是资源的start合法之一，将无法区分返回的0是否指示合法。幸运的是，由于附加程序，没有ISA资源的start从0开始。

- `void bus_delete_resource(device_t dev, int type, int rid)`

删除资源，令其未定位。

- `struct resource * bus_alloc_resource(device_t dev, int type, int *rid, u_long start, u_long end, u_long count, u_int flags)`

在start和end之间没有被其它占用的地方按count分配一个资源。不支持。如果不源尚未被置，自建它。start=0, end≈0 (全1) 的特殊意味着必须使用以前通过`bus_set_resource()`置的固定值：start和count就是它自己，end=(start+count)，情况下，如果以前资源没有定位，返回ENOENT。尽管rid通常引用，但它并不被ISA的资源分配代替置（其它可能使用不同的方法并可能修改它）。

旗帜是一个位映射，使用者感兴趣的有：

- `RF_ACTIVE` - 使得资源分配后被自动激活。
- `RF_SHAREABLE` - 资源可以同时被多个程序共享。
- `RF_TIMESHARE` - 资源可以被多个程序分时共享，也就是说，被多个程序同时分配，但任何时刻只能被其中一个激活。
- 出错返回0。被分配的可以使用 `rhand_*` 从返回的句柄获得。
- `int bus_release_resource(device_t dev, int type, int rid, struct resource *r)`
- 放弃资源，`r``bus_alloc_resource()` 返回的句柄。成功返回0，否则返回ENOENT。
- `int bus_activate_resource(device_t dev, int type, int rid, struct resource *r) int bus_deactivate_resource(device_t dev, int type, int rid, struct resource *r)`
- 激活或禁用资源。成功返回0，否则返回ENOENT。如果资源被分时共享且当前被同一程序激活，返回EBUSY。
- `int bus_setup_intr(device_t dev, struct resource *r, int flags, driver_intr_t *handler, void *arg, void **cookiep) int bus_teardown_intr(device_t dev, struct resource *r, void *cookie)`
- 分时中断处理程序与。成功返回0，否则返回ENOENT。
- `r` - 被激活的描述IRQ的资源句柄。

`flags` - 中断优先级，如下之一：

- INTR\_TYPE\_TTY - 端和其它似的字符型。使用 `spltty()` 屏蔽它。
- (INTR\_TYPE\_TTY | INTR\_TYPE\_FAST) - 入冲小的端型，而且入上的数据失很（例如老式串口）。使用 `spltty()` 屏蔽它。
- INTR\_TYPE\_BIO - 型，不包括CAM控制器上的。使用 `splbio()` 屏蔽它。
- INTR\_TYPE\_CAM - CAM (通用方法Common Access Method) 控制器。使用 `splcam()` 屏蔽它。
- INTR\_TYPE\_NET - 网接口控制器。使用 `splimp()` 屏蔽它。
- INTR\_TYPE\_MISC - 各其它。除了通过 `splhigh()` 没有其它方法屏蔽它。 `splhigh()` 屏蔽所有中断。

当中断处理程序运行，匹配其先的所有其它中断都被屏蔽，唯一的例外是 MISC，它不会屏蔽其它中断，也不会被其它中断屏蔽。

- *handler* - 指向处理程序的指，型 `driver_intr_t` 被定义为 `void driver_intr_t(void *)`
- *arg* - 处理程序的参量，特定。由处理程序将它从 `void*` 任何型。ISA 中断处理程序的旧定是使用元号作参量，新定（推）使用指向 `softc` 的指。
- *cookie[p]* - 从 `setup()` 接收的，当 `teardown()` 用于处理程序。

定了若干方法来操作源句柄(`struct resource *`)。程序写者感兴趣的有：

- `u_long rman_get_start(r)` `u_long rman_get_end(r)` 取得被分配的源的起始和结束。
- `void *rman_get_virtual(r)` 取得被激活的内存源的虚地址。

## 9.6. 内存映射

很多情况下程序和之的数据交是通过内存进行的。有可能的体：

(a) 内存位于上

(b) 内存计算机的主内存

情况(a)中，程序可能需要在上的内存与主存之间来回拷数据。为了将上的内存映射到内核的虚地址空间，上内存的物理地址和必须被定义为 `SYS_RES_MEMORY` 源。然后源就可以被分配并激活，它的虚地址通过使用 `rman_get_virtual()` 取。老的程序将函数 `pmap_mapdev()` 用于此目的，现在不再直接使用此函数。它已成为源激活的一个内部。

大多数ISA的内存配置物理地位于640KB-1MB之间的某个位置。某些ISA需要更大的内存，位于16M以下的某个位置（由于ISA上24位地址限制）。情况下，如果机器有比内存的起始地址更多的内存（即空洞，它重），必须在被使用的内存起始地址配置一个内存空洞。BIOS允许在起始于14MB或15MB配置1M的内存空洞。如果BIOS正地宣告内存空洞，FreeBSD就能正确处理它（此特性在老BIOS上可能会出问题）。

情况(b)中，只是数据的地址被送到，使用DMA通过主存中的数据。存在个限制：首先，ISA只能16MB以下的内存。其次，虚地址空间中的一面在物理地址空间中可能不，可能不得不实行分散/收集操作。子系统有些提供成的解决办法，剩下的必须由程序自己完成。

DMA内存分配使用了两个，`bus_dma_tag_t` 和 `bus_dmamap_t`。`tag` 描述了 DMA内存要求的特性。映射(map)表示按照些特性分配的内存。多个映射可以与同一。

按照特性的承而被分成型层次。子承父的所有要求，可以令其更严格，但不允许放要求。

一般地，一个单元建一个（没有父）。如果建一个，有些作父的孩子。

个需要不同要求的内存区，多个内存区都会

使用建映射的方法有。

其一，分配一大符合要求的内存（以后可以被放）。一般用于分配了与通信而存在相同的那些内存区。将的内存加到映射中非常容易：它是被看作位于当物理内存的一整。

其二，将虚内存中的任意区域加到映射中。片内存的都被，看是否符合映射的要求。如何符合留在原始位置。如果不分配一个新的符合要求的"反面(bounce page)"，用作中存。当从不符合的原始面写入数据，数据首先被拷到反面，然后从反面到。当取，数据将会从到反面，然后被拷到它不符合的原始面。原始和反面之的拷物理被称作同。一般用于次的基之上：次加缓冲区，完成，卸缓冲区。

工作在DMA内存上的函数有：

- `int bus_dma_tag_create(bus_dma_tag_t parent, bus_size_t alignment, bus_size_t boundary, bus_addr_t lowaddr, bus_addr_t highaddr, bus_dma_filter_t *filter, void *filterarg, bus_size_t maxsize, int nsegments, bus_size_t maxsegsz, int flags, bus_dma_tag_t *dmat)`

建新。成功返回0，否返回。

- *parent* - 父或者NULL，NULL用于建。
- *alignment* - 将要分配的内存区的要求。"no specific alignment"01。用于以后的`bus_dmamem_alloc()`而不是`bus_dmamap_create()`用。
- *boundary* - 物理地址界，分配内存不能穿。于"no boundary" 使用0。用于以后的`bus_dmamem_alloc()`而不是`bus_dmamap_create()`用。必2的乘方。如果以非DMA方式使用内存（也就是，DMA地址由ISA DMA控制器提供而不是自身），由于DMA硬件限制，界必不能大于64KB (64\*1024)。
- *lowaddr, highaddr* - 名字微有些。些用于限制可用于内存分配的物理地址的允。其切含根据以后不同的使用而有所不同。
  - 于`bus_dmamem_alloc()`，从0到lowaddr-1的所有地址被允，更高的地址不允使用。
  - 于`bus_dmamap_create()`，区[lowaddr; highaddr]之外的所有地址被不可。之内的地址面被函数，由它决定是否可。如果没有提供函数，整个区被不可。
  - 于ISA，正常（没有函数）：

`lowaddr = BUS_SPACE_MAXADDR_24BIT`

`highaddr = BUS_SPACE_MAXADDR`

- *filter, filterarg* - 函数及其参数。如果filterNULL，当用`bus_dmamap_create()`，整个区[lowaddr, highaddr]被不可。否，区[lowaddr; highaddr]内的个被的面的物理地址被函数，由它决定是否可。函数的原型：`int filterfunc(void *arg, bus_addr_t paddr)`。当面可以被它必须返回0，否返回非零。
- *maxsize* - 通此可以分配的最大内存（以字）。有个很估算，或者可以任意大，情况下，于ISA个可以`BUS_SPACE_MAXSIZE_24BIT`。

- *nsegments* - 支持的分散/收集段的最大数目。如果不加限制，使用当使用 `BUS_SPACE_UNRESTRICTED`。建父使用 0 个，而子指定限制。*nsegments* 等于 `BUS_SPACE_UNRESTRICTED` 的不能用于附加映射，可以将它作父。*nsetsments* 的限制大约 250-300，再高的将导致内核堆溢出（硬件无法正常支持那么多的分散/收集缓冲区）。
- *maxsegsz* - 支持的分散/收集段的最大尺寸。大于 ISA 的最大 `BUS_SPACE_MAXSIZE_24BIT`。
- *flags* - 旗的位。感兴趣的旗只有：
  - `BUS_DMA_ALLOCNOW` - 建立请求分配所有可能用到的反射面。
  - `BUS_DMA_ISA` - 比神秘的一个标志，用于 Alpha 机器。i386 机器没有定它。Alpha 机器的所有 ISA 都当使用这个标志，但似乎没有它的程序。
- *dmat* - 指向返回的新存的指。
- `int bus_dma_tag_destroy(bus_dma_tag_t dmat)`

成功返回 0，否则返回非 0。

*dmat* - 被的。
- `int bus_dmamem_alloc(bus_dma_tag_t dmat, void** vaddr, int flags, bus_dmamap_t *mapp)`

分配所描述的一内存区。被分配的内存的大小的 `maxsize`。成功返回 0，否则返回非 0。如果被用于取内存的物理地址，但在此之前必须用 `bus_dmamap_load()` 将其加。

  - *dmat* -
  - *vaddr* - 指向存的指，存空用于返回的分配区域的内核虚地址。
  - *flags* - 旗的位。唯一感兴趣的旗：
    - `BUS_DMA_NOWAIT` - 如果内存不能立即可用返回非 0。如果此标志没有设置，允许例程睡眠，直到内存可用为止。
  - *mapp* - 指向返回的新映射的存的指。
- `void bus_dmamem_free(bus_dma_tag_t dmat, void *vaddr, bus_dmamap_t map)`

放由 `bus_dmamem_alloc()` 分配的内存。目前，分配的有 ISA 限制的内存的放没有。因此，建的使用模型尽可能地保持和重用分配的区域。不要轻易地放某些区域，然后再短地分配它。并不意味着不当使用 `bus_dmamem_free()`：希望很快它就会被完整地。

  - *dmat* -
  - *vaddr* - 内存的内核虚地址
  - *map* - 内存的映射（跟 `bus_dmamem_alloc()` 返回的一）
- `int bus_dmamap_create(bus_dma_tag_t dmat, int flags, bus_dmamap_t *mapp)`

建映射，以后用于 `bus_dmamap_load()`。成功返回 0，否则返回非 0。

  - *dmat* -
  - *flags* - 理上是旗的位。但从未定任何旗，因此目前是 0。
  - *mapp* - 指向返回的新映射的存的指。

- `int bus_dmamap_destroy(bus_dma_tag_t dmat, bus_dmamap_t map)`

解除映射。成功返回0，否则返回错误。

- `dmat` - 与映射相关的DMA tag
- `map` - 将要被解除的映射

- `int bus_dmamap_load(bus_dma_tag_t dmat, bus_dmamap_t map, void *buf, bus_size_t buflen, bus_dmamap_callback_t *callback, void *callback_arg, int flags)`

将物理内存区到映射中(映射必须事先由 `bus_dmamap_create()` 或者 `bus_dmamem_alloc()` 建立)。映射的所有面都会被反向，看是否符合的要求，并将那些不符合的分配到反向面。会创建物理段描述符的数目，并将其回函数。回函数以某种方式整理个数。系统中的反向映射是受限的，因此如果需要的反向映射不能立即获得，将请求插入，当反向映射可用再调用回函数。如果回函数立即执行并返回0，如果请求被排队，等待将来执行，返回 `EINPROGRESS`。后一种情况下，与排队的回函数之同由程序决定。

- `dmat` - DMA tag
- `map` - 映射
- `buf` - 映射的内核虚地址
- `buflen` - 映射的长度
- `callback, callback_arg` - 回函数及其参数

回函数的原型：

```
void callback(void *arg, bus_dma_segment_t *seg, int nseg, int error)
```

- `arg` - 与 `bus_dmamap_load()` 的 `callback_arg` 相同。
- `seg` - 段描述符的数目
- `nseg` - 数目中的描述符个数
- `error` - 表示段数目溢出：如被 `EFBIG`，允许的最大数目的段无法容纳映射区。这种情况下数目的描述符的数目只有可的那么多。具体情况的理由由程序决定：根据希望的，程序可以将其忽略，或将映射区分多个并单独处理第二个。

段数目的每一项包含如下字段：

- `ds_addr` - 段物理地址
- `ds_len` - 段长度

- `void bus_dmamap_unload(bus_dma_tag_t dmat, bus_dmamap_t map)`

卸载映射。

- `dmat` - DMA tag
- `map` - 已加载的映射

- `void bus_dmamap_sync (bus_dma_tag_t dmat, bus_dmamap_t map, bus_dmasync_op_t op)`

与物理映射前后，将加载的映射区与其反向面执行同步。  
所有必需的数据拷贝工作。执行之前和之后必须映射区执行同步。

此函数完成原始映射区与其映射版本之

- *dmat* - *DMA*
- *map* - 已加~~的~~映射
- *op* - 要~~行~~的同~~口~~操作的~~口~~型：
- **BUS\_DMASYNC\_PREREAD** - 从~~口~~到~~口~~冲区的~~口~~操作之前
- **BUS\_DMASYNC\_POSTREAD** - 从~~口~~到~~口~~冲区的~~口~~操作之后
- **BUS\_DMASYNC\_PREWRITE** - 从~~口~~冲区到~~口~~的写操作之前
- **BUS\_DMASYNC\_POSTWRITE** - 从~~口~~冲区到~~口~~的写操作之后

当前PREREAD和POSTWRITE~~空~~操作，但将来可能会改~~口~~，因此~~口~~程序

中不能忽略它~~口~~

◦ 由bus\_dmamem\_alloc()~~得~~的内存不需要同~~口~~。

从bus\_dmamap\_load()~~中~~用回~~口~~函数之前，~~段数~~是存~~口~~在~~口~~中的。并且是按~~口~~允~~口~~的最大数目的段~~口~~先分配好的。~~口~~于i386体系~~上~~上~~口~~段数目的~~口~~限制~~口~~250-300 (内核~~口~~4KB~~去用~~~~口~~的大小，~~段数~~条目的大小~~8字~~，和~~其它必~~留出来的空~~口~~)。由于数~~口~~基于最大数目而分配，因此~~个~~必~~不能~~置成超出~~口~~需要。幸~~口~~的是，~~于~~大多数硬件而言，~~所支持的段的最大数目低很多~~。但如果~~口~~程序想~~理~~具有非常多~~分散/收集段的~~冲区，~~当一部分一部分地~~理：加~~口~~冲区的~~一部分，~~到~~到~~口~~，然后加~~口~~冲区的下一部分，如此反~~口~~。~~

~~一个~~践~~口~~是段数目可能限制~~口~~冲区的大小。如果~~口~~冲区中的所有~~面~~巧物理上不~~口~~，~~分片~~情况下支持的最大~~口~~冲区尺寸 ~~(nsegments \* page\_size)~~。例如，如果支持的段的最大数目~~10~~，~~在~~i386上可以~~保~~支持的最大~~口~~冲区大小~~40K~~。如果希望更大的~~口~~需要在~~口~~程序中使用一些特殊技巧。

如果硬件根本不支持分散/收集，或者~~口~~程序希望即使在~~重分片的~~情况下仍然支持某~~口~~冲区大小，~~解决~~法是：如果无法容~~口~~下原始~~口~~冲区，就在~~口~~程序中分配一个~~口~~的~~口~~冲区作~~中~~存~~口~~。

下面是当使用映射~~口~~的典型~~用~~序，根据~~口~~映射的具体使用而不同。字符-用于~~示~~流。

~~于~~从~~接~~到分~~口~~，~~期~~位置一直不~~口~~的~~口~~冲区：

```
bus_dmamem_alloc - bus_dmamap_load - ...use buffer... - - bus_dmamap_unload -  
bus_dmamem_free
```

~~于~~从~~口~~程序外部~~口~~去，并且~~常~~化的~~口~~冲区：

```
bus_dmamap_create -  
- bus_dmamap_load - bus_dmamap_sync(PRE...) - do transfer -  
- bus_dmamap_sync(POST...) - bus_dmamap_unload -  
...  
- bus_dmamap_load - bus_dmamap_sync(PRE...) - do transfer -  
- bus_dmamap_sync(POST...) - bus_dmamap_unload -  
- bus_dmamap_destroy
```

当加~~由~~bus\_dmamem\_alloc()~~建~~的映射~~口~~，~~口~~

~~去~~的~~口~~冲区的地址和大小必~~口~~和

bus\_dmamem\_alloc()~~中~~使用的一~~口~~。~~口~~情况下就

可以保~~口~~整个~~口~~冲区被作~~口~~一个段而映射(因而回

~~口~~可以基于此假~~口~~)，

并且~~口~~求被立即~~口~~行(永~~口~~不会返回EINPROGRESS)。~~口~~情况下回~~口~~函数

需要作的只是保存物理地址。

典型示例如下：

```

static void
alloc_callback(void *arg, bus_dma_segment_t *seg, int nseg, int error)
{
    *(bus_addr_t *)arg = seg[0].ds_addr;
}

...
int error;
struct somedata {
    ...
};

struct somedata *vsodata; /* 虚地址 */
bus_addr_t psodata; /* 物理映射的地址 */
bus_dma_tag_t tag_somedata;
bus_dmamap_t map_somedata;
...

error=bus_dma_tag_create(parent_tag, alignment,
    boundary, lowaddr, highaddr, /*filter*/ NULL, /*filterarg*/ NULL,
    /*maxsize*/ sizeof(struct somedata), /*nsegments*/ 1,
    /*maxsegsz*/ sizeof(struct somedata), /*flags*/ 0,
);
if(error)
    return error;

error = bus_dmamem_alloc(tag_somedata, , /* flags*/ 0,
    );
if(error)
    return error;

bus_dmamap_load(tag_somedata, map_somedata, (void *)vsodata,
    sizeof (struct somedata), alloc_callback,
    (void *) , /*flags*/0);

```

代码看起来有点乱，也比较复杂，但那是正确的使用方法。结果是：如果分配多个内存区域，将它们合成一个并作整体分配（如果限制允许的话）是一个很好的主意。

当添加任意缓冲区到由`bus_dmamap_create()`创建的映射时，由于回调可能被延时，因此必须采取特殊措施与回调函数同行。代码看起来像下面的样子：

```

{
    int s;
    int error;

    s = splsoftvm();
    error = bus_dmamap_load(
        dmat,
        dmamap,
        buffer_ptr,
        buffer_len,
        callback,
        /*callback_arg*/ buffer_descriptor,
        /*flags*/0);
    if (error == EINPROGRESS) {
        /*
         * 行必要的操作以保与回的同。
         * 回被保直到我行了splx()或tsleep()才会被用。
         */
    }
    splx(s);
}

```

理求的方法分是：

1. 如果通式地求已束来完成求（例如CAM求），将所有一的理放入回程序中会比，回束后会求。之后不需要太多外的同。由于流控制的原因，求列直到求完成才放可能是个好主意。
2. 如果求是在函数返回完成（例如字符上写的写求），需要在冲区描述符上置同志，并用`tsleep()`。后面当回函数被用，它将行理并同志。如果置了同志，它调出一个醒操作。在方法中，回函数或者行所由必需的理（就像前面的情况），或者在冲区描述符中存段数。回完成 后，回函数就能使用个存的段数并行所有的理。

## 9.7. DMA

ISA中Direct Memory Access (DMA)是通DMA控制器（上是它中的个，但只是无）的。为了使以前的ISA便宜，控制和地址生的都集中在DMA控制器中。幸的是，FreeBSD提供了一套函数，些函数大多把DMA控制器的繁程序藏了起来。

最情况是那些比智能的。就象PCI上的主一，它自己能生周期和内存地址。它真正从DMA控制器需要的唯一事情是仲裁。所以了此目的，它假装是从DMA控制器。当接程序，系DMA控制器需要做的唯一事情就是通用如下函数在一个DMA通道上激活模式。

