

# Linux中断处理

[www.gec-edu.org](http://www.gec-edu.org)

**中断源**——引起中断的原因，或者说发出中断请求的来源叫做中断源。一般是指由外设发出的中断请求，如：键盘中断、打印机中断等。外部中断是可以屏蔽的中断，也就是说，利用中断控制器可以屏蔽这些外部设备的中断请求。

**异常**——指因硬件出错（如复位、存储器访问等）或软件出错（如SWI、数据处理、未定义指令等）所引起的中断，也叫异常。异常中断是不可屏蔽的中断。

**中断嵌套** -- CPU在处理级别较低的中断过程中，出现了级别较高的中断请求，CPU停止执行低级别中断的处理程序而去优先处理高级别中断，等高级别中断处理完毕后，再接着执行低级别的未处理完的中断处理程序，这种中断处理方式称为中断嵌套。使用中断嵌套可以使高优先级别的中断得到及时的响应和处理。

## n 中断的作用

- 提高系统的效率,由于**CPU**速度远比外设高,轮询的效率非常低下。
- 中断使得单**CPU**的系统中可以执行多任务,它是实现现代操作系统的基础。
- 中断也是产生并发和同步问题的根源。

# ARM的中断向量表



Exception type	Mode	Normal address	High vector address
Reset	Supervisor	0x00000000	0xFFFF0000
Undefined instructions	Undefined	0x00000004	0xFFFF0004
Software interrupt (SWI)	Supervisor	0x00000008	0xFFFF0008
Prefetch Abort (instruction fetch memory abort)	Abort	0x0000000C	0xFFFF000C
Data Abort (data access memory abort)	Abort	0x00000010	0xFFFF0010
IRQ (interrupt)	IRQ	0x00000018	0xFFFF0018
FIQ (fast interrupt)	FIQ	0x0000001C	0xFFFF001C

- n 触发IRQ, CPU jump 到0x18, 同时要把irqno传入相应的寄存器, 调用一个中断通用处理函数:  
asm\_do\_IRQ(unsigned int irqno)。asm\_do\_IRQ() 这个函数根据irqno 就可以找到对应的中断描述符, 然后调用中断描述符里面的handler()了。
- n 中断向量表的初始化
- n arch/arm/kernel/entry-armv.S ——中断向量表放在这个文件里:
- n       .globl       \_\_vectors\_start
- n   \_\_vectors\_start:
- n   ARM( swi        SYS\_ERROR0)
- n   THUMB(        svc #0)
- n   THUMB(        nop )
- n   W(b)   vector\_und + stubs\_offset
- n   W(ldr) pc, .LCvswi + stubs\_offset
- n   W(b)   vector\_pabt + stubs\_offset
- n   W(b)   vector\_dabt + stubs\_offset
- n   W(b)   vector\_addrxcptn + stubs\_offset
- n   W(b)   vector\_irq + stubs\_offset
- n   W(b)   vector\_fiq + stubs\_offset
- n   .globl   \_\_vectors\_end
- n   \_\_vectors\_end:

n ARM linux内核启动时，通过start\_kernel()->trap\_init()的调用，初始化内核的中断异常向量表。

n CONFIG\_VECTORS\_BASE是一个宏，用来获取ARM异常向量的地址

n linux/arch/arm/kernel/traps.c

//复制异常中断向量表和各个异常中断对应的处理代码

n void \_\_init early\_trap\_init(void)

n {

n     unsigned long vectors = CONFIG\_VECTORS\_BASE;

n     extern char \_\_stubs\_start[], \_\_stubs\_end[];

n     extern char \_\_vectors\_start[], \_\_vectors\_end[];

n     extern char \_\_kuser\_helper\_start[], \_\_kuser\_helper\_end[];

n     int kuser\_sz = \_\_kuser\_helper\_end - \_\_kuser\_helper\_start;

n     memcpy((void \*)vectors, \_\_vectors\_start, \_\_vectors\_end - \_\_vectors\_start);

n     memcpy((void \*)vectors + 0x200, \_\_stubs\_start, \_\_stubs\_end - \_\_stubs\_start);

n     memcpy((void \*)vectors + 0x1000 - kuser\_sz, \_\_kuser\_helper\_start, kuser\_sz);

n     memcpy((void \*)KERN\_SIGRETURN\_CODE, sigreturn\_codes,

n         sizeof(sigreturn\_codes));