`void isa_dmacascade(int channel_number)`

所有一的活通程完成。当卸程序，不需要用DMA相的函数。

于的，事情反而得。使用的函数包括：

- `int isa_dma_acquire(int channel_number)`

保留一个DMA通道。成功返回0，如果通道已被保留或被其它程序保留返回EBUSY。大多数的ISA都不能共享DMA通道，因此一个函数通常在接用。源的代接口使得保留成多余，但目前仍必使用。如果不使用，后面其它 DMA例程将会panic。

- `int isa_dma_release(int channel_number)`

放先前保留的DMA通道。放通道必须不能有正在执行中的（外，放通道后必须不能再起）。

- `void isa_dmainit(int chan, u_int bouncebufsize)`

分配由特定通道使用的反冲区。求的冲区大小不能超64KB。以后，如果冲区巧不是物理的，或超出ISA可的内存，或跨越64KB的界，会自动使用反冲区。如果是使用符合上述条件的冲区（例如，由`bus_dmamem_alloc()`分配的那些），不需要用`isa_dmainit()`。但使用此函数会通DMA控制器任意数据得非常方便。

- *chan* - 通道号

- *bouncebufsize* - 以字数的反冲区的大小

- `void isa_dmastart(int flags, caddr_t addr, u_int nbytes, int chan)`

准DMA。上的之前必需用此函数来置DMA控制器。它冲区是否的且在ISA内存之内，如果不是自使用反冲区。如果需要反冲区，但反冲区没有用`isa_dmainit()`置，或于求的大小来太小，系将panic。写求且使用反冲区的情况下，数据将被自拷到反冲区。

- *flags* - 位掩，决定将要完成的操作的型。方向位B\_READ和B\_WRITE互斥。

- B\_READ - 从ISA到内存

- B\_WRITE - 从内存写到ISA上

- B\_RAW - 如果置DMA控制器将会住冲区，并在结束后自动重新初始化它自己，再次重同一冲区（当然，程序可能起的一个之前改冲区中的数据）。如果没有置，参数只一次有效，在起下一次之前必须再次用`isa_dmastart()`。只有在不使用反冲区使用B\_RAW才有意。

- *addr* - 冲区的虚地址

- *nbytes* - 冲区度。必须小于等于64KB。不允许度：因 DMA控制器将会理解64KB，而内核代把它理解0，那就就会致不可的效果。于通道号等于和高于4的情况，度必需偶数，因些通道次2字。奇数度情况下，最后一个字不被。

- *chan* - 通道号

- `void isa_dmadone(int flags, caddr_t addr, int nbytes, int chan)`

告完成后，同内存。如果是使用反冲区的操作，将数据从反冲区拷到原始冲区。参量与`isa_dmastart()`的相同。允使用B\_RAW志，但它一点也不会影`isa_dmadone()`。

- `int isa_dmastatus(int channel_number)`

返回当前中剩余的字数。在`isa_dmastart()`中置了B\_READ的情况下，返回的数字一定不会等于零。束它会被自动位到冲区的度。正式的用法是在信号指示已完成剩余的字数。如果字数不0，此次可能有。

- `int isa_dmastop(int channel_number)`

放回当前的值并返回剩余未读的字节数。

## 9.8. xxx\_isa\_probe

该函数探测是否存在。如果程序支持自定义配置的某些部分（如中断向量或内存地址），自定义必须在此例程中完成。

对于任意其他值，如果不能读到值，或者读到但丢失，或者产生某些其他值，当返回一个正的值。如果不存在必返回 `ENXIO`。其他值可能表示其他条件。零或意味着成功。大多数程序返回零表示成功。

当PnP支持多个接口时使用返回。例如，不同程序支持老的兼容接口和新的高接口。两个程序都将。在探例程中返回高接口的程序得先（句，返回0的程序具有最高的优先级，返回1的其次，返回-2的更后，等等）。同时，支持老接口的将被老程序处理（其当从探例程中返回1），而同时也支持新接口的将由新程序处理（其当从探例程中返回0）。

描述符xxx\_softc由系统在用探例程之前分配。如果探例程返回0，描述符会被系统取消分配。因此如果出于安全地取消分配。如果探成功完成，描述符将由系统保存并在以后的例程xxx\_isa\_attach()。如果程序返回0，就不能保证它将获得最高优先级且其接例程会被调用。因此情况下它也必须在返回前放所有的资源，并在需要的时候在接例程中重新分配它。当xxx\_isa\_probe()返回0，在返回前放资源也是一个好主意，而且中矩的情况下，允许程序在从探例程返回0和接例程的执行之间保持资源。

典型的探例程以取得描述符和元号开始：

```
struct xxx_softc *sc = device_get_softc(dev);
int unit = device_get_unit(dev);
int pnperror;
int error = 0;

sc-dev = dev; /* 接回来 */
sc-unit = unit;
```

然后调用PnP。这是通过一个包含PnP ID列表的表进行的。此表包含所有程序支持的PnP ID和以人工形式列出的一些ID的型号的描述。

```
pnperror=ISA_PNP_PROBE(device_get_parent(dev), dev,
xxx_pnp_ids); if(pnperror == ENXIO) return ENXIO;
```

`ISA_PNP_PROBE` 的如下：如果（元）没有被作PnP到，返回`ENOENT`。如果被作PnP到，但到的ID不匹配表中的任一ID，返回`ENXIO`。最后，如果能支持PnP且匹配表中的一个ID，返回0，并且由`device_set_desc()`从表中取得适当的描述行置。

如果程序支持PnP，情况看起来如下：

```
if(pnperror != 0)
    return pnperror;
```

对于不支持PnP的程序不需要特殊处理，因为程序会调用空的PnP ID表，且如果在PnP上使用会得到ENXIO。

探测例程通常至少需要某些最少量的资源，如I/O端口号，来调用并探测。对于不同的硬件，程序可能会自己调用其他必需的资源。PnP中的所有资源由PnP子系统先配置，因此程序不需要自己调用它们。

通常探测所需要的最少信息就是端口号。然后某些分配从哪里开始：

```
sc-port0 = bus_get_resource_start(dev,
    SYS_RES_IOPORT, 0 /*rid*/); if(sc-port0 == 0) return ENXIO;
```

低端口地址被保存在softc中，以便将来使用。如果需要经常使用低端口，每次调用资源函数将会慢得无法忍受。如果我们没有得到低端口，返回即可。相反，一些低端口程序相当聪明，会探测所有可能的低端口，如下：

```

/* 此处所有可能的基I/O端口地址表 */
static struct xxx_allports {
    u_short port; /* 端口地址 */
    short used; /* 旗：此端口是否已被其他单元使用 */
} xxx_allports = {
    { 0x300, 0 },
    { 0x320, 0 },
    { 0x340, 0 },
    { 0, 0 } /* 表结束 */
};

...
int port, i;
...

port = bus_get_resource_start(dev, SYS_RES_IOPORT, 0 /*rid*/);
if(port !=0 ) {
    for(i=0; xxx_allports[i].port!=0; i++) {
        if(xxx_allports[i].used || xxx_allports[i].port != port)
            continue;

        /* 到了 */
        xxx_allports[i].used = 1;
        /* 在已知端口上探 */
        return xxx_really_probe(dev, port);
    }
    return ENXIO; /* 端口无法用或已被使用 */
}

/* 在需要猜端口的时候才会到这儿 */
for(i=0; xxx_allports[i].port!=0; i++) {
    if(xxx_allports[i].used)
        continue;

    /* 已被使用 - 即使我在此端口什么也没有
     * 至少我以后不会再次探
     */
    xxx_allports[i].used = 1;

    error = xxx_really_probe(dev, xxx_allports[i].port);
    if(error == 0) /* 在那个端口找到一个 */
        return 0;
}
/* 探所有可能的地址，但没有可用的 */
return ENXIO;

```

当然，做些事情通常使用程序的 **identify()** 例程。但可能有一个正当的理由来说明为什么在函数 **probe()** 中完成更好：如果 **probe()** 会做一些其他敏感操作。探例程按旗 **sensitive** 排序：敏感首先被探测，然后是其他。但 **identify()** 例程在所有探测之前被用，因此它不会考敏感并可能乱些。

在，我得到起始端口以后就需要置端口数（PnP除外），因为内核在配置文件中没有这个信息。

```
if(pnperror /* 只非PnP */
    bus_set_resource(dev, SYS_RES_IOPORT, 0, sc-port0,
    XXX_PORT_COUNT)0)
    return ENXIO;
```

最后分配并激活一片端口地址空间（特殊start和end意思是“使用我通过bus\_set\_resource() 置的那些”）：

```
sc-port0_rid = 0;
sc-port0_r = bus_alloc_resource(dev, SYS_RES_IOPORT,
port0_rid,
/*start*/ 0, /*end*/ ~0, /*count*/ 0, RF_ACTIVE);

if(sc-port0_r == NULL)
    return ENXIO;
```

在可以端口映射的寄存器后，我就可以以某种方式向写入数据并是否如我期望的那样作出反应。如果没有，表明可能其他的在个地址上，或者个地址上根本没有。

通常程序直到接例程才会置中断处理函数。之前我替代以模式进行探测，超以DELAY()。探测例程必须不能永久挂起，上的所有等待必须在超内完成。如果不在段内，可能出故障或配置，程序必须返回，当定超间隔，一些例外以保可靠：尽管假定DELAY()在任何机器上都延相同数量的，但随具体CPU的不同，此函数是有一定的差幅度。

如果探测例程真的想中断是否真的工作，它可以也配置和探测中断。但不建议。

```
/* 以重依赖于具体的方式 */
if(error = xxx_probe_ports(sc))
    goto bad; /* 返回前放源 */
```

依于所的切型号，函数xxx\_probe\_ports()也可能置描述。但如果只支持一型号，也可以硬的形式完成。当然，于PnP，PnP支持从表中自置描述。

```
if(pnperror)
    device_set_desc(dev, "Our device model 1234");
```

探测例程当或者通过取配置寄存器来所有源的，或者由用式置。我将假定一个板上内存的例子。探测例程当尽可能是非入式的，分配和其余源功能性的工作就可以更好地留接例程来做。

内存地址可以在内核配置文件中指定，或者某些可以在非易失性配置寄存器中先配置。如果做法均可用却不同，那当用？可能用在内核配置文件中明置地址，但他知道自己在干什么，当先使用个。一个的例子可能是的：

```

/* 首先拿出配置地址 */
sc-mem0_p = bus_get_resource_start(dev, SYS_RES_MEMORY, 0 /*rid*/);
if(sc-mem0_p == 0) { /* 没有, 用没指定 */
    sc-mem0_p = xxx_read_mem0_from_device_config(sc);

    if(sc-mem0_p == 0)
        /* 从配置寄存器也到不了哪儿 */
        goto bad;
} else {
    if(xxx_set_mem0_address_on_device(sc) 0)
        goto bad; /* 不支持那地址 */
}

/* 就像端口, 置内存大小,
 * 于某些, 内存大小不是常数,
 * 而当从配置寄存器中取, 以不同的不同型号
 * 一个用把内存大小置“msize”配置源,
 * 由ISA自理
 */
if(pnperror) { /*非PnP */
    sc-mem0_size = bus_get_resource_count(dev, SYS_RES_MEMORY, 0 /*rid*/);
    if(sc-mem0_size == 0) /* 用没有指定 */
        sc-mem0_size = xxx_read_mem0_size_from_device_config(sc);

    if(sc-mem0_size == 0) {
        /* 假定是也非常老的一型号, 没有自配置特性,
         * 用也没有偏好置, 因此假定最低要求的情况
         * (当然, 真将根据程序而不同)
         */
        sc-mem0_size = 8*1024;
    }

    if(xxx_set_mem0_size_on_device(sc) 0)
        goto bad; /*不支持那个大小 */

    if(bus_set_resource(dev, SYS_RES_MEMORY, /*rid*/0,
                        sc-mem0_p, sc-mem0_size)0)
        goto bad;
} else {
    sc-mem0_size = bus_get_resource_count(dev, SYS_RES_MEMORY, 0 /*rid*/);
}

```

似, 很容易出IRQ和DRQ所用的源。

如果一切正常, 然后就可以放所有源并返回成功。

```

xxx_free_resources(sc);
return 0;

```

最后，处理棘手情况。所有资源应当在返回前被释放。我利用其中一个  
以前被零化，因此我能看出是否分配了某些资源：如果分配有些资源的描述符非零。

事由：softc在xxx我

```
bad:  
  
    xxx_free_resources(sc);  
    if(error)  
        return error;  
    else /* 一切未知 */  
        return ENXIO;
```

是完整的探例程。资源的释放从多个地方完成，因此将它放到一个函数中，看起来可能像下面的样子：

```

static void
    xxx_free_resources(sc)
        struct xxx_softc *sc;
{
    /* 释放所有资源，如果非0释放 */

    /* 中断处理函数 */
    if(sc-intr_r) {
        bus_teardown_intr(sc-dev, sc-intr_r, sc-intr_cookie);
        bus_release_resource(sc-dev, SYS_RES_IRQ, sc-intr_rid,
            sc-intr_r);
        sc-intr_r = 0;
    }

    /* 我分配的所有内存 */
    if(sc-data_p) {
        bus_dmamap_unload(sc-data_tag, sc-data_map);
        sc-data_p = 0;
    }
    if(sc-data) { /* sc-data_map等于0有可能合法 */
        /* the map will also be freed */
        bus_dmamem_free(sc-data_tag, sc-data, sc-data_map);
        sc-data = 0;
    }
    if(sc-data_tag) {
        bus_dma_tag_destroy(sc-data_tag);
        sc-data_tag = 0;
    }

    ... 如果有，释放其他的映射和 ...

    if(sc-parent_tag) {
        bus_dma_tag_destroy(sc-parent_tag);
        sc-parent_tag = 0;
    }

    /* 释放所有资源 */
    if(sc-mem0_r) {
        bus_release_resource(sc-dev, SYS_RES_MEMORY, sc-mem0_rid,
            sc-mem0_r);
        sc-mem0_r = 0;
    }
    ...
    if(sc-port0_r) {
        bus_release_resource(sc-dev, SYS_RES_IOPORT, sc-port0_rid,
            sc-port0_r);
        sc-port0_r = 0;
    }
}

```

## 9.9. xxx\_isa\_attach

如果探例程返回成功并且系统接那个程序， 接例程 将程序接到系。如果探例程返回0， 接例程期望 接收完整的softc，此由探例程置。同时，如果探例程 返回0，它可能期望0个接例程当在将来的某点被用。如果 探例程返回0，程序可能不会作此假。

如果成功完成， 接例程返回0， 否返回000。

接例程的跟探例程相似， 将一些常用数据取到一些更容易的量中。

```
struct xxx_softc *sc = device_get_softc(dev);
int unit = device_get_unit(dev);
int error = 0;
```

然后分配并激活所需源。由于端口通常在从探返回前就 被放，因此需要重新分配。我希望探例程已当地置了 所有的源，并将它保存在softc中。如果探例程留下了 一些被分配的源，就不需要再次分配（重新分配被0000）。

```
sc-port0_rid = 0;
sc-port0_r = bus_alloc_resource(dev, SYS_RES_IOPORT, port0_rid,
    /*start*/ 0, /*end*/ ~0, /*count*/ 0, RF_ACTIVE);

if(sc-port0_r == NULL)
    return ENXIO;

/* 板上内存 */
sc-mem0_rid = 0;
sc-mem0_r = bus_alloc_resource(dev, SYS_RES_MEMORY, mem0_rid,
    /*start*/ 0, /*end*/ ~0, /*count*/ 0, RF_ACTIVE);

if(sc-mem0_r == NULL)
    goto bad;

/* 取得虚地址 */
sc-mem0_v = rman_get_virtual(sc-mem0_r);
```

DMA求通道(DRQ)以相似方式被分配。使用 isa\_dma\*()函数族行初始化。例如：

```
isa_dmacascade(sc-drq0);
```

中断求(IRQ)有点特殊。除了分配以外，程序的中断理 函数也当与它。在古老的ISA程序中，由系统中断理 函数的参量是元号。但在代程序中，按照定，建000 指向softc的指。一个很重要的原因在于当softc被分配后， 从softc取得元号很容易，而从元号取得softc很困。同时，一个 定也使得用于不同的用程序看起来一，并允它共享代： 000有其自己的探， 接， 分和其他相的例程，而它之可以共享大的程序代。

```
sc-intr_rid = 0;
sc-intr_r = bus_alloc_resource(dev, SYS_RES_MEMORY, intr_rid,
    /*start*/ 0, /*end*/ ~0, /*count*/ 0, RF_ACTIVE);

if(sc-intr_r == NULL)
    goto bad;

/*
 * 假定XXX_INTR_TYPE的定义依于程序的类型,
 * 例如INTR_TYPE_CAM用于CAM的程序
 */
error = bus_setup_intr(dev, sc-intr_r, XXX_INTR_TYPE,
    (driver_intr_t *) xxx_intr, (void *) sc, intr_cookie);
if(error)
    goto bad;
```

如果程序需要与内存进行DMA，内存应当按前述方式分配：

```

error=bus_dma_tag_create(NULL, /*alignment*/ 4,
    /*boundary*/ 0, /*lowaddr*/ BUS_SPACE_MAXADDR_24BIT,
    /*highaddr*/ BUS_SPACE_MAXADDR, /*filter*/ NULL, /*filterarg*/ NULL,
    /*maxsize*/ BUS_SPACE_MAXSIZE_24BIT,
    /*nsegments*/ BUS_SPACE_UNRESTRICTED,
    /*maxsegsz*/ BUS_SPACE_MAXSIZE_24BIT, /*flags*/ 0,
    parent_tag);
if(error)
    goto bad;

/* 很多东西是从父继承而来
 * 假设sc-data指向存储共享数据的区域，例如一个缓冲区可能是：
 * struct {
 *     u_short rd_pos;
 *     u_short wr_pos;
 *     char bf[XXX_RING_BUFFER_SIZE]
 * } *data;
 */
error=bus_dma_tag_create(sc-parent_tag, 1,
    0, BUS_SPACE_MAXADDR, 0, /*filter*/ NULL, /*filterarg*/ NULL,
    /*maxsize*/ sizeof(* sc-data), /*nsegments*/ 1,
    /*maxsegsz*/ sizeof(* sc-data), /*flags*/ 0,
    data_tag);
if(error)
    goto bad;

error = bus_dmamem_alloc(sc-data_tag, data, /* flags*/ 0,
    data_map);
if(error)
    goto bad;

/* 在data_p的情况下，xxx_alloc_callback()只是将物理地址
 * 保存到作为其参数传递的指针中。
 * 参看于内存映射一节中的内容。
 * 其可以像这样：
 *
 * static void
 * xxx_alloc_callback(void *arg, bus_dma_segment_t *seg,
 *     int nseg, int error)
 * {
 *     *(bus_addr_t *)arg = seg[0].ds_addr;
 * }
 */
bus_dmamap_load(sc-data_tag, sc-data_map, (void *)sc-data,
    sizeof (* sc-data), xxx_alloc_callback, (void *) data_p,
    /*flags*/0);

```

分配了所有的资源后，应当被初始化。初始化可能包括所有特性，确保它起作用。

```
if(xxx_initialize(sc)  0)
    goto bad;
```

功能子系将自己在控制台上打印由探测例程设置的描述。但  
外信息，也是可能的，例如：

如果程序想打印一些于的

```
device_printf(dev, "has on-card FIFO buffer of %d bytes\n", sc-fifosize);
```

如果初始化例程遇到任何错误，建返回之前打印有信息。

接例程的最后一是将功能子系接到内核中的功能子系。完成

个的精方式依于程序的型：字符、

、网、CAM SCSI等等。

个的精方式依于程序的型：字符、

如果所有均工作正常返回成功。

```
error = xxx_attach_subsystem(sc);
if(error)
    goto bad;

return 0;
```

最后，处理棘手情况。返回前，所有源当被取消分配。

我利用一个事：softc我

之前被零化，因此我能看出是否分配了某些源：如果分配它描述符非零。

我利用一个事：softc我

```
bad:

xxx_free_resources(sc);
if(error)
    return error;
else /* exact error is unknown */
    return ENXIO;
```

就是接例程的全部。

## 9.10. xxx\_isa\_detach

如果程序中存在个函数，且程序被可加载模块，该程序具有被卸载的能力。如果硬件支持拔，  
是一个很重要的特性。但ISA不支持拔，因此这个特性于ISA不是特别重要。卸载程序的能力可能在有用，但很多情况下只有在老版本的程序莫名其妙地住系  
的情况下才需要安装新版本的程序，并且无论如何都需要重，使得花精力写分例程有些不得。  
一个宣称卸允在用于生的机器上升程序的点看起来似乎更多的只是理而已。升程序是一危  
的操作，决不应当在用于生的机器上行（并且当系行于安全模式也是不被允  
的）。然而，出于完整性考，是会提供分例程。

如果程序成功分，分例程返回0，否返回。

分离是直接的。要做的第一件事情就是将程序从内核子系分离。如果当前正打着，程序有以下几个：拒分离或者抵制并行分离。用这种方式取决于特定内核子系并行抵制的能力和程序作者的偏好。通常抵制似乎是更好的。

```
struct xxx_softc *sc = device_get_softc(dev);
int error;

error = xxx_detach_subsystem(sc);
if(error)
    return error;
```

下一，程序可能希望硬件到某一致的状态。包括停止任何将要执行的，禁用DMA通道和中断以避免破坏内存。对于大多数程序而言，正是例程所做的，因此如果程序中包括例程，我只要用它就可以了。

```
xxx_isa_shutdown(dev);
```

最后放所有源并返回成功。

```
xxx_free_resources(sc);
return 0;
```

## 9.11. xxx\_isa\_shutdown

当系统要的时候用此例程。通过它使硬件进入某一致的状态。对于大多数ISA而言不需要特殊操作，因此这个函数并非真正必需，因为不管多重设备会被重新初始化。但有些必须按特定顺序，以保在重后能被正确地放到（对于很多使用私有设备的特有用）。很多情况下，在寄存器中禁用DMA和中断，并停止将要执行的是个好主意。一切操作取决于硬件，因此我无法在此说明。

## 9.12. xxx\_intr

当收到来自特定的中断就会用中断处理函数。ISA不支持中断共享（某些特殊情况例外），因此如果中断处理函数被用，几乎可以相信中断是来自其。然而，中断处理函数必须保证寄存器并保中断是由它的产生的。如果不是，中断处理函数应当返回。

ISA程序的旧定是取元号作参数。在已，当用bus\_setup\_intr()新程序接收任何在接例程中他指定的参数。根据新定，它应当是指向softc的指。因此中断处理函数通常像下面那开始：

```
static void
xxx_intr(struct xxx_softc *sc)
{
```

它行在由bus\_setup\_intr()的中断型参数指定的中断先上。意味着禁用所有其他同类型的中断和所有事件中断。

为了避免争，中断处理例程通常写成循环形式：

```
while(xxx_interrupt_pending(sc)) {  
    xxx_process_interrupt(sc);  
    xxx_acknowledge_interrupt(sc);  
}
```

中断处理函数必须只向内核答中断，但不能向中断控制器答，后者由系统管理。

# Chapter 10. PCI

本章将介绍FreeBSD中PCI上的读写程序而提供的机制。

## 10.1. 探测与连接

本节的信息是关于PCI代理如何迭代通过未连接的设备，并查看新加载的kld是否会连接其中一个。

### 10.1.1. 示例程序源代码(mypci.c)

```
/*
 * 与PCI函数进行交互的KLD
 *
 * Murray Stokely
 */

#include sys/param.h          /* kernel.h中使用的定义 */
#include sys/module.h
#include sys/system.h
#include sys/errno.h
#include sys/kernel.h          /* 模块初始化中使用的类型 */
#include sys/conf.h            /* cdevsw等 */
#include sys/uio.h             /* uio等 */
#include sys/malloc.h
#include sys/bus.h              /* pci用到的头、原型 */

#include machine/bus.h
#include sys/rman.h
#include machine/resource.h

#include dev/pci/pcivar.h     /* 为了使用get_pci宏! */
#include dev/pci/pcireg.h

/* softc保存我们自己的数据。 */
struct mypci_softc {
    device_t    my_dev;
    struct cdev *my_cdev;
};

/* 函数原型 */
static d_open_t      mypci_open;
static d_close_t     mypci_close;
static d_read_t      mypci_read;
static d_write_t     mypci_write;

/* 字符入口点 */

static struct cdevsw mypci_cdevsw = {
    .d_version =      D_VERSION,
```

```

    .d_open = mypci_open,
    .d_close = mypci_close,
    .d_read = mypci_read,
    .d_write = mypci_write,
    .d_name = "mypci",
};

/*
 * 在cdevsw例程中，我通过实体cdev中的成员si_drv1指出我的softc。
 * 当我建立/dev/my， 在我的已附着的例程中，
 * 我设置一个变量指向我的softc。
 */
int
mypci_open(struct cdev *dev, int oflags, int devtype, struct thread *td)
{
    struct mypci_softc *sc;

    /* Look up our softc. */
    sc = dev-si_drv1;
    device_printf(sc-my_dev, "Opened successfully.\n");
    return (0);
}

int
mypci_close(struct cdev *dev, int fflag, int devtype, struct thread *td)
{
    struct mypci_softc *sc;

    /* Look up our softc. */
    sc = dev-si_drv1;
    device_printf(sc-my_dev, "Closed.\n");
    return (0);
}

int
mypci_read(struct cdev *dev, struct uio *uio, int ioflag)
{
    struct mypci_softc *sc;

    /* Look up our softc. */
    sc = dev-si_drv1;
    device_printf(sc-my_dev, "Asked to read %d bytes.\n", uio-uio_resid);
    return (0);
}

int
mypci_write(struct cdev *dev, struct uio *uio, int ioflag)
{
    struct mypci_softc *sc;

```

```

/* Look up our softc. */
sc = dev_si_drv1;
device_printf(sc-my_dev, "Asked to write %d bytes.\n", uio-uio_resid);
return (0);
}

/* PCI支持函数 */

/*
 * 将某个置的与个程序支持的相比。
 * 如果相符，置描述字符并返回成功。
 */
static int
mypci_probe(device_t dev)
{
    device_printf(dev, "MyPCI Probe\nVendor ID : 0x%x\nDevice ID : 0x%x\n",
        pci_get_vendor(dev), pci_get_device(dev));

    if (pci_get_vendor(dev) == 0x11c1) {
        printf("We've got the Winmodem, probe successful!\n");
        device_set_desc(dev, "WinModem");
        return (BUS_PROBE_DEFAULT);
    }
    return (ENXIO);
}

/* 只有当探成功才用接函数 */

static int
mypci_attach(device_t dev)
{
    struct mypci_softc *sc;

    printf("MyPCI Attach for : deviceID : 0x%x\n", pci_get_devid(dev));

    /* Look up our softc and initialize its fields. */
    sc = device_get_softc(dev);
    sc-my_dev = dev;

    /*
     * Create a /dev entry for this device. The kernel will assign us
     * a major number automatically. We use the unit number of this
     * device as the minor number and name the character device
     * "mypciunit".
     */
    sc-my_cdev = make_dev(mypci_cdevsw, device_get_unit(dev),
        UID_ROOT, GID_WHEEL, 0600, "mypci%u", device_get_unit(dev));
    sc-my_cdev-si_drv1 = sc;
    printf("Mypci device loaded.\n");
    return (0);
}

```

```

}

/* 分离。 */

static int
mypci_detach(device_t dev)
{
    struct mypci_softc *sc;

    /* Teardown the state in our softc created in our attach routine. */
    sc = device_get_softc(dev);
    destroy_dev(sc->my_cdev);
    printf("Mypci detach!\n");
    return (0);
}

/* 系统周期在sync之后用。 */

static int
mypci_shutdown(device_t dev)
{
    printf("Mypci shutdown!\n");
    return (0);
}

/*
 * 挂起例程。
 */
static int
mypci_suspend(device_t dev)
{
    printf("Mypci suspend!\n");
    return (0);
}

/*
 * 恢复（重新启动）例程。
 */
static int
mypci_resume(device_t dev)
{
    printf("Mypci resume!\n");
    return (0);
}

static device_method_t mypci_methods[] = {
    /* 接口 */
    DEVMETHOD(device_probe,      mypci_probe),

```

```

DEVMETHOD(device_attach,    mypci_attach),
DEVMETHOD(device_detach,   mypci_detach),
DEVMETHOD(device_shutdown, mypci_shutdown),
DEVMETHOD(device_suspend,  mypci_suspend),
DEVMETHOD(device_resume,   mypci_resume),

{ 0, 0 }
};

static devclass_t mypci_devclass;

DEFINE_CLASS_0(mypci, mypci_driver, mypci_methods, sizeof(struct mypci_softc));
DRIVER_MODULE(mypci, pci, mypci_driver, mypci_devclass, 0, 0);

```

## 10.1.2. 示例程序的Makefile

```

# 程序mypci的Makefile

KMOD= mypci
SRCS= mypci.c
SRCS+= device_if.h bus_if.h pci_if.h

.include bsd.kmod.mk

```

如果将上面的源文件和 Makefile 放入一个目录，可以运行 make 编译示例程序。如果有，可以运行 make load 将程序装入到当前正在运行的内核中，而 make unload 可在装入后卸载程序。

## 10.1.3. 更多源

- [PCI Special Interest Group](#)
- PCI System Architecture, Fourth Edition by Tom Shanley, et al.

## 10.2. 源

FreeBSD 从父辈求源提供了一面向对象的机制。几乎所有都是某类型的（PCI, ISA, USB, SCSI 等等）的孩子成员，并且有些需要从他的父辈取源（例如内存段, 中断, 或者 DMA 通道）。

### 10.2.1. 基地址寄存器

为了 PCI 做些有用的事情，需要从 PCI 配置空间取 Base Address Registers (BARs)。取 BAR 的 PCI 特定的被抽象在函数 `bus_alloc_resource()` 中。

例如，一个典型的程序可能在 `attach()` 函数中有些类似于下面的东西：

```

sc-bar0id = PCIR_BAR(0);
sc-bar0res = bus_alloc_resource(dev, SYS_RES_MEMORY, sc-bar0id,
                                0, ~0, 1, RF_ACTIVE);
if (sc-bar0res == NULL) {
    printf("Memory allocation of PCI base register 0 failed!\n");
    error = ENXIO;
    goto fail1;
}

sc-bar1id = PCIR_BAR(1);
sc-bar1res = bus_alloc_resource(dev, SYS_RES_MEMORY, sc-bar1id,
                                0, ~0, 1, RF_ACTIVE);
if (sc-bar1res == NULL) {
    printf("Memory allocation of PCI base register 1 failed!\n");
    error = ENXIO;
    goto fail2;
}
sc-bar0_bt = rman_get_bustag(sc-bar0res);
sc-bar0_bh = rman_get_bushandle(sc-bar0res);
sc-bar1_bt = rman_get_bustag(sc-bar1res);
sc-bar1_bh = rman_get_bushandle(sc-bar1res);

```

两个基址寄存器的句柄被保存在 `softc` 中，以便以后可以使用它向写入。

然后就能使用些句柄与 `bus_space_*` 函数一起写寄存器。例如，程序可能包含如下的快捷函数，用来取板子特定的寄存器：

```

uint16_t
board_read(struct ni_softc *sc, uint16_t address)
{
    return bus_space_read_2(sc-bar1_bt, sc-bar1_bh, address);
}

```

似的，可以用下面的函数写寄存器：

```

void
board_write(struct ni_softc *sc, uint16_t address, uint16_t value)
{
    bus_space_write_2(sc-bar1_bt, sc-bar1_bh, address, value);
}

```

些函数以8位、16位和32位的版本存在，当相地使用 `bus_space_{read|write}_{1|2|4}`。

在 FreeBSD 7.0 和更高版本中，可以用 `bus_*` 函数来代替 `bus_space_*`。`bus_*` 函数使用的参数是 `struct resource *` 指针，而不是 bus tag 和句柄。因此，就可以将 `softc` 中的 bus tag 和 bus 句柄全部成批量去掉，并将 `board_read()` 函数改写为：



```
uint16_t  
board_read(struct ni_softc *sc, uint16_t address)  
{  
    return (bus_read(sc->bar1res, address));  
}
```

### 10.2.2. 中断

中断按照和分配内存资源相似的方式从面向对象的接口代为分配。首先，必须从父接口分配IRQ资源，然后必须设置中断处理函数来处理每个IRQ。

再一次，来自 `attach()` 函数的例子比文字更具说明性。

```
/* 取得IRQ资源 */  
  
sc->irqid = 0x0;  
sc->irqres = bus_alloc_resource(dev, SYS_RES_IRQ, (sc->irqid),  
    0, ~0, 1, RF_SHAREABLE | RF_ACTIVE);  
if (sc->irqres == NULL) {  
    printf("IRQ allocation failed!\n");  
    error = ENXIO;  
    goto fail3;  
}  
  
/* 在我应当设置中断处理函数 */  
  
error = bus_setup_intr(dev, sc->irqres, INTR_TYPE_MISC,  
    my_handler, sc, (sc->handler));  
if (error) {  
    printf("Couldn't set up irq\n");  
    goto fail4;  
}
```

在接口的分例程中必须注意一些问题。必须停止所有的中断流，并移除中断处理函数。一旦 `bus_teardown_intr()` 返回，知道的中断处理函数不会再被调用，并且所有可能已执行了。所有中断处理函数的流程都已经返回。由于此函数可以睡眠，调用此函数时必须不能有任何互斥体。

### 10.2.3. DMA

本节已过时，只是由于历史原因而给出。处理一些旧的不当方法是使用 `bus_space_dma*()` 函数。当更新一个以反映那些用法时，这段就可能被去掉。然而，目前API不断有些变化，因此一旦它被固定下来后，更新一个来反映那些改动就很好了。

在PC上，想执行主控DMA的外设必物理地址，由于  
仍是个。幸的是，有个函数，`vtophys()`可以帮助我。

FreeBSD使用虚内存并且只管理虚地址，

```
#include <vm/vm.h>
#include <vm/pmap.h>

#define vtophys(virtual_address) (...)
```

然而这个解决办法在alpha上有点不一，并且我真正想要的是一个称`vtobus()`的函数。

```
#if defined(__alpha__)
#define vtobus(va)    alpha_XXX_dmamap((vm_offset_t)va)
#else
#define vtobus(va)    vtophys(va)
#endif
```

#### 10.2.4. 取消分配资源

取消`attach()`期间分配的所有资源非常重要。必须小心谨慎，即使在失败的条件下也要保证取消分配那些正确的東西，当它的程序去掉后系统仍然可以使用。

# Chapter 11. 通用方法SCSI控制器

## 11.1. 提要

本文假定读者对FreeBSD的驱动程序和SCSI有大致了解，本文中很多信息是从以下程序中：

- ncr (/sys/pci/ncr.c) 由Wolfgang Staelmeier and Stefan Esser写
- sym (/sys/dev/sym/sym\_raid.c) 由Gerard Roudier写
- aic7xxx (/sys/dev/aic7xxx/aic7xxx.c) 由Justin T. Gibbs写

和从CAM的代码本身（作者 Justin T. Gibbs, /sys/cam/\*）中摘出。当一些解决方法看起来更具创造性，并且基本上是从 Justin T. Gibbs 的代码中一字不差地摘出，我将其称为“recommended”。

本文以代码例子进行说明。尽管有例子中包含很多注释，并且看起来很像真正代码，但它仍然只是代码。代码是通过以一种可理解的方式来展示概念。对于真正的代码，其它方法可能更模块化，并且更加高效。文档也硬件进行抽象，对于那些会模糊我们所要展示的概念的注释，或被包含在作者手册的其他章节中已有描述的注释也做同样的处理。这些通常以注释具有描述性名字的函数、注释或语句的形式展示。幸运的是，具有完整的完整例子，包括所有注释，可以在真正的代码程序中找到。

## 11.2. 通用基础

CAM代表通用方法（Common Access Method）。它以SCSI方式地址I/O。这就允许将通用驱动程序和控制I/O的代码分离来：例如磁盘程序能同时控制SCSI、IDE、且/或任何其他上的磁盘，磁盘程序部分不必修改新的I/O而重写（或拷贝修改）。一个最重要的活动体是：

- 外部模块 - 外部（磁盘，磁带，CD-ROM等）的代码
- SCSI接口模块(SIM) - 连接到I/O，如SCSI或IDE，的主机接口驱动器代码。

外部程序从OS接收请求，将它们转换为SCSI命令序列并将这些SCSI命令送到SCSI接口模块。SCSI接口模块将这些命令发送到硬件（或者如果硬件不是SCSI，而是例如IDE，也要将这些SCSI命令转换为硬件的native命令）。

由于这儿我感兴趣的是编写SCSI驱动器代码，从此开始我将从SIM的角度考虑所有的事情。

典型的SIM代码需要包括如下的CAM相关的文件：

```
#include cam/cam.h
#include cam/cam_ccb.h
#include cam/cam_sim.h
#include cam/cam_xpt_sim.h
#include cam/cam_debug.h
#include cam/scsi/scsi_all.h
```

一个SIM代码必须做的第一件事情是向CAM子系统注册它自己。  
和以后的 xxx\_ 用于指唯一的一个程序名字前缀）期完成。

在代码的xxx\_attach()函数（此  
xxx\_attach()函数自身由系统自配置代码用，我在此不描述这部分代码。

需要好几段来完成：首先需要分配与SIM相关的队列：

```
struct cam_devq *devq;

if(( devq = cam_simq_alloc(SIZE) )==NULL) {
    error; /* 一些清理的代码 */
}
```

此段 `SIZE` 要分配的队列的大小，它能包含的最大请求数目。它是 SIM 程序在 SCSI 上能并行处理的请求数目。一般可以如下估算：

```
SIZE = NUMBER_OF_SUPPORTED_TARGETS * MAX_SIMULTANEOUS_COMMANDS_PER_TARGET
```

下一节我的SIM创建描述符：

```
struct cam_sim *sim;

if(( sim = cam_sim_alloc(action_func, poll_func, driver_name,
    softc, unit, max_dev_transactions,
    max_tagged_dev_transactions, devq) )==NULL) {
    cam_simq_free(devq);
    error; /* 一些清理代码 */
}
```

注意如果我不能创建SIM描述符，我也无法放 `devq`，因为我无法做任何其他事情，而且我想分配内存。

如果SCSI上有多条SCSI总线，每条都需要它自己的 `cam_sim`。

一个有趣的是，如果SCSI有不只一条SCSI总线我应该怎么做，一个需要一个`devq`还是两条SCSI总线？在CAM代码的注释中给出的答案是：任一方式均可，由程序的作者决定。

参数：

- `action_func` - 指向程序 `xxx_action` 函数的指针。

```
static void xxx_action ( struct cam_sim *sim union ccb *ccb);
struct cam_sim *sim, union ccb *ccb ;
```

- `poll_func` - 指向程序 `xxx_poll()` 函数的指针。

```
static void xxx_poll ( struct cam_sim *sim);
struct cam_sim *sim ;
```

- `driver_name` - 程序的名字，例如 "ncr" 或 "wds"。
- `softc` - 指向一个SCSI程序的内部描述符的指针。这个指针以后被程序用来取私有数据。

- unit - 控制器元号，例如于控制器 "wds0" 的此数字将为0。
- max\_dev\_transactions - 无模式下每个SCSI目标的最大并行 (simultaneous) 事务数。这个数一般几乎是等于1，只有非 SCSI 才可能例外。此外，如果程序希望执行一个事务的同准一个事务，可以将其置为2，但似乎不得加属性。
- max\_tagged\_dev\_transactions - 同的东西，但是在模式下。这是SCSI在上起多个事务的方式：一个事务被给予一个唯一的ID，并被送到。当完成某些事务，它将结果同一起送回来，SCSI 配器（和程序）就能知道这个事务完成了。此参数也被设置为最大深度。它取决于SCSI 配器的能力。

最后我注册与我的SCSI配器的SCSI。

```
if(xpt_bus_register(sim, bus_number) != CAM_SUCCESS) {
    cam_sim_free(sim, /*free_devq*/ TRUE);
    error; /* 一些清理代码 */
}
```

如果一条SCSI有一个devq（即，我将有多条的看作多个，每个有一条），口号是00，否SCSI上的条当有不同的号。条需要它自己独的cam\_sim。

之后我的控制器完全挂接到CAM系。在 devq的可以被：在所有以后从CAM出来的用中将以 sim参数，devq可以由它输出。

CAM一些事件提供了框架。有些事件来自底（SIM程序），有些来自外程序，有一些来自CAM子系本身。任何程序都可以某些类型的事件注册回，那些事件发生它就会被通知。

事件的一个典型例子就是定位。个事和事件以 "path" 的方式区分它所作用的。目特定的事件通常在与行事理期生。因此那个事的路径可以被重用 来告此事件（是安全的，因事件路径的拷贝是在事件告例程中行的，而且既不会被 deallocate 也不作一）。在任何时刻，包括中断例程中，分配路径也是安全的，尽管那会致某些意外，并且方法可能存在的一个问题是巧那可能没有空内存。于定位事件，我需要定位包括上所有在内的通配符路径。我就能提前 以后的定位事件建路径，避免以后内存不足的：

```
struct cam_path *path;

if(xpt_create_path(path, /*periph*/NULL,
                    cam_sim_path(sim), CAM_TARGET_WILDCARD,
                    CAM_LUN_WILDCARD) != CAM_REQ_CMP) {
    xpt_bus_deregister(cam_sim_path(sim));
    cam_sim_free(sim, /*free_devq*/TRUE);
    error; /* 一些清理代码 */
}

softc-wpath = path;
softc-sim = sim;
```