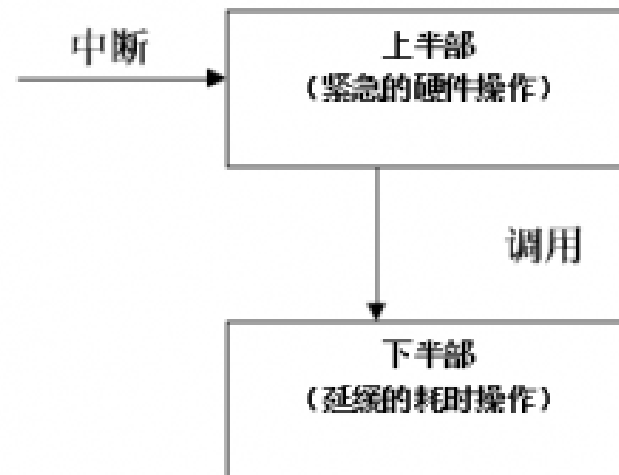
n     flush\_icache\_range(vectors, vectors + PAGE\_SIZE);

n     modify\_domain(DOMAIN\_USER, DOMAIN\_CLIENT);

n }

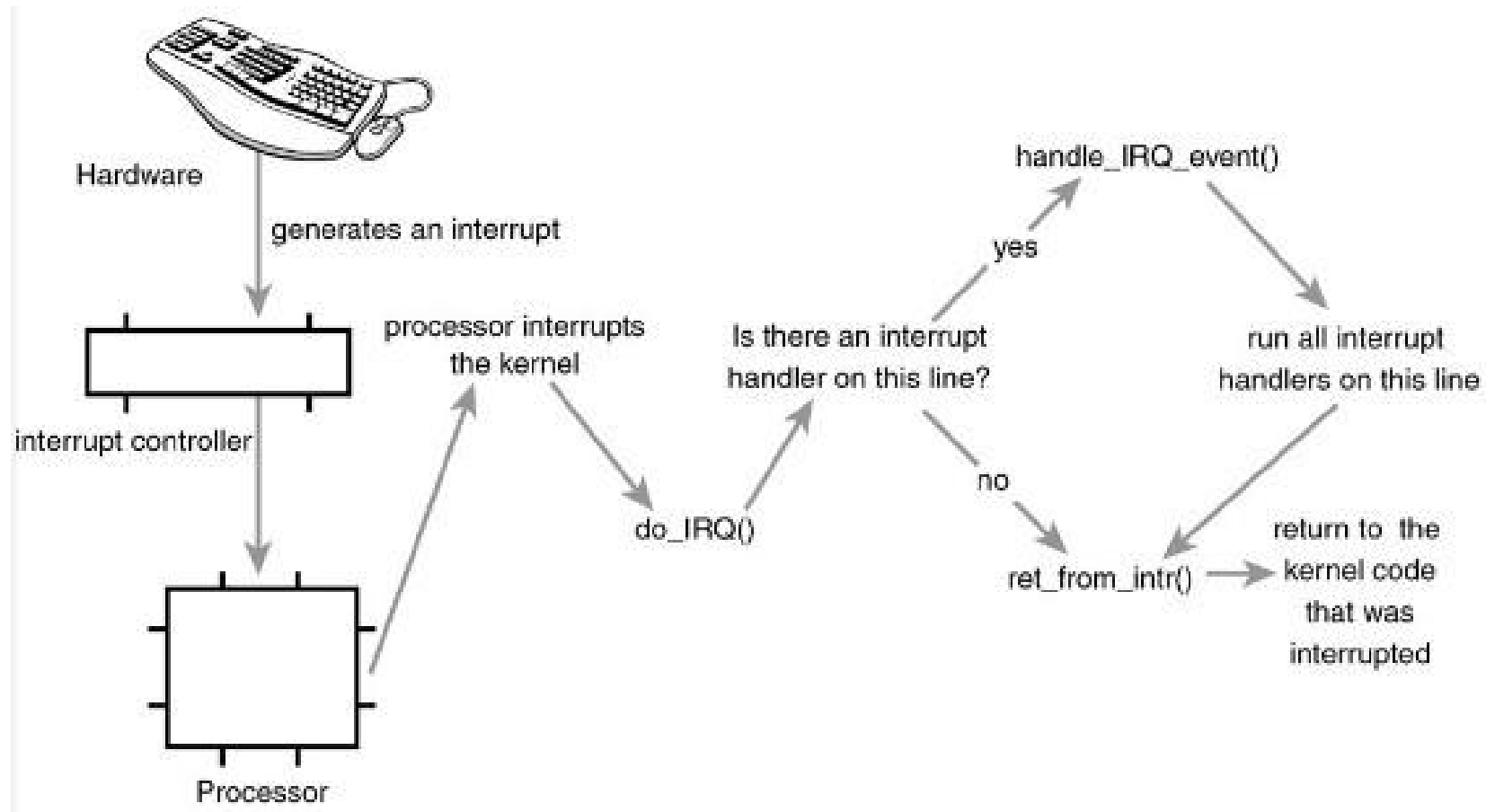
- n **CONFIG\_VECTORS\_BASE**是一个宏，用来获取ARM异常向量的地址，该宏在/arch/arm/include/asm/system.h中定义：
- n `#define CPU_ARCH_ARMv5 4`
- n `#define CR_V (1 << 13)`
- n `extern unsigned long cr_no_alignment;`
- n `extern unsigned long cr_alignment;`
- n `#if __LINUX_ARM_ARCH__ >= 4`
- n `#define vectors_high() (cr_alignment & CR_V)`
- n `#else`
- n `#define vectors_high() (0)`
- n `#endif`
- n 对于ARMv4以下的版本，这个地址固定为0；
- n ARMv4及其以上的版本，ARM异常向量表的地址受协处理器CP15的c1寄存器中V位(bit[13])控制：
- n 如果V=0，则异常向量表的地址为:0x00000000~0x0000001C；
- n 如果V=1，则异常向量表的地址为:0xFFFF0000~0xFFFF001C。