正如所看到的，路径包括：

- 外部程序的ID（由于我一个也没有，故此空）
- SIM程序的ID（`cam_sim_path(sim)`）
- 所有的SCSI目标号（`CAM_TARGET_WILDCARD`的意思指 "所有devices"）
- 子所有的SCSI LUN号（`CAM_LUN_WILDCARD`的意思指 "所有LUNs"）

如果程序不能分配一个路径，它将不能正常工作，因此那种情况下我卸除（dismantle）那个SCSI。

我在`softc`中保存路径指以便以后使用。之后我保存sim的（或者如果我愿意，也可以在从`xxx_probe()`退出时）。

就是最低要求的初始化所需要做的一切。为了把事情做正无，剩下下一个。

于SIM程序，有一个特殊感兴趣的事件：何目被找不到。情况下位与个的SCSI商可能是个好主意。因此我个事件向CAM注册一个回。通常类型的请求来请求CAM控制上的CAM操作，请求就被到CAM：（注：参看下面示例代码和原文）

```
struct ccb_setasync csa;

xpt_setup_ccb(csa.ccb_h, path, /*先*/5);
csa.ccb_h.func_code = XPT_SASYNC_CB;
csa.event_enable = AC_LOST_DEVICE;
csa.callback = xxx_async;
csa.callback_arg = sim;
xpt_action((union ccb *)csa);
```

在我看一下`xxx_action()`和`xxx_poll()`的程序入口点。

```
static void xxx_action ( struct cam_sim *sim union ccb *ccb);
struct cam_sim *sim, union ccb *ccb ;
```

CAM子系的请求采取某些操作。Sim描述了请求的SIM，CCB请求本身。CCB代表"CAM Block"。它是很多特定的组合，一个例子某些类型的事描述参数。所有这些例子共享着参数公共部分的CCB部分。（注：一段不很准确，自行参考原文）

CAM既支持SCSI控制器工作于发起者(initiator)("normal")模式，也支持SCSI控制器工作于目标(target)（模式）模式。这儿我只考虑发起者模式有部分。

定义了几个函数和宏（句，方法）来从sim中公共数据：

- `cam_sim_path(sim)` - 路径ID（参上面）
- `cam_sim_name(sim)` - sim的名字
- `cam_sim_softc(sim)` - 指向softc（程序私有数据）的指
- `cam_sim_unit(sim)` - 元号
- `cam_sim_bus(sim)` - ID

为了，`xxx_action()`可以使用些函数得到元号和指向它的softc的指。

请求的类型被存储在 `ccb->ccb_h.func_code`。因此，通常 `xxx_action()` 由一个大的switch语句成：

```
struct xxx_softc *softc = (struct xxx_softc *) cam_sim_softc(sim);
struct ccb_hdr *ccb_h = ccb->ccb_h;
int unit = cam_sim_unit(sim);
int bus = cam_sim_bus(sim);

switch(ccb_h->func_code) {
    case ...:
        ...
    default:
        ccb_h->status = CAM_REQ_INVALID;
        xpt_done(ccb);
        break;
}
```

从default case语句部分可以看出（如果收到未知命令），命令的返回值被置到 `ccb->ccb_h.status` 中，并且通常用 `xpt_done(ccb)` 将整个CCB返回到CAM中。

`xpt_done()`不必从 `xxx_action()` 中调用：例如I/O请求可以在SIM程序和/或它的SCSI控制器中排队。（注：它指I/O请求？）然后，当POST(post)一个中断信号，指示自此请求的理由已结束，`xpt_done()` 可以从中断理由例程中被调用。

通常上，CCB状态不是通过一个返回值，而是始终有某种状态。CCB被调用 `xxx_action()` 例程前，其取得状态 `CAM_REQ_INPROG`，表示其正在执行中。`/sys/cam/cam.h` 中定义了数量级的状态，它能非常完全地表示请求的状态。更有趣的是，状态上是一个枚举状态（低6位）和一些可能出现的附加“或”（似）旗位（高位）的“位或(bitwise or)”。枚举会在以后更详细地。它的一些可以在概要（Errors Summary section）中找到。可能的状态旗：

- `CAM_DEV_QFRZN` - 当处理CCB时，如果SIM程序得到一个重叠（例如，程序不能完成或违反了SCSI）, 它将调用 `xpt_freeze_simq()` 请求队列，把此队列的其他已入但尚未被处理的CCB返回到CAM队列，然后将所有CCB置为旗并调用 `xpt_done()`。这个旗将使得CAM子系统处理后解队列。
- `CAM_AUTOSNS_VALID` - 如果返回条件，且CCB中未置旗 `CAM_DIS_AUTOSENSE`，SIM程序必须自行REQUEST SENSE命令来从队列抽取sense（扩展信息）数据。如果成功，sense数据将被保存在CCB中且置此旗。
- `CAM_RELEASE_SIMQ` - 似乎 `CAM_DEV_QFRZN`，但用于SCSI控制器自身出错（或资源短缺）的情况。此后控制器的所有请求将被 `xpt_freeze_simq()` 停止。SIM程序克服短缺情况，并返回置了此旗的CCB通知CAM后，控制器队列将被重新启动。
- `CAM_SIM_QUEUED` - 当SIM将一个CCB放入其请求队列时当置此旗（或当CCB出队但尚未返回CAM时去掉）。在此旗没有在CAM代码的任何地方使用，因此其目的纯粹用于判断）。

函数 `xxx_action()` 不允许睡眠，因此所有源队列的所有同必须通过SIM或队列来完成。除了前述的旗外，CAM子系统提供了函数 `xpt_release_simq()` 和 `xpt_release_devq()` 来直接解队列，而不必将CCB移到CAM。

CCB部分包含如下字段：

- `path` - 请求的路径ID

- *target\_id* - 目标设备ID
- *target\_lun* - 目标LUN ID
- *timeout* - 命令的超时间隔，以毫秒为单位
- *timeout\_ch* - 一个SIM程序存储器超时管理函数的方便之所（CAM子系统自身并不以此作任何假设）
- *flags* - 有关请求的各个信息位
- *spriv\_ptr0, spriv\_ptr1* - SIM程序保留私用的字段（例如接到SIM总线或SIM私有控制口）；通常上，它们会合存在：spriv\_ptr0和spriv\_ptr1具有void \*型(void \*), spriv\_field0和spriv\_field1具有unsigned long型，sim\_priv.entries[0].bytes和sim\_priv.entries[1].bytes与它们合的其他形式大小一致的字节数，sim\_priv.bytes是一个倍大小的数据

使用CCB的SIM私有字段的建方法是将它定一些有意的名字，  
的名字，就像下面这样：

```
#define ccb_some_meaningful_name    sim_priv.entries[0].bytes
#define ccb_hcb spriv_ptr1 /* 用于硬件控制 */
```

最常的起始模式的请求是：

- *XPT\_SCSI\_IO* - 执行I/O操作

结合ccb的"struct ccb\_scsio csio"示例用于参数。它是：

- *cdb\_io* - 指向SCSI命令缓冲区的指针或缓冲区本身
- *cdb\_len* - SCSI命令长度
- *data\_ptr* - 指向数据缓冲区的指针（如果使用分散/集中会更简单）
- *dxfer\_len* - 待传输数据的长度
- *sglist\_cnt* - 分散/集中段的数量
- *scsi\_status* - 返回SCSI状态的地方
- *sense\_data* - 命令返回保存SCSI sense信息的缓冲区（情况下，如果没有设置CCB的标志CAM\_DIS\_AUTOSENSE，假定SIM程序会自行REQUEST SENSE命令）
- *sense\_len* - 缓冲区的长度（如果正好大于sense\_data的大小，SIM程序必须悄悄地采用较小的）（注：一点改动，参考原文及代码）
- *resid, sense\_resid* - 如果数据或SCSI sense返回，它们就是返回的剩余（未读）数据的数量。它看起来并不是特别有意，因此当很计算的情况下（例如，读取SCSI控制器FIFO缓冲区中的字节数），使用近似也完全可以。对于成功完成的IO，它们必须被置为0。
- *tag\_action* - 使用的IO的操作：

  - CAM\_TAG\_ACTION\_NONE - 操作不使用IO
  - MSG\_SIMPLE\_Q\_TAG, MSG\_HEAD\_OF\_Q\_TAG, MSG\_ORDERED\_Q\_TAG - 等于适当的IO信息（/sys/cam/scsi/scsi\_message.h）；输出类型，SIM程序必须自己设置

处理请求的通常如下：

要做的第一件事情是确保可能的争条件，确保命令位于队列中不会被中止：

```
struct ccb_scsio *csio = ccb->csio;

if ((ccb_h->status & CAM_STATUS_MASK) != CAM_REQ_INPROG) {
    xpt_done(ccb);
    return;
}
```

我也不知道我的控制器完全支持什么：

```
if(ccb_h->target_id > OUR_MAX_SUPPORTED_TARGET_ID
|| cch_h->target_id == OUR_SCSI_CONTROLLERS_OWN_ID) {
    ccb_h->status = CAM_TID_INVALID;
    xpt_done(ccb);
    return;
}
if(ccb_h->target_lun > OUR_MAX_SUPPORTED_LUN) {
    ccb_h->status = CAM_LUN_INVALID;
    xpt_done(ccb);
    return;
}
```

然后分配我物理要求所需的数据（如相的硬件控制等）。如果我不能分配SIM队列，  
下我有一个挂起的操作，返回CCB，要求CAM将CCB重新入。以后当源可用，必通过  
返回其状态中置CAM\_SIMQ\_RELEASE位的ccb来解SIM队列。否，如果所有正常，  
将CCB与硬件控制（HCB）连接，并将其志已入。

```
struct xxx_hcb *hcb = allocate_hcb(softc, unit, bus);

if(hcb == NULL) {
    softc->flags |= RESOURCE_SHORTAGE;
    xpt_freeze_simq(sim, /*count*/1);
    ccb_h->status = CAM_REQUEUE_REQ;
    xpt_done(ccb);
    return;
}

hcb->ccb = ccb; ccb_h->ccb_hcb = (void *)hcb;
ccb_h->status |= CAM_SIM_QUEUEDED;
```

从CCB中提取目数据到硬件控制。是否要求我分配一个，如果是生成一个唯一的并制造SCSI信息。SIM程序也与商定彼此支持的最大速度、同速率和偏移。

```

hcb-target = ccb_h-target_id; hcb-lun = ccb_h-target_lun;
generate_identify_message(hcb);
if( ccb_h-tag_action != CAM_TAG_ACTION_NONE )
    generate_unique_tag_message(hcb, ccb_h-tag_action);
if( !target_negotiated(hcb) )
    generate_negotiation_messages(hcb);

```

然后设置SCSI命令。可以在CCB中以多种有趣的方式指定命令的存取。有些方式由CCB中的旗帜指定。命令缓冲区可以包含在CCB中或者用指针指向，后者情况下指针可以指向物理地址或虚地址。由于硬件通常需要物理地址，因此我选择是将地址映射为物理地址。

不太相同的提示：通常这是通过使用 `vtophys()` 来完成的，但由于特殊的Alpha架构之故，除了PCI（它在占SCSI控制器的大多数）程序向Alpha架构的可移植性，必须替代以 `vtobus()` 来完成。[IMHO 提供两个独特的函数 `vtop()` 和 `ptobus()`，而 `vtobus()` 只是它们的一个替代，做更好得多。] 在寻求物理地址的情况下，返回值有状态 `CAM_REQ_INVALID` 的CCB是可以的，当前的程序就是那样做的。但也可能像一个例子（程序中应当有不正确的更直接做法）中那样 Alpha特定的代码片段。如果需要物理地址也能映射回虚地址，但那代价很大，因此我不那样做。

```

if(ccb_h-flags & CAM_CDB_POINTER) {
    /* CDB is a pointer */
    if(!(ccb_h-flags & CAM_CDB_PHYS)) {
        /* CDB是指向虚地址的 */
        hcb-cmd = vtobus(csio-cdb-io.cdb_ptr);
    } else {
        /* CDB是指向物理地址的 */
        #if defined(__alpha__)
            hcb-cmd = csio-cdb-io.cdb_ptr | alpha_XXX_dmamap_or ;
        #else
            hcb-cmd = csio-cdb-io.cdb_ptr ;
        #endif
    }
} else {
    /* CDB在ccb缓冲区中 */
    hcb-cmd = vtobus(csio-cdb-io.cdb_bytes);
}
hcb-cmdlen = csio-cdb-len;

```

现在是在设置数据的时候了，又一次，可以在CCB中以多种有趣的方式指定数据存取，有些方式由CCB中的旗帜指定。首先我得到数据的方向。最简单的情况是没有数据需要映射的情况：

```

int dir = (ccb_h-flags & CAM_DIR_MASK);

if (dir == CAM_DIR_NONE)
    goto end_data;

```

然后我将数据在一个chunk中（是在分散/集中列表中，并且是

物理地址还是虚地址。

SCSI控制器可能只能管理有限数目有限程度的  
大。如果请求到到个限制我就返回。我  
使用一个特殊函数返回CCB，并在一个地方管理HCB资源短缺。增加chunk的函数是  
的可以参看SCSI命令(CDB)的  
的，此我不入它的。由于地址翻  
理的描述。如果某些物体于特定的太困难或不可能，返回状态 **CAM\_REQ\_INVALID** 是可以的。  
上，在的CAM代码中似乎儿也没有使用分散/集中能力。但至少必须一个非分散虚  
冲区的情况，CAM中情况用得很多。

```

int rv;

initialize_hcb_for_data(hcb);

if((!(ccb_h->flags & CAM_SCATTER_VALID)) {
    /* 没有散列区 */
    if(!(ccb_h->flags & CAM_DATA_PHYS)) {
        rv = add_virtual_chunk(hcb, csio-data_ptr, csio-dxfer_len, dir);
    }
} else {
    rv = add_physical_chunk(hcb, csio-data_ptr, csio-dxfer_len,
dir);
}
} else {
    int i;
    struct bus_dma_segment *segs;
    segs = (struct bus_dma_segment *)csio-data_ptr;

    if ((ccb_h->flags & CAM_SG_LIST_PHYS) != 0) {
        /* SG列表指针是物理的 */
        rv = setup_hcb_for_physical_sg_list(hcb, segs, csio-sglist_cnt);
    } else if (!(ccb_h->flags & CAM_DATA_PHYS)) {
        /* SG缓冲区指针是虚的 */
        for (i = 0; i < csio-sglist_cnt; i++) {
            rv = add_virtual_chunk(hcb, segs[i].ds_addr,
                segs[i].ds_len, dir);
            if (rv != CAM_REQ_CMP)
                break;
        }
    } else {
        /* SG缓冲区指针是物理的 */
        for (i = 0; i < csio-sglist_cnt; i++) {
            rv = add_physical_chunk(hcb, segs[i].ds_addr,
                segs[i].ds_len, dir);
            if (rv != CAM_REQ_CMP)
                break;
        }
    }
}

if(rv != CAM_REQ_CMP) {
    /* 如果成功添加了一chunk, 我们希望add_*_chunk()函数返回
     * CAM_REQ_CMP, 如果要求太大 (太多字节或太多chunks)
     * 则返回CAM_REQ_TOO_BIG, 其他情况下返回CAM_REQ_INVALID。
     */
    free_hcb_and_ccb_done(hcb, ccb, rv);
    return;
}
end_data:

```

如果一个CCB不允许断开连接, 我们就将信息到hcb :

```

if(ccb_h->flags & CAM_DIS_DISCONNECT)
    hcb_disable_disconnect(hcb);

```

如果控制器能完全自己执行REQUEST SENSE命令，它也应当将标志CAM\_DIS\_AUTOSENSE的值置为0。可以在CAM子系统不想REQUEST SENSE时阻止自REQUEST SENSE。

剩下的唯一事情是设置超时，将我自己的hcb置为硬件并返回，余下的由中断处理函数（或超时处理函数）完成。

```

ccb_h->timeout_ch = timeout(xxx_timeout, (caddr_t) hcb,
    (ccb_h->timeout * hz) / 1000); /* 将毫秒转换为滴答数 */
put_hcb_into_hardware_queue(hcb);
return;

```

这儿是返回CCB的函数的一个可能实现：

```

static void
free_hcb_and_ccb_done(struct xxx_hcb *hcb, union ccb *ccb, u_int32_t
status)
{
    struct xxx_softc *softc = hcb->softc;

    ccb->ccb_h.ccb_hcb = 0;
    if(hcb != NULL) {
        untimeout(xxx_timeout, (caddr_t) hcb, ccb->ccb_h.timeout_ch);
        /* 我们要释放hcb，因此资源短缺也就不存在了 */
        if(softc->flags & RESOURCE_SHORTAGE) {
            softc->flags = ~RESOURCE_SHORTAGE;
            status |= CAM_RELEASE_SIMQ;
        }
        free_hcb(hcb); /* 同时从任何内部列表中移除hcb */
    }
    ccb->ccb_h.status = status |
        (ccb->ccb_h.status & ~(CAM_STATUS_MASK | CAM_SIM_QUEUED));
    xpt_done(ccb);
}

```

- *XPT\_RESET\_DEV* - 送SCSI "BUS DEVICE RESET"消息到所有

除了外部CCB中没有数据外，其中最令人感兴趣的参数是target\_id。由于控制器硬件，硬件控制器就像XPT\_SCSI\_IO请求中那样被创建（参看XPT\_SCSI\_IO请求的描述）并被送到控制器，或者立即编程向SCSI控制器发送RESET消息到所有，或者请求可能只是不被支持（并返回状态CAM\_REQ\_INVALID）。而且请求完成后，所有的已断开连接(disconnected)的事必被中止（可能在中断例程中）。

而且所有的当前商在位会丢失，因此它也可能被清除。

或者清除可能被延时，因为不管所有请求将在下一次请求重新商。

- *XPT\_RESET\_BUS* - 送RESET信号到SCSI

CCB中并不参量，唯一感兴趣的参量是由指向sim的指的SCSI。

最小会忘上所有SCSI商，并返回状态 CAM\_REQ\_CMP。

恰当的上会外位SCSI（可能也位SCSI控制器）并 将所有在硬件列中的和断接的那些正被理的CCB的完成状态 CAM\_SCSI\_BUS\_RESET。像：

```

int targ, lun;
struct xxx_hcb *h, *hh;
struct ccb_trans_settings neg;
struct cam_path *path;

/* SCSI命令位可能会花很长时间，情况下当使用中断或超时来判断
 * 命令是否完成。但为了，我这儿假命令很快。
 */
reset_scsi_bus(softc);

/* 对所有入口的CCB */
for(h = softc->first_queued_hcb; h != NULL; h = hh) {
    hh = h->next;
    free_hcb_and_ccb_done(h, h->ccb, CAM_SCSI_BUS_RESET);
}

/* 商的（清除操作后的）干净，我宣告一个 */
neg.bus_width = 8;
neg.sync_period = neg.sync_offset = 0;
neg.valid = (CCB_TRANS_BUS_WIDTH_VALID
             | CCB_TRANS_SYNC_RATE_VALID | CCB_TRANS_SYNC_OFFSET_VALID);

/* 对所有断开连接的CCB和干净的商（注：干净=clean） */
for(targ=0; targ = OUR_MAX_SUPPORTED_TARGET; targ++) {
    clean_negotiations(softc, targ);

    /* 如果可能宣告事件 */
    if(xpt_create_path(path, /*periph*/NULL,
                        cam_sim_path(sim), targ,
                        CAM_LUN_WILDCARD) == CAM_REQ_CMP) {
        xpt_async(AC_TRANSFER_NEG, path, neg);
        xpt_free_path(path);
    }

    for(lun=0; lun = OUR_MAX_SUPPORTED_LUN; lun++)
        for(h = softc->first_discon_hcb[targ][lun]; h != NULL; h = hh) {
            hh=h->next;
            free_hcb_and_ccb_done(h, h->ccb, CAM_SCSI_BUS_RESET);
        }
}

ccb->ccb_h.status = CAM_REQ_CMP;
xpt_done(ccb);

/* 告事件 */
xpt_async(AC_BUS_RESET, softc->wpath, NULL);
return;

```

将SCSI命令位作为函数来可能是好主意，因为如果事情出了差错，它会被超时函数作最后的宣告重用。

- XPT\_ABORT - 中止指定的CCB

参数在包含ccb的例"struct ccb\_abort cab" 中。其中唯一的参数字段是：

*abort\_ccb* - 指向被中止的ccb的指针

如果不支持中断就返回CAM\_UA\_ABORT。这也是最小化所有可用的易方式，任何情况下都返回CAM\_UA\_ABORT。

困难方式才是真正地真正地要求。首先用到SCSI事例的中止：

```
struct ccb *abort_ccb;
abort_ccb = ccb->cab.abort_ccb;

if(abort_ccb->ccb_h.func_code != XPT_SCSI_IO) {
    ccb->ccb_h.status = CAM_UA_ABORT;
    xpt_done(ccb);
    return;
}
```

然后需要在我例的例中找到一个CCB。可以通过遍历所有硬件来完成：

控制列表，与一个CCB的控制

```
struct xxx_hcb *hcb, *h;

hcb = NULL;

/* 我假设softc-first_hcb是与此例的所有HCB的列表元素,
 * 包括那些入待处理的、硬件正在处理的和断开接的那些。
 */
for(h = softc->first_hcb; h != NULL; h = h->next) {
    if(h->ccb == abort_ccb) {
        hcb = h;
        break;
    }
}

if(hcb == NULL) {
    /* 我的例中没有的CCB */
    ccb->ccb_h.status = CAM_PATH_INVALID;
    xpt_done(ccb);
    return;
}

hcb=found_hcb;
```

在我来看一下HCB当前的物理状态。它可能或呆在例中正等待  
例中，或已断开接并等待命令结果，或者例上已由硬件完成但尚未被例完成。为了保证我不会与硬件发生争条件，我将HCB中止(aborted)，如果一个

被送到SCSI例，或此例正在  
或者例上已由硬件完成但尚未被例完成。为了保证我不会  
HCB要被送到SCSI例的，

SCSI控制器将会看到一个旗并跳过它。

```
int hstatus;

/* 此示一个函数，有需要特殊操作才能使得一个旗硬件可
 */
set_hcb_flags(hcb, HCB_BEING_ABORTED);

abort_again:

hstatus = get_hcb_status(hcb);
switch(hstatus) {
case HCB_SITTING_IN_QUEUE:
    remove_hcb_from_hardware_queue(hcb);
    /* 行 */
case HCB_COMPLETED:
    /* 是一的情况 */
    free_hcb_and_ccb_done(hcb, abort_ccb, CAM_REQ_ABORTED);
    break;
}
```

如果CCB此正在队列中，我一般会以某硬件相的方式信号 SCSI控制器，通知它我希望中止当前的队列。SCSI控制器会置 SCSI ATTENTION信号，并当目其行后发送ABORT消息。我也位超时，以保目不会永远睡眠。如果命令不能在某个合理的队列，如 10秒内中止，超例程就会行并位整个SCSI队列。由于命令会在某个合理的队列后被中止，因此我可以在将中止请求返回，当作成功完成，并将被中止的CCB队列中止（但没有将它完成）。

```
case HCB_BEING_TRANSFERRED:
    untimeout(xxx_timeout, (caddr_t) hcb, abort_ccb_ccb_h.timeout_ch);
    abort_ccb_ccb_h.timeout_ch =
        timeout(xxx_timeout, (caddr_t) hcb, 10 * hz);
    abort_ccb_ccb_h.status = CAM_REQ_ABORTED;
    /* 要求控制器中止CCB，然后产生一个中断并停止
     */
    if(signal_hardware_to_abort_hcb_and_stop(hcb) == 0) {
        /* 呀，我没有得与硬件的争条件，在我中止
         * 之前它就脱队列，再一次
         * （注：脱=getoff）
         */
        goto abort_again;
    }

    break;
```

如果CCB位于断接的列表中，将它置中止请求，并在硬件队列的前端将它重新插入。位超时，并告中止请求完成。

```

case HCB_DISCONNECTED:
    untimeout(xxx_timeout, (caddr_t) hcb, abort_ccb_ccb_h.timeout_ch);
    abort_ccb_ccb_h.timeout_ch =
        timeout(xxx_timeout, (caddr_t) hcb, 10 * hz);
    put_abort_message_into_hcb(hcb);
    put_hcb_at_the_front_of_hardware_queue(hcb);
    break;
}
ccb_ccb_h.status = CAM_REQ_CMP;
xpt_done(ccb);
return;

```

就是由于ABORT请求的全部，尽管只有一个。由于ABORT消息清除LUN上所有正在执行中的事，我必须将LUN上所有其他活事中止。那当在中断例程中完成，且在中止事之后。

将CCB中止作函数来可能是个很好的主意，因为如果I/O事超一个函数能被重用。唯一的不同是超事将超请求返回状态CAM\_CMD\_TIMEOUT。于是XPT\_ABORT的case句就会很小，像下面这样：

```

case XPT_ABORT:
    struct ccb *abort_ccb;
    abort_ccb = ccb->cab.abort_ccb;

    if(abort_ccb_ccb_h.func_code != XPT_SCSI_IO) {
        ccb_ccb_h.status = CAM_UA_ABORT;
        xpt_done(ccb);
        return;
    }
    if(xxx_abort_ccb(abort_ccb, CAM_REQ_ABORTED) == 0)
        /* no such CCB in our queue */
        ccb_ccb_h.status = CAM_PATH_INVALID;
    else
        ccb_ccb_h.status = CAM_REQ_CMP;
    xpt_done(ccb);
    return;

```

- *XPT\_SET\_TRAN\_SETTINGS* - 式设置SCSI队列设置的

在结合ccb的例"struct ccb\_trans\_setting cts" 中的参数：

- *valid* - 位掩，指示当更新那些设置：
- *CCB\_TRANS\_SYNC\_RATE\_VALID* - 同步速率
- *CCB\_TRANS\_SYNC\_OFFSET\_VALID* - 同步位移
- *CCB\_TRANS\_BUS\_WIDTH\_VALID* - 宽度
- *CCB\_TRANS\_DISC\_VALID* - 置用/禁用断接
- *CCB\_TRANS\_TQ\_VALID* - 置用/禁用队列的排列
- *flags* - 由部分组成，元参数和子操作。元参数：

- *CCB\_TRANS\_DISC\_ENB* - 断开连接
- *CCB\_TRANS\_TAG\_ENB* - 标记的排队
- 子操作:
  - *CCB\_TRANS\_CURRENT\_SETTINGS* - 改变当前的厂商参数
  - *CCB\_TRANS\_USER\_SETTINGS* - 读希望的用户参数
  - *sync\_period, sync\_offset* - 自解约的，如果 $sync\_offset == 0$ 求同步模式
  - *bus\_width* - 宽度，以位（而不是字）



参考原文和源码

支持三个厂商参数，用设置和当前设置。用设置在SIM中程序中  
通常只是一片内存，供上存储（并在以后恢复）其对于参数的一些主要。设置用参数并不会导致重新厂商速率。但当SCSI控制器厂商，它必须永远不能设置高于用参数的，因此它实际上是上限。

当前设置，正如其名字所示，指当前的。改变它意味着下一次必须重新厂商参数。又一次，有些"new current settings" 并没有被假定限制用于上，它只是用作厂商的起始值。此外，它必须受SCSI控制器的能力限制：例如，如果SCSI控制器有8位宽，而要求设置16位宽，在发送前参数必须被悄悄地截取为8位。

一个需要注意的是就是宽度和同一个参数是相对于而言的，而断开连接和用同一参数是相对于lun而言的。

建立的目的是保持3个厂商参数（宽度和同）：

- *user* - 用的一，如上
- *current* - 现生效的那些
- *goal* - 通过"current"参数所求的那些

代码看起来像：

```

struct ccb_trans_settings *cts;
int targ, lun;
int flags;

cts = ccb-cts;
targ = ccb_h-target_id;
lun = ccb_h-target_lun;
flags = cts->flags;
if(flags & CCB_TRANS_USER_SETTINGS) {
    if(flags & CCB_TRANS_SYNC_RATE_VALID)
        softc->user_sync_period[targ] = cts->sync_period;
    if(flags & CCB_TRANS_SYNC_OFFSET_VALID)
        softc->user_sync_offset[targ] = cts->sync_offset;
    if(flags & CCB_TRANS_BUS_WIDTH_VALID)
        softc->user_bus_width[targ] = cts->bus_width;

    if(flags & CCB_TRANS_DISC_VALID) {
        softc->user_tflags[targ][lun] = ~CCB_TRANS_DISC_ENB;
        softc->user_tflags[targ][lun] |= flags & CCB_TRANS_DISC_ENB;
    }
    if(flags & CCB_TRANS_TQ_VALID) {
        softc->user_tflags[targ][lun] = ~CCB_TRANS_TQ_ENB;
        softc->user_tflags[targ][lun] |= flags & CCB_TRANS_TQ_ENB;
    }
}
if(flags & CCB_TRANS_CURRENT_SETTINGS) {
    if(flags & CCB_TRANS_SYNC_RATE_VALID)
        softc->goal_sync_period[targ] =
            max(cts->sync_period, OUR_MIN_SUPPORTED_PERIOD);
    if(flags & CCB_TRANS_SYNC_OFFSET_VALID)
        softc->goal_sync_offset[targ] =
            min(cts->sync_offset, OUR_MAX_SUPPORTED_OFFSET);
    if(flags & CCB_TRANS_BUS_WIDTH_VALID)
        softc->goal_bus_width[targ] = min(cts->bus_width, OUR_BUS_WIDTH);

    if(flags & CCB_TRANS_DISC_VALID) {
        softc->current_tflags[targ][lun] = ~CCB_TRANS_DISC_ENB;
        softc->current_tflags[targ][lun] |= flags & CCB_TRANS_DISC_ENB;
    }
    if(flags & CCB_TRANS_TQ_VALID) {
        softc->current_tflags[targ][lun] = ~CCB_TRANS_TQ_ENB;
        softc->current_tflags[targ][lun] |= flags & CCB_TRANS_TQ_ENB;
    }
}
ccb->ccb_h.status = CAM_REQ_CMP;
xpt_done(ccb);
return;

```

此后当下一次要处理I/O请求，它会判断其是否需要重新协商，  
target\_negotiated(hcb)。它可以如下所示：

例如通过用函数

```

int
target_negotiated(struct xxx_hcb *hcb)
{
    struct softc *softc = hcb->softc;
    int targ = hcb->targ;

    if( softc->current_sync_period[targ] != softc->goal_sync_period[targ]
    || softc->current_sync_offset[targ] != softc->goal_sync_offset[targ]
    || softc->current_bus_width[targ] != softc->goal_bus_width[targ] )
        return 0; /* FALSE */
    else
        return 1; /* TRUE */
}

```

重新商些后, 果必同当前和目的(goal)参数, 于以后的I/O事  
当前和目的参数将相同, 且 target\_negotiated()会返回TRUE。当初始化(在  
xxx\_attach()中) 当前商必被初始化最窄同模式, 目的和当前必被初始化  
控制器所支持的最大。 (注: 原文可能有, 此未改)

- XPT\_GET\_TRAN\_SETTINGS - 得SCSI置的

此操作XPT\_SET\_TRAN\_SETTINGS的逆操作。用通旗CCB\_TRANS\_CURRENT\_SETTINGS  
或CCB\_TRANS\_USER\_SETTINGS (如果同置有程序返回当前置 所求而得的数据填充CCB例  
"struct ccb\_trans\_setting cts". \*

XPT\_CALC\_GEOMETRY - 算磁的(BIOS) (geometry)

参量在合ccb的例"struct ccb\_calc\_geometry ccg" 中:

- block\_size - 入, 以字的大小 (也称扇区)
- volume\_size - 入, 以字的卷大小
- cylinders - 出, 0柱面
- heads - 出, 0磁
- secs\_per\_track - 出, 0磁道的0扇区

如果返回的与SCSI控制器BIOS所想象的差很大, 并且SCSI 控制器上的磁被作可引的, 系  
可能无法。从aic7xxx 程序中摘取的典型示例 :

```

    struct ccb_calc_geometry *ccg;
    u_int32_t size_mb;
    u_int32_t secs_per_cylinder;
    int extended;

    ccg = ccb->ccg;
    size_mb = ccg->volume_size
        / ((1024L * 1024L) / ccg->block_size);
    extended = check_cards_EEPROM_for_extended_geometry(softc);

    if (size_mb > 1024 && extended) {
        ccg->heads = 255;
        ccg->secs_per_track = 63;
    } else {
        ccg->heads = 64;
        ccg->secs_per_track = 32;
    }
    secs_per_cylinder = ccg->heads * ccg->secs_per_track;
    ccg->cylinders = ccg->volume_size / secs_per_cylinder;
    ccb->ccb_h.status = CAM_REQ_CMP;
    xpt_done(ccb);
    return;

```

给出了一般思路，精算依赖于特定BIOS的癖好(quirk)。如果 BIOS没有提供方法置EEPROM 中的"extended translation" 旗，此旗通常当假定等于1。其他流行有：

128 heads, 63 sectors - Symbios控制器
16 heads, 63 sectors - 老式控制器

一些系BIOS和SCSI BIOS会相互争，不定，例如Symbios 875/895 SCSI和Phoenix BIOS的合在系加会出128/63，而当冷或后会是255/63。

- *XPT\_PATH\_INQ* - 路径，句，得SIM程序和SCSI控制器（也称HBA - 主机接口器）的特性。

特性在合ccb的例"struct ccb\_pathinq cpi" 中返回：

- version\_num - SIM程序号，当前所有程序使用1
- hba\_inquiry - 控制器所支持特性的位掩：
- PI\_MDP\_ABLE - 支持MDP消息（来自SCSI3的一些东西？）
- PI\_WIDE\_32 - 支持32位SCSI
- PI\_WIDE\_16 - 支持16位SCSI
- PI\_SDTR\_ABLE - 可以商同速率
- PI\_LINKED\_CDB - 支持链接的命令
- PI\_TAG\_ABLE - 支持命令
- PI\_SOFT\_RST - 支持位（硬位和软位在SCSI中是互斥的）

- target\_sprt - 目標模式支持的旗幟，如果不支持置0
- hba\_misc - 控制器特性：
  - PIM\_SCANHILO - 从高ID到低ID的扫描
  - PIM\_NOREMOVE - 可移除但不包括在扫描之列
  - PIM\_NOINITIATOR - 不支持发起者角色
  - PIM\_NOBUSRESET - 用禁用初始BUS RESET
  - hba\_eng\_cnt - 神秘的HBA引脚数，与一些东西有关，当前是置0
  - vuhiba\_flags - 供应商唯一的旗幟，当前未用
- max\_target - 最大支持的目標ID（8位0007，16位00015，光通道0127）
- max\_lun - 最大支持的LUN ID（老的SCSI控制器07，新的063）
- async\_flags - 安装的物理函数的位掩码，当前未用
- hpath\_id - 子系统中最高的路径ID，当前未用
- unit\_number - 控制器单元号，cam\_sim\_unit(sim)
- bus\_id - 总线号，cam\_sim\_bus(sim)
- initiator\_id - 控制器自己的SCSI ID
- base\_transfer\_speed - 窄带的名义速率，以KB/s，对于SCSI等于3300
- sim\_vid - SIM程序的供应商ID，以0结束的字符串，包含尾0在内的最大长度SIM\_IDLEN
- hba\_vid - SCSI控制器的供应商ID，以0结束的字符串，包含尾0在内的最大长度HBA\_IDLEN
- dev\_name - 程序名字，以0尾的字符串，包含尾0在内的最大长度DEV\_IDLEN，等于cam\_sim\_name(sim)

置字符串字段的建方法是使用strncpy，如：

```
strncpy(cpi->dev_name, cam_sim_name(sim), DEV_IDLEN);
```

置这些后将状态置CAM\_REQ\_CMP，并将CCB完成。

### 11.3. 中断

```
static void xxx_poll ( struct cam_sim *sim);
struct cam_sim *sim ;
```

函数用于当中断子系统不起作用（例如，系统崩坏或正在重建子系统）模中断。CAM子系统在使用函数前置适当的中断。因此它所需做全部的只是调用中断例程（或其他方法，例程来执行操作，而中断例程只是调用例程）。那为什么要麻烦弄出一个独立的函数来？是由于不同的决定。xxx\_poll例程取cam\_sim的指针作参数，而PCI中断例程按照普通决定取的是指向xxx\_softc的指针，ISA中断例程只是取总线号，因此例程一般看起来像：

```
static void
xxx_poll(struct cam_sim *sim)
{
    xxx_intr((struct xxx_softc *)cam_sim_softc(sim)); /* for PCI device */
}
```

or

```
static void
xxx_poll(struct cam_sim *sim)
{
    xxx_intr(cam_sim_unit(sim)); /* for ISA device */
}
```

## 11.4. 事件

如果建立了事件，当指定函数。

```
static void
ahc_async(void *callback_arg, u_int32_t code, struct cam_path *path, void *arg)
```

- callback\_arg - 注册回调提供的
- code - 事件型
- path - 事件作用于其上的
- arg - 事件特定的参数

一型事件的，AC\_LOST\_DEVICE，看起来如下：

```

struct xxx_softc *softc;
struct cam_sim *sim;
int targ;
struct ccb_trans_settings neg;

sim = (struct cam_sim *)callback_arg;
softc = (struct xxx_softc *)cam_sim_softc(sim);
switch (code) {
case AC_LOST_DEVICE:
    targ = xpt_path_target_id(path);
    if(targ = OUR_MAX_SUPPORTED_TARGET) {
        clean_negotiations(softc, targ);
        /* send indication to CAM */
        neg.bus_width = 8;
        neg.sync_period = neg.sync_offset = 0;
        neg.valid = (CCB_TRANS_BUS_WIDTH_VALID
            | CCB_TRANS_SYNC_RATE_VALID | CCB_TRANS_SYNC_OFFSET_VALID);
        xpt_async(AC_TRANSFER_NEG, path, neg);
    }
    break;
default:
    break;
}

```

## 11.5. 中断

中断例程的执行类型依赖于SCSI控制器所接到的外设的类型（PCI, ISA等等）。

SIMOO程序的中断例程运行在中断splcam上。因此应当在程序中使用splcam()来同中断例程与程序剩余部分的活动（由于能察多处理器的程序，事情更要有趣，但此我忽略这种情况）。本文中的代码忽略了同。代码一定不能忽略它。一个的方法就是在插入其他例程的入口点splcam()，并在返回时将它置位，从而用一个大的边界区保护它。为了保中断是会被恢复，可以定义一个包装函数，如：

```

static void
xxx_action(struct cam_sim *sim, union ccb *ccb)
{
    int s;
    s = splcam();
    xxx_action1(sim, ccb);
    splx(s);
}

static void
xxx_action1(struct cam_sim *sim, union ccb *ccb)
{
    ... process the request ...
}

```

方法而且健壮，但它存在的问题是中断可能会被阻塞相当多的事件，会影响系统性能。另一方面，`spl()`函数族有相当高的开销，因此大量的小的I/O区可能也不好。

中断例程处理的情况和其中最重要依赖于硬件。我考虑 "典型(typical)" 情况。

首先，我假设上是否遇到了SCSI忙（可能由同一SCSI总线上的一块SCSI控制器引起）。如果所有输入的和输出接的请求，报告事件并重新初始化我的SCSI控制器。初始化期间控制器不会输出一个忙位，所以我十分重要的，否同一SCSI总线上的两个控制器可能会一直来回地位下去。控制器致命挂起的情况可以在同一地方处理，但可能需要发送RESET信号到SCSI总线来定位与SCSI总线的连接状态。

```
int fatal=0;
struct ccb_trans_settings neg;
struct cam_path *path;

if( detected_scsi_reset(softc)
|| (fatal = detected_fatal_controller_error(softc)) ) {
    int targ, lun;
    struct xxx_hcb *h, *hh;

    /* 处理所有输入的CCB */
    for(h = softc->first_queued_hcb; h != NULL; h = hh) {
        hh = h->next;
        free_hcb_and_ccb_done(h, h->ccb, CAM_SCSI_BUS_RESET);
    }

    /* 要报告的厂商的干 */
    neg.bus_width = 8;
    neg.sync_period = neg.sync_offset = 0;
    neg.valid = (CCB_TRANS_BUS_WIDTH_VALID
                 | CCB_TRANS_SYNC_RATE_VALID | CCB_TRANS_SYNC_OFFSET_VALID);

    /* 处理所有断开接的CCB和厂商 */
    for(targ=0; targ = OUR_MAX_SUPPORTED_TARGET; targ++) {
        clean_negotiations(softc, targ);

        /* report the event if possible */
        if(xpt_create_path(path, /*periph*/NULL,
                           cam_sim_path(sim), targ,
                           CAM_LUN_WILDCARD) == CAM_REQ_CMP) {
            xpt_async(AC_TRANSFER_NEG, path, neg);
            xpt_free_path(path);
        }
    }

    for(lun=0; lun = OUR_MAX_SUPPORTED_LUN; lun++)
        for(h = softc->first_discon_hcb[targ][lun]; h != NULL; h = hh) {
            hh=h->next;
            if(fatal)
                free_hcb_and_ccb_done(h, h->ccb, CAM_UNREC_HBA_ERROR);
            else
                free_hcb_and_ccb_done(h, h->ccb, CAM_SCSI_BUS_RESET);
        }
}
```

```

        }

/* 告事件 */
xpt_async(AC_BUS_RESET, softc->wpath, NULL);

/* 重新初始化可能花很多时间，某些情况下应当由单一中断信号
 * 指示初始化否完成，或在超时 - 但为了我假
 * 初始化真的很快
 */
if(!fatal) {
    reinitialize_controller_without_scsi_reset(softc);
} else {
    reinitialize_controller_with_scsi_reset(softc);
}
schedule_next_hcb(softc);
return;
}

```

如果中断不是由控制器的条件引起的，很可能当前硬件控制输出了。由于硬件，可能有非HCB相关的事件，此我指示不考虑它。然后我分析一个HCB发生了什么：

```

struct xxx_hcb *hcb, *h, *hh;
int hcb_status, scsi_status;
int ccb_status;
int targ;
int lun_to_freeze;

hcb = get_current_hcb(softc);
if(hcb == NULL) {
    /* 或者丢失(stray)的中断，或者某些东西重置，
     * 或者是硬件相关的某些东西
     */
    // 行必要的清理;
    return;
}

targ = hcb->target;
hcb_status = get_status_of_current_hcb(softc);

```

首先我HCB是否完成，如果完成我就返回的SCSI状态。

```

if(hcb_status == COMPLETED) {
    scsi_status = get_completion_status(hcb);
}

```

然后看这个状态是否与REQUEST SENSE命令有关，如果有详细地清理一下它。

```

if(hcb->flags & DOING_AUTOSENSE) {
    if(scsi_status == GOOD) { /* autosense成功 */
        hcb->ccb->ccb_h.status |= CAM_AUTOSNS_VALID;
        free_hcb_and_ccb_done(hcb, hcb->ccb, CAM_SCSI_STATUS_ERROR);
    } else {
        autosense_failed:
            free_hcb_and_ccb_done(hcb, hcb->ccb, CAM_AUTOSENSE_FAIL);
    }
    schedule_next_hcb(softc);
    return;
}

```

否命令自身已完成，把更多注意力放在上。如果一个CCB 没有禁用auto-sense并且命令同sense数据失，执行REQUEST SENSE 命令接收那些数据。

```

hcb->ccb->csio.scsi_status = scsi_status;
calculate_residue(hcb);

if( (hcb->ccb->ccb_h.flags & CAM_DIS_AUTOSENSE)==0
    ( scsi_status == CHECK_CONDITION
        || scsi_status == COMMAND_TERMINATED ) ) {
    /* 启动SENSE */
    hcb->flags |= DOING_AUTOSENSE;
    setup_autosense_command_in_hcb(hcb);
    restart_current_hcb(softc);
    return;
}
if(scsi_status == GOOD)
    free_hcb_and_ccb_done(hcb, hcb->ccb, CAM_REQ_CMP);
else
    free_hcb_and_ccb_done(hcb, hcb->ccb, CAM_SCSI_STATUS_ERROR);
schedule_next_hcb(softc);
return;
}

```

属于商事件的一个典型事情：从SCSI目标（回答我的商企或 由目标起的）接收到的商消息，或目标无法商（拒我的商消息 或不回答它）。

```

switch(hcb_status) {
case TARGET_REJECTED_WIDE_NEG:
    /* 恢复到8-bit */
    softc->current_bus_width[targ] = softc->goal_bus_width[targ] = 8;
    /* 告事件 */
    neg.bus_width = 8;
    neg.valid = CCB_TRANS_BUS_WIDTH_VALID;
    xpt_async(AC_TRANSFER_NEG, hcb->ccb->ccb_h.path_id, neg);
    continue_current_hcb(softc);
    return;
}

```

```

case TARGET_ANSWERED_WIDE_NEG:
{
    int wd;

    wd = get_target_bus_width_request(softc);
    if(wd == softc->goal_bus_width[targ]) {
        /* 可接受的回答 */
        softc->current_bus_width[targ] =
            softc->goal_bus_width[targ] = neg.bus_width = wd;

        /* 告事件 */
        neg.valid = CCB_TRANS_BUS_WIDTH_VALID;
        xpt_async(AC_TRANSFER_NEG, hcb->ccb.ccb_h.path_id, neg);
    } else {
        prepare_reject_message(hcb);
    }
}
continue_current_hcb(softc);
return;
case TARGET_REQUESTED_WIDE_NEG:
{
    int wd;

    wd = get_target_bus_width_request(softc);
    wd = min(wd, OUR_BUS_WIDTH);
    wd = min(wd, softc->user_bus_width[targ]);

    if(wd != softc->current_bus_width[targ]) {
        /* 度改了 */
        softc->current_bus_width[targ] =
            softc->goal_bus_width[targ] = neg.bus_width = wd;

        /* 告事件 */
        neg.valid = CCB_TRANS_BUS_WIDTH_VALID;
        xpt_async(AC_TRANSFER_NEG, hcb->ccb.ccb_h.path_id, neg);
    }
    prepare_width_nego_response(hcb, wd);
}
continue_current_hcb(softc);
return;
}

```

然后我用与前面相同的办法处理auto-sense期可能出的任何。否，我再一次入。

```

if(hcb->flags & DOING_AUTOSENSE)
    goto autosense_failed;

switch(hcb->status) {

```

我考虑的下一事件是末期的接断，一个事件在ABORT或RESET消息之后被看作是正常的，其他情况下是非正常的。

BUS

DEVICE

```
case UNEXPECTED_DISCONNECT:
    if(requested_abort(hcb)) {
        /* 中止影目和LUN上的所有命令，因此将那个目和LUN上的
         * 所有断接的HCB也中止
        */
        for(h = softc-first_discon_hcb[hcb-target][hcb-lun];
            h != NULL; h = hh) {
            hh=h-next;
            free_hcb_and_ccb_done(h, h-ccb, CAM_REQ_ABORTED);
        }
        ccb_status = CAM_REQ_ABORTED;
    } else if(requested_bus_device_reset(hcb)) {
        int lun;

        /* 位影那个目上的所有命令，因此将那个目和LUN上的
         * 所有断接的HCB位
        */

        for(lun=0; lun = OUR_MAX_SUPPORTED_LUN; lun++)
            for(h = softc-first_discon_hcb[hcb-target][lun];
                h != NULL; h = hh) {
                hh=h-next;
                free_hcb_and_ccb_done(h, h-ccb, CAM_SCSI_BUS_RESET);
            }

        /* 送事件 */
        xpt_async(AC_SENT_BDR, hcb-ccb-ccb_h.path_id, NULL);

        /* 是CAM_RESET_DEV请求本身，它完成了 */
        ccb_status = CAM_REQ_CMP;
    } else {
        calculate_residue(hcb);
        ccb_status = CAM_UNEXP_BUSFREE;
        /* request the further code to freeze the queue */
        hcb-ccb-ccb_h.status |= CAM_DEV_QFRZN;
        lun_to_freeze = hcb-lun;
    }
    break;
```

如果目拒接受，我就通知CAM，并返回此LUN的所有命令：

```

case TAGS_REJECTED:
    /* 告事件 */
    neg.flags = 0 ~CCB_TRANS_TAG_ENB;
    neg.valid = CCB_TRANS_TQ_VALID;
    xpt_async(AC_TRANSFER_NEG, hcb->ccb.ccb_h.path_id, neg);

    ccb_status = CAM_MSG_REJECT_REC;
    /* 请求后面的代码冻结 */
    hcb->ccb.ccb_h.status |= CAM_DEV_QFRZN;
    lun_to_freeze = hcb->lun;
    break;

```

然后我再一些其他情况，处理(processing)基本上局限于置CCB状态：

```

case SELECTION_TIMEOUT:
    ccb_status = CAM_SEL_TIMEOUT;
    /* request the further code to freeze the queue */
    hcb->ccb.ccb_h.status |= CAM_DEV_QFRZN;
    lun_to_freeze = CAM_LUN_WILDCARD;
    break;
case PARITY_ERROR:
    ccb_status = CAM_UNCOR_PARITY;
    break;
case DATA_OVERRUN:
case ODD_WIDE_TRANSFER:
    ccb_status = CAM_DATA_RUN_ERR;
    break;
default:
    /*以通用方法处理所有其他*/
    ccb_status = CAM_REQ_CMP_ERR;
    /* 请求后面的代码冻结 */
    hcb->ccb.ccb_h.status |= CAM_DEV_QFRZN;
    lun_to_freeze = CAM_LUN_WILDCARD;
    break;
}

```

然后我是否需要重到需要入队列，直到它得到方可解冻，如果是那就来处理：

```

if(hcb->ccb->ccb_h.status == CAM_DEV_QFRZN) {
    /* 例 */
    xpt_freeze_devq(ccb->ccb_h.path, /*count*/1);

    /* 重新入队所有命令，将它们返回CAM */

    for(h = softc->first_queued_hcb; h != NULL; h = hh) {
        hh = h->next;

        if(targ == h->targ
            (lun_to_freeze == CAM_LUN_WILDCARD || lun_to_freeze == h->lun) )
            free_hcb_and_ccb_done(h, h->ccb, CAM_REQUEUE_REQ);
    }
}

free_hcb_and_ccb_done(hcb, hcb->ccb, ccb_status);
schedule_next_hcb(softc);
return;

```

包括通用中断处理，尽管特定处理器可能需要某些附加处理。

## 11.6. 中断

当执行I/O请求很多事情可能出错。可以在CCB状态中非常尽地  
为了完整起见此列出典型条件的建的一个：

- *CAM\_RESRC\_UNAVAIL* - 某些资源不可用，并且当其可用SIMO程序不能产生事件。一个例子就是某些控制器内部硬件资源，当其可用控制器不会在其产生中断。
- *CAM\_UNCOR\_PARITY* - 产生不可恢复的奇偶校验
- *CAM\_DATA\_RUN\_ERR* - 数据外溢或未预期的数据状态(phase)（不在一个方向上而不是CAM\_DIR\_MASK指定的方向），或由于输出奇数程度
- *CAM\_SEL\_TIMEOUT* - 产生超时（目标不响应）
- *CAM\_CMD\_TIMEOUT* - 生成命令超时（超时函数执行）
- *CAM\_SCSI\_STATUS\_ERROR* - 生成返回的错误
- *CAM\_AUTOSENSE\_FAIL* - 生成且REQUEST SENSE命令失败
- *CAM\_MSG\_REJECT\_REC* - 收到MESSAGE REJECT消息
- *CAM\_SCSI\_BUS\_RESET* - 收到SCSI复位
- *CAM\_REQ\_CMP\_ERR* - 出现"不可能(impossible)"SCSI状态(phase) 或者其他怪异事情，或者如果唯一的信息不可用只是通用错误
- *CAM\_UNEXP\_BUSFREE* - 出现未预期的断开连接
- *CAM\_BDR\_SENT* - BUS DEVICE RESET消息被送到目标
- *CAM\_UNREC\_HBA\_ERROR* - 不可恢复的主机总线适配器错误
- *CAM\_REQ\_TOO\_BIG* - 请求大于控制器太大

- *CAM\_REQQUEUE\_REQ* - 此命令当被重新插入以保持事务的次序性。典型地出现在下列时刻：SIMIO出了当命令队列的命令，并且必须在SIMIO上将命令的其他入口放回到XPT队列。某些典型的典型情况有命令超时、命令超时和其他类似情况。某些情况下出错的命令返回状态来指示，此命令和其他没有被送到命令被重新插入。
- *CAM\_LUN\_INVALID* - SCSI控制器不支持命令中的LUN ID
- *CAM\_TID\_INVALID* - SCSI控制器不支持命令中的目标ID

## 11.7. 超时管理

当HCB的超时期，命令就应当被中止，就像处理XPT\_ABORT命令一样。唯一区别在于被中止的命令的返回状态是CAM\_CMD\_TIMEOUT，而不是CAM\_REQ\_ABORTED（这就是为什么中止的命令最好由函数来完成）。但有一个可能的情况：如果中止命令自己出了麻烦呢？情况下应当对位SCSI命令，就像处理XPT\_RESET\_BUS命令一样（并且将其重命名为函数，从一个地方使用的想法也适用于这儿）。而且如果命令位命令出了问题，我应当对位整个SCSI命令。因此最超时函数看起来像下面这样：

```
static void
xxx_timeout(void *arg)
{
    struct xxx_hcb *hcb = (struct xxx_hcb *)arg;
    struct xxx_softc *softc;
    struct ccb_hdr *ccb_h;

    softc = hcb->softc;
    ccb_h = hcb->ccb->ccb_h;

    if(hcb->flags & HCB_BEING_ABORTED
       || ccb_h->func_code == XPT_RESET_DEV) {
        xxx_reset_bus(softc);
    } else {
        xxx_abort_ccb(hcb->ccb, CAM_CMD_TIMEOUT);
    }
}
```

当我中止一个命令，同一目标/LUN的所有其他断开连接的命令也会被中止。因此出了一个问题，我应当返回它的状态CAM\_REQ\_ABORTED还是CAM\_CMD\_TIMEOUT？当前的程序使用CAM\_CMD\_TIMEOUT。看起来符合逻辑，因为如果一个命令超时，很可能已经出了某些的很糟的事情，因此如果它没有被乱它自己应当超时。

# Chapter 12. USB

## 12.1. 介绍

通用串行总线(USB)是将设备接到个人计算机的一种新方法。它突出了双向通信的特色，并且其充分考虑到了正逐渐智能化和需要与host进行更多交互的方面。USB的支持包含在当前所有芯片中，因此在最近制造的PC中都可用。苹果(Apple)引入了USB的iMac硬件制造商生产的USB版本的是一个很大的激励。未来的PC将指定PC上的所有老接器将由一个或多个USB接器取代，提供通用的即插即用能力。USB硬件的支持在NetBSD的相当早期就有了，它是由Lennart Augustsson在NetBSD项目中的。代码已被移植到FreeBSD上，我目前拿着一个底共享代码。USB子系统的到来，许多USB的特性很重要。

Lennart Augustsson已经完成了NetBSD项目中USB支持的大部分。十分感谢他们的工作量。也十分感谢Ardy和Dirk对本文稿的贡献和校对。

- 直接接到计算机上的端口，或者接到称作集中器的设备，形成树型连接。
- 可在任何时间接或断开。
- 可以挂起自身并触发host系统的重新投入运行。
- 由于设备可由host供电，因此host文件必须跟踪每个集中器的电源状态。
- 不同类型的需要不同的服务质量，并且同一端口可以连接最多126个设备，就需要恰当地调度线路上的数据以充分利用12Mbps的可用带宽。(USB 2.0超过400Mbps)
- 智能化并包含很容易读到的关于自身的信息。

USB子系统以及接到它的驱动程序受已实现或将要实现的果(Apple)通过使得通用程序可从其操作系统MacOS中获得，而且不鼓励使用新的API。强烈推行基于标准的程序。本章整理基本信息以便在FreeBSD/NetBSD中USB的当前有个基本的了解。然而，建议将下面参考中提及的相关与本章同读。

### 12.1.1. USB的驱动

FreeBSD中的USB支持可被分为三部分。最底层包含主控器，向硬件及其调度提供一个通用接口。它支持硬件初始化，执行读写操作，管理已完成/失败的请求。一个主控器驱动程序控制一个虚拟hub，以硬件无方式提供对控制机器背面根端口的寄存器的访问。

中间层处理连接和断开，基本初始化，驱动程序的通信通道(管道)和电源管理。一个驱动程序也控制默管道和其上方的需求。

顶层包含支持特定的操作系统的各个驱动程序。有些驱动程序除了默管道外的其他管道上使用的驱动。它们也有额外功能，使得内核或用户空间是可见的。它们使用暴露出来的USB程序接口(USBPI)。

## 12.2. 主控器

主控器(HC)控制包上层的驱动。使用1毫秒的间隔，在帧开始时，主控器生成一个帧开始(SOF, Start of Frame)包。

SOF包用于同步的帧开始和跟踪包的数目。包在主机中被生成，或由host到设备(out)，或由设备到host(in)。

包是由host发起（host-initiated）。因此一条USB包只能有一个host。每个包的头部都有一个状态段，数据接收者可以在其中返回ACK（答接收），NAK（重置），STALL（停止条件）或什么也没有（混乱数据段，不可用或已断开）。USB [specification](#)的第8.5节更详细地解除了包的头部。USB上可以输出四中不同类型的包：控制(control)，大(bulk)，中断(interrupt)和同(isochronous)。它们的特性和他们的特性在下面描述（管道'子节中）。

USB上的包和程序的大型被主控器或HC程序分割为多个包。

到默端点的包（控制）有些特殊。它由一个或三个段组成：SETUP，数据(DATA, 可选) 和状态(STATUS)。设置(set-up)包被送到。如果存在数据段，数据包的方向在设置包中指出。状态段中的方向与数据段期间的方向相反，或者当没有数据段时IN。主控器硬件也提供寄存器，用于保存根端口的当前状态和自从状态改变寄存器最后一次位以来所发生的更改。USB[2]建议使用一个虚拟hub来提供一些寄存器的。虚拟hub必须符合第11章中列出的hub规范。它必须提供一个默管道使得包可以发送。它返回标准和hub特定的一组描述符。它也应当提供一个中断管道用来报告其端口发生的变化。当前可用的主控器有4个：[通用主控器接口](#)(UHCI；英特尔)和[开放主控器接口](#)(OHCI；康柏，微软，国家半导体)。UHCI的通信要求主控器程序的提供完整的深度，从而减少了硬件的复杂性。OHCI类型的控制器自身提供一个更抽象的接口来完成很多工作，从而更加独立。

### 12.2.1. UHCI

UHCI主控器总共有1024个指向数据的列表。它理解不同的数据类型：TD描述符(TD)和QH列表(QH)。一个TD表示与端点进行通信的一个包。QH是将一些TD(和QH)分成的一组方法。

一个由一个或多个包组成。UHCI程序将大的分割成多个包。除同外，每个都会分配一个QH。对于类型的，都有一个与此类型对应的QH，所有这些QH都会被集中到一个QH上。由于有固定的延迟需求，同必须首先执行，它是通信列表中的直接引用的。最后的同TD引用那一的中断的QH。中断的所有QH指向控制的QH，控制的QH又指向大类的QH。下面的表格出了一个形概：

导致下面的深度会在中执行。控制器从列表中取得当前中的所有的同(isochronous)包执行TD。一些TD的最后一个引用那一的中断的QH。然后主控器将从那个QH下行到各个中断的QH。完成那一列后，中断的QH会将控制器指向到所有控制的QH。它将在那儿等待深度的所有子列，然后是在大QH中排列的所有。为了方便管理已完成或失败的，硬件会在末尾产生不同的中断。在的最后一个TD中，HC程序设置Interrupt-On-Completion位来完成的一个中断。如果TD到了其最大数，就产生中断。如果在TD中设置短包位，且小于所设置的包深度(的包)，就会此中断以通知控制器程序已完成。输出一个已完成或产生是主控器程序的任务。当中断服务例程被调用，它将定位所有已完成的并用它的回。更尽的描述看 [UHCI specification](#)。

### 12.2.2. OHCI

OHCI主控器编程要容易得多。控制器假有一端点(endpoint)可用，并知道中不同类型的深度优先和排序。主控器使用的主要数据是端点描述符(ED)，它上面接着一个TD描述符的列。ED包含端点所允许的最大的包大小，控制器硬件完成包的分割。每次都会更新指向数据缓冲区的指针，当起始和止指针相等，TD就退到完成队列(done-queue)。四种类型的端点各有其自己的列。控制和大(bulk)端点分在它们自己的队列。中断ED在中排队，在中的深度决定了它执行的深度。

列表 中断 同(isochronous) 控制 大(bulk)

主控器在其中执行的顺序看起来如下。控制器首先执行非周期性控制和大环列，最后可到HCD程序设置的一个限制。然后以低5位作为中断ED上深度的那一行中的索引，执行那个号的中断。在那个的末尾，同ED被接，并随后被遍历。同TD包含了当行其中的第一个的号。所有周期性的命令行以后，控制和大环列再次被遍历。中断服务例程会被周期性地调用，来处理完成的行，每个命令用回，并重新调度中断和同端点。

更尽的描述看 [OHCI specification](#)。服务，即中，提供了以可控的方式进行行，并随着由不同程序和服务所使用的源。此物理下面几方面：

- 命令配置信息
- 与命令行通信的管道
- 探测和连接，以及从分离(detach)。

## 12.3. USB命令信息

### 12.3.1. 命令配置信息

一个命令提供了不同的配置信息。一个命令具有一个或多个配置，探测/连接期从其中一个定一个。配置提供功率和要求。一个配置中可以有多个接口。接口是端点的集(collection)。例如，USB声器可以有一个音频接口（音），和旋钮(knob)、号(dial)和按的接口（HID）。一个配置中的所有接口可以同时有效，并可被不同的程序连接。一个接口可以有通用接口，以提供不同量的服务参数。例如，在照相机中，用来提供不同的大小以及秒数。

一个接口中可以指定0或多个端点。端点是与命令行通信的方向点。它提供缓冲区来存放从而来的，或外出到的数据。一个端点在配置中有唯一地址，即端点号加上其方向。默认端点，即端点0，不是任何接口的一部分，并且在所有配置中可用。它由服务管理，并且命令程序不能直接使用。

Level 0 Level 1 Level 2 Slot 0

Slot 3 Slot 2 Slot 1

(只展示了32个槽中的4个)

多次化配置信息在命令中通过准的一组描述符来描述（参看 [USB\[2\]第9.6章](#)）。它可以通Get Descriptor Request来求。服务存一些描述符以避免在USB上进行不必要的。这些描述符的通函数是用来提供的。

- 描述符：对于的通用信息，如供应商，产品和修订ID，支持的、子和通用的，，默认端点的最大包大小等。
- 配置描述符：此配置中的接口数，支持的挂起和恢复能力，以及功率要求。
- 接口描述符：接口、子和通用的，接口通用配置的数目和端点数目。
- 端点描述符：端点地址、方向和类型，支持的最大包大小，如果是中断类型的端点包括频率。默认端点（端点0）没有描述符，而且从不被列入接口描述符中。
- 字符串描述符：在其他描述符中会某些字段提供字符串索引。它可被用来取描述性字符串，可能以多语言的形式提供。

说明(specification)可以添加它自己的描述符类型，有些描述符也可以通过GetDescriptor Request来得。

管道与端点的通信，流到所用的管道。程序将到端点的提交到管道，并提供（失败）失败或完成用的回，或等待完成（同）。到端点的在管道中被串行化。或者完成，或者失败，或者超（如果设置了超时）。由于有类型的超时。超时的可能由于USB上的超时（毫秒）。有些超时被失败，可能是由于中断接引起的。唯一超时在文件中，当没有在指定的秒内完成触。是由于的包否定答引起的。其原因是由于没有准备好接收数据，缓冲区欠或超，或。

如果管道上的大于端点描述符中指定的最大包大小，主控器（OHCI）或HCD程序（UHCI）将按最大包大小分割，最后一个包可能小于最大包的大小。

有时候来返回少于所求的数据并不是个。例如，到制解器的大in可能请求200字的数据，但制解器那只有5个字可用。程序可以设置短包（SPD）志。它允许主控器即使在的数据量少于所求的数据量的情况下也接受包。这个志只在in中有效，因将要被送到的数据量是事先知道的。如果程序中出不可恢复的，管道会被停。接受或送更多数据以前，程序需要定停的原因，并通常在默管道上发送清除端点挂起请求（clear endpoint halt device request）来清除端点停条件。

有四种不同类型的端点和的管道：-

- 控制管道/默管道：一个有一个控制管道，接到默端点（端点0）。此管道请求和的数据。默管道和其他管道上的的区别在于所使用的，这在USB[2]中描述。一些请求用于位和配置。所有必须支持USB[2]的第9章中提供的一基本命令。管道上支持的命令可以通扩展，以支持外的功能。
- 大(bulk)管道：是USB与原始媒体的等价物。
- 中断管道：host向发送数据请求，如果没有应答，将NAK（否定答）数据包。中断按建管道指定的率被度。
- 同管道：些管道用于具有固定延的同数据，例如或音流，但不一定。当前中已有类型的某些支持。当期出，或者由于，例如缺乏缓冲区空来存入的数据而引起的否定答包（NAK），控制、大和中断管道中的包会被重。而同包在失或包NAK不会重，因那可能反同束。

所需度的可用性在管道的建期被算。在1毫秒的内行度。中的分配由USB的第5.6决定。同和中断被消耗中多90%的。控制和大包在所有同和中断包之后行度，并将消耗所有剩余。

于度和回收的更多信息可以在USB[2]的第5章，UHCI[3]的第1.3，OHCI[4]的3.4.2中到。

## 12.4. 的探和接

集中器(hub)通知新已接后，服务端口加(switch on)，提供100mA的流。此位于其默状态，并听地址0。服务会通过默管道取各描述符。此后它将向发送Set Address请求，将从默地址(地址0)移。可能有多个程序支持此。例如，一个制解器可能通过AT兼容接口支持ISDN TA。然而，特定型号的ISDN配器的程序可能提供比更好的支持。为了支持的活性，探会返回先，指示他的支持。支持品的特定版本会具有最高。如果一个配置内有多个接口，也可能多个程序会接到一个。一个程序只需支持所有接口的一个子集。

新接的探程序，首先探特定的程序。

如果没有，探代替在所有支持的配置上重探

程序，直到在一个配置中接到一个程序。为了支持不同接口上使用多个程序的，探针会在一个配置中的所有尚未被程序声明(claim)的接口上重行。超出集中器功率算的配置会被忽略。接期，程序当把初始化到当状态，但不能位，因为那会使得将它自己从上断开，并重新探测程序。为了避免消耗不必要的，不当在声明中断管道，而当延分配管道，直到打开文件并真的使用数据。当文件，管道也当被再次，尽管可能仍然接着。

### 12.4.1. 断接(disconnect)和分离(detach)

程序与行任何事期，当期会接收到。USB的的支持并鼓励在任何点及断接。程序当保当不在做正的事情。

此外，断接(disconnect)后又重新接(reconnect)的不会被重新接(reattach)相同的例。将来当更多的支持序列号（参看描述符），或出其他定的方法的时候，情况可能会改。

断接是由集中器在到集中器程序的中断包中信号通知(signal)的。状态信息指示个端口接改。接到那个端口上的的所有程序共用的分方法被用，被底清理。如果端口状态指示同一个已接(connect)到那个端口，探和接的程将被。位将在集中器上生一个断接序列，并将按上面所述行理。

## 12.5. USB程序的信息

USB没有定除默管道外其他管道上使用的。方面的信息可以从各来源得。最准的来源是USB主[1]上的者部分。从些面上可以得到数目不断的程序。些指定从程序角度看起来兼容当，它需要提供的基本功能和通信通道上使用的。USB[2]包括了集中器的描述。人机界面(HID)的已建出来，以迎合、数字入板、条形器、按钮、旋(手柄knob)、等的要求。一个例子是用于大容量存储的。完整的完整列表参看USB主[1]的者部分。

然而，多的程序信息没有被公布。于所用的信息可能可以从制造的公司得。一些公司会在之前要求签署保密(Non-Disclosure Agreement, NDA)。大多数情况下，会阻止将程序放源代。

一个信息的很好来源是Linux程序源代，因很多公司已始他们的提供Linux下的程序。系那些程序作者他自己的信息来源是一个好主意。

例子：人机界面。人机界面，如、鼠标、数字入板、按钮、等的被其他程序引用，并在很多中使用。

例如，音声器提供到数模器的端点，可能提供外管道用于麦克。它前面的按钮和号在独的接口中提供HID端点。器控制也是如此。通可用的内核和用空的，与HID程序或通用程序一起可以直接地建些接口的支持。一个可以作在一个配置中的多个接口由不同的程序的例子，一个是一便宜的，有老的鼠标接口。为了避免在USB集中器包括一个硬件而致的成本上升，制造商将从背面的PS/2端口接收到的鼠标数据与来自的按钮合成在同一个配置中的一个独的接口。鼠标和程序各自接到适当的接口，并分配到个独立端点的管道。

例子：固件下。已出来的多是基于通用目的处理器，并将外的USB核心加入其中。由于程序的和USB的固件仍然非常新，多需要在接(connect)之后下固件。

下面的非常明直接。通供商和品ID自身。第一个程序探并接到它，并将固件下到其中。此后自己位，程序分。短的停之后宣布它在上的存在。将改其供商/品

/版本的ID以反映其提供有固件的事□，因此□一个□□程序 将探□它并□接(attach)到它。

□些□型的□□的一个例子是基于EZ-USB的ActiveWire I/O板。□个 芯片有一个通用固件下□器。下□到ActiveWire板子上的固件改□版本ID。 然后它将□行EZ-USB芯片的USB部分的□□位，从USB□□上断□，并再次 重新□接。

例子：大容量存□□。□大容量存□□的支持主要□□□有的 □□□建。Iomega USB Zip□□器是基于SCSI版本的□□器。SCSI命令和 状态信息被包装到□□中，在大□(bulk)管道上□□到/来自□□，在USB□□上模□SCSI控制器。ATAPI和UFI命令以相似的方式被支持。

大容量存□□□支持□□不同□型的□命令□的包装。最初的□□ 基于通□默□管道□送命令和状态信息，使用大□□□在host和□□之□ 移□数据。在□□基□上□□出□一□方法，□□方法基于包装命令和 状态，并在大□out和in端点上□送它□。□□精□地指定了何□必□ □生什□，以及在□□到□□条件下□□做什□。□□些□□□写 □□程序的最大挑□是□□基于USB的□□，□□它□合已有的□□大容量存□□□ 的支持。CAM提供了□□子，以相当直接了当的方式来完成□□。ATAPI就 没有□□□□了，因□□史上IDE接口从未有□多□不同的表□方式。

来自Y-E Data的□□USB□□的支持也不是那□直□，因□□□了一套 新的命令集。

# Chapter 13. Newbus

特·感·Matthew N. Dodd, Warner Losh, Bill Paul, Doug Rabson, Mike Smith, Peter Wemm and Scott Long.

本章·解·了Newbus·框架。

## 13.1. ····程序

### 13.1.1. ····程序的目的

···程序是·件·件，它在内核·于外···（例如，磁·、网···配·）的通用··和外···的····之·提供了接口。···程序接口(DDI)是内核与···程序·件之·定·的接口。

### 13.1.2. ····程序的··型

在UNIX®那个·代，FreeBSD也从中延·而来，定·了四··型的··：

- 程序
- 字符···程序
- 网···程序
- 程序

··以使用固定大小的[数据]·的方式·行。··型的··程序依·所·的··冲区·存(buffer cache)，其目的·是在内存中的·用区域·存··的数据·。··冲区·存常常基于后台写(write-behind)，·意味着数据在内存中被修改后，当系··行其周期性·磁·刷新·才会被同·到磁·，从而·化写操作。

### 13.1.3. 字符··

然而，在FreeBSD 4.0版本以及后·版本中，···和字符··的区··得不存在了。

## 13.2. Newbus概·

Newbus·了一·基于抽象·的新型···，可以在FreeBSD 3.0中看到···的介·，当·Alpha的移植被·入到代··中。直到4.0它才成···程序使用的默·系·。其目的是·主机··系·提供··操作系·的各···和··的互··提供更加·面向··象的方法。

其主要特性包括：

- 接
- 程序容易模··化
- 

最·著的改·之一是从平面和特殊系·演···布局。

··留的是"根"·，它作·父·，所有其他·挂接在它上面。·于·个·，通常"根"·只有·个孩子，其上·接着·如host-to-PCI·等·西。·于x86，··"根"··· "nexus"·，·于Alpha，Alpha的各·不同型号有不同的···

, 不同的硬件芯片，包括 *lca*, *apecs*, *cia*和*tsunami*。

Newbus上下文中的表示系中的一个硬件体。例如，一个PCI被 表示一个Newbus。系中的任何一个可以有孩子；有孩子的通常被 称"bus"。系中常用的例子就是 ISA和PCI，他们各自管理接到ISA和PCI上的列表。

通常，不同类型的之的接被表示 " "，它的孩子就是它所接的 。一个例子就是PCI-to-PCI，它在父PCI上被 表示`pcibN`，而用它的孩子 `-pciN`表示接在它上面的 。布局化了PCI的，允公共代码用于和接的 。

Newbus中的个求它的父来其映射源。父接着求 它的父，直到到nexus。因此，基本上nexus是Newbus系中唯一知道所有源的部分。



ISA可能想在`0x230`映射其IO端口，因此它向其父求，情况下是ISA。ISA将它交PCI-to-ISA，PCI-to-ISA 接着求PCI，PCI到host-to-PCI ，最后到nexus。向上 渡的美之在于可以有空来求。`0x230`IO端口 的求在MIPS机器上可以被PCI成`0xb000230`的内存映射。

源分配可以在的任何地方加以控制。例如，在很多Alpha平台上， ISA中断与PCI中断是独管理的， ISA中断的源分配是由Alpha的ISA管理。在IA-32上，ISA和PCI中断都由nexus管理。于移植， 内存和端口地址空间由个体管理 - 在IA-32上是nexus，在Alpha（例如， CIA或tsunami）上是相的芯片程序。

为了简化内存和端口映射源的， Newbus整合了NetBSD的 `bus_space` API。他提供了唯一的API来代替inb/outb 和直接内存写。做的在于一个程序就可以使用内存映射寄存器或端口映射寄存器（有些硬件支持者）。

支持被合并到了源分配机制中。分配源， 程序可以从源 中取的`bus_space_tag_t`和`bus_space_handle_t`。

Newbus也允在用于此目的的文件中定接口方法。些是 .m文件，可以在src/sys 目中到。

Newbus系的核心是可展的"基于象程(object-based programming)"的模型。系中的个具有它所支持的一个方法表。系和其他使用些方法来控制并求服。所支持的不同方法 被定多个"接口"。"接口"只是 的一个相的方法。

在Newbus系中， 方法是通系中的各程序提供的。当 自配置(auto-configuration)期被接(attach) 到程序，它使用程序声明的方法表。以后可以从其程序 分(detach)，并 重新接(re-attach)到具有新方法表的新程序。就 允程序的替，而 替于程序的非常有用。

接口通与文件系中用于定vnode操作的言相似的接口定言来描述。接口被保存在方法文件中（通常命名`foo_if.m`）。

## 例 6. Newbus的方法

```
# Foo 子系的程序 (注...) 

INTERFACE foo

METHOD int doit {
    device_t dev;
};

# 如果没有通过DEVMETHOD()提供一个方法， DEFAULT将会被使用的方法

METHOD void doit_to_child {
    device_t dev;
    driver_t child;
} DEFAULT doit_generic_to_child;
```

当接口被声明后，它生成一个头文件 "foo\_if.h"，其中包含函数声明：

```
int FOO_DOIT(device_t dev);
int FOO_DOIT_TO_CHILD(device_t dev, device_t child);
```

伴随自动生成的头文件，也会生成一个源文件 "foo\_if.c"；其中包含一些函数的实现。这些函数用于在对象方法表中相应函数的位置并调用那个函数。

系统定义了两个主要接口。第一个基本接口被称为 "device"，并包括与所有设备相关的所有方法。 "device" 接口中的方法包括 "探测(probe)"、"连接(attach)" 和 "分离(detach)"，它们用来控制硬件的连接，以及 "shutdown"、"挂起(suspend)" 和 "恢复(resume)"，它们用于事件通知。

另一个，更加通用的接口是 "bus"。此接口包含的方法用于处理有孩子的设备，包括向特定的设备发送信息，事件通知 (child\_detached, driver\_added) 和资源管理 (alloc\_resource, activate\_resource, deactivate\_resource, release\_resource)。

"bus" 接口中的很多方法只对某些孩子设备有效。有些方法通常使用前两个参数指定提供服务的设备和请求服务的子设备。除了简化代码的程序代码，有些方法中的很多都有访问者(accessor)函数，后者函数用来调用父设备并调用父设备上的方法。例如，方法 BUS\_TEARDOWN\_INTR(device\_t dev, device\_t child, ...) 可以使用函数 bus\_teardown\_intr(device\_t child, ...) 来调用。

系统中的某些驱动程序提供了对外接口以提供特定功能的服务。例如，PCI驱动程序定义了 "pci" 接口，此接口有两个方法 read\_config 和 write\_config，用于访问 PCI 设备的配置寄存器。

## 13.3. Newbus API

由于Newbus API非常庞大，本章努力将它文档化。本文档的下一版本会带来更多信息。

### 13.3.1. 源代码目录中的重要位置

src/sys/[arch]/[arch] - 特定机器的 内核代码位于这个目录。例如 i386 或 SPARC64。

src/sys/dev/[bus] - 支持特定 [bus] 的代码位于这个目录。

src/sys/dev/pci - PCI 代码支持代码位于这个目录。

src/sys/[isa | pci] - PCI/ISA 程序 位于这个目录。FreeBSD 4.0 版本中，PCI/ISA 支持代码 去存在于这个目录中。

### 13.3.2. 重要的结构和类型定义

`devclass_t` - 这是指向 `struct devclass` 的指针的类型定义。

`device_method_t` - 与 `kobj_method_t` 相同（参看 `src/sys/kobj.h`）。

`device_t` - 这是指向 `struct device` 的指针的类型定义。`device_t` 表示系统中的设备。它是内核对象。参看 `src/sys/sys/bus_private.h`。

`driver_t` - 这是一个类型定义，它引用 `struct driver`。`driver` 是一个 `device` 的内核对象；它也保存着程序的私有数据。

#### driver\_t 定义

```
struct driver {
    KOBJ_CLASS_FIELDS;
    void     *priv;           /* 程序私有数据 */
};
```

`device_state_t` 是一个枚举型，即 `device_state`。它包含 Newbus 在自配置前后 可能的状态。

#### 状态 `device_state_t`

```
/*
 * src/sys/sys/bus.h
 */
typedef enum device_state {
    DS_NOTPRESENT, /* 未探测或探测失败 */
    DS_ALIVE,       /* 探测成功 */
    DS_ATTACHED,    /* 使用了连接方法 */
    DS_BUSY         /* 已打包 */
} device_state_t;
```

# Chapter 14. 声音子系

## 14.1. 介绍

FreeBSD声音子系清晰地将通用声音处理与特定的分隔开来。使得更容易加入新旧的支持。

pcm(4)框架是声音子系的中心部分。它主要下面的文件：

- 一个到数字化声音和混音器函数的公用接口（read, write, ioctls）。ioctl命令集合兼容老的OSS或Voxware接口，允许多媒体应用程序不加修改地移植。
- 处理声音数据的公共代码（格式，虚通道）。
- 一个唯一的文件接口，与硬件特定的声音接口模块接口
- 某些通用硬件接口（ac97）或共享的硬件特定代码（例如：ISA DMA例程）的额外支持。

特定声音的支持是通过硬件特定的程序来完成的，有些程序提供通道和混音器接口，进入到通用pcm代码中。

本章中，pcm将指声音程序的中心，通用部分，是比硬件特定的模块而言的。

初期的程序作者当然希望从有模块开始，并使用那些代码作为参考。但是，由于声音代码十分漂亮，也基本上免除了注释。本文简述出框架接口的一个概貌，并回答改写有代码可能遇到的一些问题。

作之外的途径，或者除了从一个可工作的例子开始的方法之外，

可以从<http://people.FreeBSD.org/~cg/template.c>到一个注释的程序模板。

## 14.2. 文件

除/usr/src/sys/sys/soundcard.h中的公共 ioctl接口定义外，所有的相代码当前(FreeBSD 4.4)位于 /usr/src/sys/dev/sound/。

在/usr/src/sys/dev/sound/下面，pcm/目录中保存着中心代码，而isa/和pci/目录中有ISA和PCI板的程序。

## 14.3. 探测，连接等

声音程序使用与任何硬件程序模块相同的方法探测和连接（即）。可能希望看一下手册中ISA或PCI章节的内容来获取更多信息。

然而，声音程序在某些方面又有些不同：

- 它自己声明pcm，有一个私有`struct snddev_info`：

```

static driver_t xxx_driver = {
    "pcm",
    xxx_methods,
    sizeof(struct snddev_info)
};

DRIVER_MODULE(snd_xxxpci, pci, xxx_driver, pcm_devclass, 0, 0);
MODULE_DEPEND(snd_xxxpci, snd_pcm, PCM_MINVER, PCM_PREFVER,PCM_MAXVER);

```

大多数声音程序需要存于其的附加私有信息。私有数据通常在接例程中分配。其地址通常用 `pcm_register()` 和 `mixer_init()` 用 `pcm`。后面 `pcm` 将此地址作用声音程序接口的参数回来。

- 声音程序的接例程当通常用 `mixer_init()` 向 `pcm` 声明它的 MIXER 或 AC97 接口。于 MIXER 接口，会接着引起用 `xxxmixer_init()`。
- 声音程序的接例程通常用 `pcm_register(dev, sc, nplay, nrec)` 向 `pcm` 声明其通用 CHANNEL 配置，其中 `sc` 是数据的地址，在 `pcm` 以后的用中将会用到它，`nplay` 和 `nrec` 是播放和音通道的数目。
- 声音程序的接例程通常用 `pcm_addchan()` 声明它的个通道象。会在 `pcm` 中建立起通道合成，并接着会引起用 `xxxchannel_init()`（注：参考原文）。
- 声音程序的分例程在放其源之前当用 `pcm_unregister()`。

有可能的方法来理非PnP：

- 使用 `device_identify()` 方法（例：sound/isa/es1888.c）。`device_identify()` 方法在已知地址探硬件，如果支持的就会建一个新的 `pcm`，一个 `pcm` 接着会被到 `probe/attach`。
- 使用定制内核配置的方法，`pcm` 置适当的 hints（例：sound/isa/mss.c）。

`pcm` 程序当 `device_suspend`，`device_resume` 和 `device_shutdown` 例程，源管理和模卸就能正确地作用。

## 14.4. 接口

`pcm` 核心与声音程序之的接口以 内核象的叫法来定。

声音程序通常提供主要的接口：*CHANNEL* 以及 *MIXER* 或 *AC97*。

*AC97* 是一个很小的硬件（寄存器写）接口，由程序 *AC97* 解器的硬件来。情况下，的 *MIXER* 接口由 `pcm` 中共享的 *AC97* 提供。

### 14.4.1. CHANNEL 接口

#### 14.4.1.1. 函数参数的通常注意事

声音程序通常用一个私有数据来描述他的，程序所支持的播放和音数据通道各有一个。

于所有的 CHANNEL 接口函数，第一个参数是一个不透明的指。

第二个参数是指向私有的通道数据的指针，`channel_init()`是个例外，它的指针指向私有  
(并返回由pcm以后使用的通道指针)。

#### 14.4.1.2. 数据操作概要

对于声音数据，pcm核心与声音程序是通过一个由`struct snd_dbuf`描述的共享内存区域进行通信的。

`struct snd_dbuf`是pcm私有的，声音程序通常用前者函数(`sndbuf_getxxx()`)来获得感兴趣的。

共享内存区域的大小等于`sndbuf_getsize()`，并被分割为大小固定，且等于`sndbuf_getblk()`字节的很多块。

当播放时，常量的机制如下(将意思反过来就是录音)：

- pcm首先填充缓冲区，然后以参数PCMTRIG\_START用声音程序的`xxxchannel_trigger()`。
- 声音程序接着安排以`sndbuf_getblk()`字节大小，重新将整个内存区域(`sndbuf_getbuf()`,  
`sndbuf_getsize()`)映射到。对于每个映射调用pcm函数`chn_intr()`(通常在中断发生)。
- `chn_intr()`安排将新数据拷贝到那些数据已映射到(在空)的区域，并由`snd_dbuf`进行适当的更新。

#### 14.4.1.3. channel\_init

用`xxxchannel_init()`来初始化一个播放和录音通道。它用从声音程序的直接例程中起。(参看[探针和  
连接一](#))。

```
static void *
xxxchannel_init(kobj_t obj, void *data,
                 struct snd_dbuf *b, struct pcm_channel *c, int dir).
{
    struct xxx_info *sc = data;
    struct xxx_chinfo *ch;
    ...
    return ch;
}
```

b通道`struct snd_dbuf`的地址。它应当在函数中通过用`sndbuf_alloc()`来初始化。所用的缓冲区大小通常是典型的大小的一个较小的倍数。`c`是pcm通道控制块的指针。是个不透明指针。函数应当将它保存到局部通道块中，在后面用pcm函数(例如：`chn_intr(c)`)会使用它。`dir`指示通道方向(`PCMDIR_PLAY`或`PCMDIR_REC`)。函数应当返回一个指针，此指针指向用于控制此通道的私有区域。它将作参数被映射到其他通道接口的用。

#### 14.4.1.4. channel\_setformat

`xxxchannel_setformat()`应当按特定通道，特定声音格式设置硬件。

```

static int
xxxchannel_setformat(kobj_t obj, void *data, u_int32_t format).
{
    struct xxx_chinfo *ch = data;
    ...
    return 0;
}

```

`format` AFMT\_XXX `value`之一 (soundcard.h)。

#### 14.4.1.5. channel\_setspeed

`xxxchannel_setspeed()`按指定的取速度设置通道硬件，并返回可能的速度。

```

static int
xxxchannel_setspeed(kobj_t obj, void *data, u_int32_t speed)
{
    struct xxx_chinfo *ch = data;
    ...
    return speed;
}

```

#### 14.4.1.6. channel\_setblocksize

`xxxchannel_setblocksize()`设置大小，是pcm与声音程序，以及声音程序与之的位的大小。周期，次大小的数据后，声音程序都用pcm的 `chn_intr()`。

大多数程序只注意块的大小，因当开始使用个。

```

static int
xxxchannel_setblocksize(kobj_t obj, void *data, u_int32_t blocksize)
{
    struct xxx_chinfo *ch = data;
    ...
    return blocksize;
}

```

函数返回可能的大小。如果大小真的化了，情况下当用 `sndbuf_resize()` 整缓冲区的大小。

#### 14.4.1.7. channel\_trigger

`xxxchannel_trigger()`由 pcm 来控制程序中的操作。

```

static int
xxxchannel_trigger(kobj_t obj, void *data, int go).
{
    struct xxx_chinfo *ch = data;
    ...
    return 0;
}

```

`go` 定义当前用的工作。可能的值：



如果程序使用ISA DMA，应当在DMA上执行工作前用 `sndbuf_isadma()`，并清理 DMA 芯片一方的事情。

#### 14.4.1.8. channel\_getptr

`xxxchannel_getptr()` 返回缓冲区中当前的缓冲。它通常由 `chn_intr()` 用，而且也是为什么 `pcm` 知道它应当往哪儿送新数据。

#### 14.4.1.9. channel\_free

用 `xxxchannel_free()` 来释放通道资源，例如当程序卸载，并且如果通道数据是分配的，或者如果不使用 `sndbuf_alloc()` 行缓冲区分配，应当调用这个函数。

#### 14.4.1.10. channel\_getcaps

```

struct pcmchan_caps *
xxxchannel_getcaps(kobj_t obj, void *data)
{
    return xxx_caps;
}

```

这个例程返回指向（通常静态的）`pcmchan_caps` 的指针（在 `sound/pcm/channel.h` 中定义）。这个结构保存着最小和最大采样率和被接受的声音格式。任何声音程序都可以作一个例子。

#### 14.4.1.11. 更多函数

`channel_reset()`, `channel_resetdone()` 和 `channel_notify()` 用于特殊目的，未与威士士 (Cameron Grant) 行探之前不应当在程序中用它。

不成使用 `channel_setdir()`。

### 14.4.2. MIXER 接口

#### 14.4.2.1. mixer\_init

`xxxmixer_init()` 初始化硬件，并告诉 `pcm` 什么混音器可用来播放和录音。

```

static int
xxxmixer_init(struct snd_mixer *m)
{
    struct xxx_info    *sc = mix_getdevinfo(m);
    u_int32_t v;

    [初始化硬件]

    [置播放混音器中v中当的位].
    mix_setdevs(m, v);
    [置音混音器中v中当的位]
    mix_setrecdevs(m, v)

    return 0;
}

```

置一个整数中的位，并用 `mix_setdevs()` 和 `mix_setrecdevs()` 来告 pcm 存在什么。

混音器的位定可以在 `soundcard.h` 中找到。（`SOUND_MASK_XXX` 和 `SOUND_MIXER_XXX` 移位）。

#### 14.4.2.2. mixer\_set

`xxxmixer_set()` 混音器置音量 (level)。

```

static int
xxxmixer_set(struct snd_mixer *m, unsigned dev,
             unsigned left, unsigned right).
{
    struct sc_info *sc = mix_getdevinfo(m);
    [置音量(left)]
    return left | (right < 8);
}

```

被指定 `SOUND_MIXER_XXX` 在 [0-100] 之指定音量。零当静音。由于硬件(音量)可能不匹配入比例，会出某些整，例程返回如上面所示的值 (0-100 内)。

#### 14.4.2.3. mixer\_setrecsrc

`xxxmixer_setrecsrc()` 定音源。

```

static int
xxxmixer_setrecsrc(struct snd_mixer *m, u_int32_t src).
{
    struct xxx_info *sc = mix_getdevinfo(m);

    [看src中的非零位, 置硬件]

    [更新src反映操作]
    return src;
}

```

期望的音由一个位域指定。返回置用来音的。一些程序只能置一个音。如果出错，函数当返回-1。

#### 14.4.2.4. mixer\_uninit, mixer\_reinit

`xxxmixer_uninit()`当保不会出任何声音，并且如果可能当混音器硬件断开。

`xxxmixer_reinit()`当保混音器硬件加， 并且恢所有不受`mixer_set()`或`mixer_setrecsrc()`控制的置。

### 14.4.3. AC97接口

AC97由有AC97解器的程序。它只有三个方法：

- `xxxac97_init()`返回到的 ac97解器的数目。
- `ac97_read()`与 `ac97_write()`写指定的寄存器。

The AC97接口由 pcm中的AC97代码来执行高操作。参看 sound/pci/maestro3.c或sound/pci/下很多其他内容作例。

# Chapter 15. PC Card

本章将介绍FreeBSD如何处理PC Card或CardBus的驱动程序而提供的机制。但目前本文只提到了如何向已有的pccard程序中添加驱动程序。

## 15.1. 添加驱动

向所支持的pccard列表中添加新驱动的步骤与系统在FreeBSD 4中使用的方法不同了。在以前的版本中，需要在 /etc 中的一个文件来列出驱动。从 FreeBSD 5.0 开始，驱动程序知道它支持什么。在内核中有一个受支持的表，驱动程序用它来接线。

### 15.1.1. 概述

可以有多种方法来处理PC Card，它们都基于上的 CIS 信息。第一种方法是使用制造商和产品的数字 ID。第二种方法是使用人可读的字符串，字符串也是包含在 CIS 中。PC Card 使用集中式数据和一些宏来提供一个易用的模式，驱动程序的编写者很容易地指定匹配其驱动程序的 ID。

一个很普遍的情况是，某个公司的一款 PC Card 产品会参考原厂，然后把这个产品外的公司，以便在市场上出售。那些公司改掉原厂，把向他们的目标客户群或地理区域出售产品，并将他们自己的名字放到其中。然而所有的都有自己的改名，即使做任何修改，这些修改通常也微乎其微。然而，除了美化他们自己版本的品牌，有些供应商常常会把他自己的名字放入 CIS 空白的可读字符串中，却不会改掉制造商和产品的 ID。

对于以上情况，由于 FreeBSD 来使用数字 ID 可以减少工作量。同时也可能会将 ID 加入到系统程序中所来的属性。必须注意的是，真正制造者，尤其当提供原本的供应商在中心数据中已有不同的 ID。Linksys, D-Link 和 Netgear 是经常出售相同产品的几个美国制造商。相同的 ID 可能在日本以及 Buffalo 和 Corega 的名字出售。然而，有些 ID 常常具有相同的制造商和产品 ID。

PC Card 在其中心数据 /sys/dev/pccard/pccarddevs 中保存了 ID 的信息，但不包含那个驱动程序与它相关的信息。它也提供了一套宏，以允许在驱动程序用来声明驱动的表中容易地构建条目。

最后，某些非常低端的驱动根本不包含制造商 ID。有些驱动需要使用可读 CIS 字符串来匹配它们。如果我不需要急功近利有多好，但对于某些非常低端却非常流行的 CD-ROM 播放器来说却是必需的。通常应当避免使用数字方法，但本节中是列出了很多，因为它们是在驱动到 PC Card 商的 OEM 本之前加入的，应当先使用数字方法。

### 15.1.2. pccarddevs 的格式

pccarddevs 文件有四部分。第一部分使用它的一些供应商列出了制造商 ID。本节按数字排序。下一节包含了所有供应商使用的所有产品，包括它们的产品 ID 号和描述字符串。描述字符串通常不会被使用（相反，即使我可以匹配数字版本号，我仍然基于人可读的 CIS 置位的描述）。然后使用字符串匹配方法的那些部分重叠的东西。最后，文件任何地方可以使用 C 格式的注释。

文件的第一部分包含供应商 ID。保持列表按数字排序。此外，除了能有一个通用清晰的保存地来方便地保存一些信息，我与 NetBSD 共享此文件，因此不要对此文件的任何更改。例如：

vendor FUJITSU	0x0004	Fujitsu Corporation
vendor NETGEAR_2	0x000b	Netgear
vendor PANASONIC	0x0032	Matsushita Electric Industrial Co.
vendor SANDISK	0x0045	Sandisk Corporation

展示了几个供应商ID。很巧的是NETGEAR\_2 行上是NETGEAR从其的OEM，那些提供支持的作者那并不知道 NETgear使用的是人的ID。有些条目相当直接易行。行上都有供应商 字来指示本行的。也有供应商的名字。名字将会在pccarddevs文件 的后面重出，名字也会用在程序的匹配表中，因此保持它的短小 并且是有效的C符。有一个供应商的十六进制数字ID。不要添加0xffffffff或0xffff形式的ID， 因它保留ID（前者是'空ID集合'，而后者有会在量其差的中看到，用来指示none）。最后有于公司的描述字符串。个字符串 在FreeBSD中除了用于注目的外并没有被使用。

文件的第二包含品. 如在下面例子中看到的:

```
/* Allied Telesis K.K. */
product ALLIEDTELESIS LA_PCM    0x0002 Allied Telesis LA-PCM

/* Archos */
product ARCHOS ARC_ATAPI    0x0043 MiniCD
```

格式与供应商的那些行相似。其中有品字。然后是供应商名字， 由上面重而来。后面跟着品名字，此名字在程序中使用，且当 是一个有效C符，但可以以数字。然后是的十六进制品ID。供应商通常0xffffffff和 0xffff有相同的定。最后是于自身的字符串 描述。由于FreeBSD的pccard程序会从人可的CIS条目建一个 字符串，因此个字符串在FreeBSD中通常不被使用，但某些CIS条目不能 足要求的情况下可能使用。品按制造商的字母序排序，然后再按 品ID的数字排序。个制造商条目前有一条C注，条目之有一个空行。

第三很象前面的供应商一，但所由的制造商ID -1。-1在FreeBSD pccard 代中意味着"匹配任何东西"。由于它是C符，它的名字必唯一。除此之外格式等同于文件的第一。

最后一包含那些必用字符串匹配的。一的格式与通用 的格式有点不同：

```
product ADDTRON AWP100      { "Addtron", "AWP-100spWirelessspPCMCIA",
"Versionsp1.02", NULL }
product ALLIEDTELESIS WR211PCM { "AlliedspTelesisspK.K.", "WR211PCM", NULL, NULL }
Allied Telesis WR211PCM
```

我已熟悉了品字，后跟供应商名字，然后再跟的名字， 就象在文件第二中那。然而，之后就与那格式不同了。有一个 {}分，后跟几个字符串。些字符串CIS\_INFO三元中定的 供应商，品和外信息。些字符串被生成 pccarddevs.h的程序，将 sp替 空格。空条目意味着条目的部分当被忽略。在我 的例子中 有一个的条目。除非的操作来至重要，否不正当在其中包含版本号。有供应商在个字段中会有许多不同版本，有些版本 都能工作，情况下那些信息只会那些有相似的人在FreeBSD中 更以使用。有当供应商出于市考（可用性，价格等等），希望出售同一品牌的很多不同部分，也是必要的。如果，在那些 供应商仍然保持相同的制造商/品的少情况下，能否区分至 重要. 此不能使用正表式匹配。

### 15.1.3. 探针例程示例

要获得如何向所支持的列表中添加，就必须得很多程序。都有的探针和/或匹配例程。由于也有一些提供了一个兼容，在FreeBSD 5.x中有一点。由于只是window-dressing不同，从而出了一个理想化的版本。

```
static const struct pccard_product wi_pccard_products[] = {
    PCMCIA_CARD(3COM, 3CRWE737A, 0),
    PCMCIA_CARD(BUFFALO, WLI_PCM_S11, 0),
    PCMCIA_CARD(BUFFALO, WLI_CF_S11G, 0),
    PCMCIA_CARD(TDK, LAK_CD011WL, 0),
    { NULL }
};

static int
wi_pccard_probe(dev)
    device_t    dev;
{
    const struct pccard_product *pp;

    if ((pp = pccard_product_lookup(dev, wi_pccard_products,
        sizeof(wi_pccard_products[0]), NULL)) != NULL) {
        if (pp->pp_name != NULL)
            device_set_desc(dev, pp->pp_name);
        return (0);
    }
    return (ENXIO);
}
```

这儿我有一个可以匹配少数几个的pccard探针例程。如上面所提到，名字可能不同（如果不是`foo_pccard_probe()`就是`foo_pccard_match()`）。函数`pccard_product_lookup()`是一个通用函数，它遍历表并返回指向它所匹配的第一的指针。一些程序可能使用一个机制来将某些的附加信息映射到程序的其它部分，因此表中可能有些体。唯一的要求就是如果有不同的表，表的第一个元素是`pccard_product`。

察一下表`wi_pccard_products`就会，所有条目都是`PCMCIA_CARD(foo, bar, baz)`的形式。`foo`部分来自`pccarddevs`的制造商ID。`bar`部分产品。`baz`此所期望的功能号。许多`pccards`可以有多个功能，需要有办法区分功能1和功能0。可以看一下`PCMCIA_CARD_D`，它包括了来自`pccarddevs`文件的描述。也可以看看`PCMCIA_CARD2`和`PCMCIA_CARD2_D`，当需要按“使用默认描述”和“从`pccarddevs`中取得”做法，同匹配CIS字符串和制造商号就会用到它。

### 15.1.4. 将它合在一起

因此，加了一个新加新，必须行下面。首先，必须从得信息。完成个最容易的方法就是将入到PC Card或CF槽中，并出`devinfo -v`。可能会看到一些似下面的东西：

```
cbb1 pnpinfo vendor=0x104c device=0xac51 subvendor=0x1265 subdevice=0x0300  
class=0x060700 at slot=10 function=1  
    cardbus1  
    pccard1  
        unknown pnpinfo manufacturer=0x026f product=0x030c cisvendor="BUFFALO"  
        cisproduct="WLI2-CF-S11" function_type=6 at function=0
```

作出的一部分。制造商和产品品的数字ID。而cisvender和  
产品品的字符串。

cisproductCIS中提供的描述本

由于我首先想先使用数字，因此首先创建基于此的条目。  
看到的供应商BUFFALO，它已有一个条目了：

```
vendor BUFFALO      0x026f  BUFFALO (Melco Corporation)
```

就可以了。一个一个条目，但我没有。但我：

```
/* BUFFALO */  
product BUFFALO WLI_PCM_S11 0x0305  BUFFALO AirStation 11Mbps WLAN  
product BUFFALO LPC_CF_CLT 0x0307  BUFFALO LPC-CF-CLT  
product BUFFALO LPC3_CLT   0x030a  BUFFALO LPC3-CLT Ethernet Adapter  
product BUFFALO WLI_CF_S11G 0x030b  BUFFALO AirStation 11Mbps CF WLAN
```

就可以向pccarddevs中添加：

```
product BUFFALO WLI2_CF_S11G 0x030c  BUFFALO AirStation ultra 802.11b CF
```

目前，需要一个手来重新生成pccarddevs.h，用来将这些字符串到客程序。在程序中使用它之前必须完成下面：

```
# cd src/sys/dev/pccard  
# make -f Makefile.pccarddevs
```

一旦完成了这些，就可以向程序中添加了。只是一个添加一行的操作：

```
static const struct pccard_product wi_pccard_products[] = {  
    PCMCIA_CARD(3COM, 3CRWE737A, 0),  
    PCMCIA_CARD(BUFFALO, WLI_PCM_S11, 0),  
    PCMCIA_CARD(BUFFALO, WLI_CF_S11G, 0),  
+    PCMCIA_CARD(BUFFALO, WLI_CF2_S11G, 0),  
    PCMCIA_CARD(TDK, LAK_CD011WL, 0),  
    { NULL }  
};
```

注意，我在我添加的行前面包含了`+`，但`只是`用来`一行`。不要把它添加到`程序中`。一旦`添加了`  
`行`，就`可以重新``内核或模`，并`着看它是否能`。如果它`出``并能工作`，`提交`  
`。如果它不工作`，`出``它工作所需要的`并`提交一个`。如果它根本`，那`可能做`了什`，  
`当重新``一`。

如果`是一个FreeBSD源代`的committer，并且所有`看起来都`正常工作，`当把`些改`提交到`  
`中`。然而有些小技巧的`需要考`。首先，`必`提交pccarddevs文件到`中`。完成后，`必`重新`生`pccarddevs.h  
`并将它作`一次提交来提交（`是`了`保正`的`$FreeBSD$`会留在后面的文件中）。最后，`需要把`其它`提交到`程序。

### 15.1.5. 提交新`...`

很多人直接把新`的条目`送`作者`。`不要那`做。`将它`作`PR`来提交，并将PR号`送`作者用于`。`  
`保条目`不会`失`。提交`PR`，`丁`中没有必要包含pccarddevs.h的diff，`因`那些`可以重新`生。包含`的描述和客`程序的`丁`是必要`的`。如果`不知道名字`，使用OEM99作`名字`，作者将会`后相`地`整`OEM99。提交者不`当`提交OEM99，而`到`最高的OEM条目并提交高于那个`的一个`。

## 部分 III: 附录

## 参考目

[1] Marshall Kirk McKusick、Keith Bostic、Michael J Karels和John S Quarterman. 版 © 1996 Addison-Wesley Publishing Company, Inc.. 0-201-54979-4. Addison-Wesley Publishing Company, Inc.. *The Design and Implementation of the 4.4 BSD Operating System.* 1-2.