- Linux把中断处理分为两半部分，也就是所谓的上下半部的处理。
- 上半部处理非常紧急的事情，如从硬件读数据，恢复硬件的状态等，这部分要非常快完成，在这部分处理中，对应的中断被屏蔽。
- 下半部处理指把中断处理中不非常紧急的处理延后到一个合适的时间执行，如把读到的数据放进队列，唤醒等待的进程等。
- 上半部与下半部的处理的主要区别在于中断是否被屏蔽，下半部主要由上半部安装调度。





# 中断处理



# struct irq\_desc



```
struct irq_desc {
    struct irq_data          irq_data;
    struct timer_rand_state  *timer_rand_state;
    unsigned int __percpu    *kstat_irqs;
    irq_flow_handler_t       handle_irq;
#ifdef CONFIG_IRQ_PREFLOW_FASTEOI
    irq_preflow_handler_t    preflow_handler;
#endif
    struct irqaction          *action;          /* IRQ action list */
    unsigned int              status_use_accessors;
    unsigned int              core_internal_state__do_not_mess_with_it;
    unsigned int              depth;            /* nested irq disables */
    unsigned int              wake_depth;       /* nested wake enables */
    unsigned int              irq_count;        /* For detecting broken IRQs */
    unsigned long             last_unhandled; /* Aging timer for unhandled count */
    unsigned int              irqs_unhandled;
    raw_spinlock_t            lock;
#ifdef CONFIG_SMP
    const struct cpumask      *affinity_hint;
    struct irq_affinity_notify *affinity_notify;
#endif
#ifdef CONFIG_GENERIC_PENDING_IRQ
    cpumask_var_t             pending_mask;
#endif
    unsigned long             threads_oneshot;
    atomic_t                  threads_active;
    wait_queue_head_t         wait_for_threads;
#ifdef CONFIG_PROC_FS
    struct proc_dir_entry     *dir;
#endif
    const char                *name;
} ____cacheline_internodealigned_in_smp;
```

## struct irq\_desc



- n 在include/linux/irq.h中定义
- n 在内核中，每个中断向量都有相应的有一个irq\_desc结构体。所有这样的描述符组织在一起形成irq\_desc[NR\_IRQS]数组。
- n irq\_data: 主要用来保存中断号irq和chip相关的数据；
- n struct irq\_data {
  - n unsigned int irq; //中断号
  - n unsigned int node;
  - n unsigned int state\_use\_accessors;
  - n struct irq\_chip \*chip; //底层中断控制器描述符
  - n void \*handler\_data;
  - n void \*chip\_data;
  - n struct msi\_desc \*msi\_desc;
  - n #ifdef CONFIG\_SMP
  - n cpumask\_var\_t affinity;
  - n #endif
  - n };

## struct irq\_desc



- n `*kstat_irqs`: 用于系统中断的统计计数。
- n `handle_irq`: 是个函数指针，指向该IRQ线的公共服务程序。
- n `action`: 针对一某具体设备的中断处理对象，设备驱动程序会通过`request_irq`向其中挂载设备特定的中断处理函数。它指向一个单链表，该单链表是由该中断线上所有中断服务程序（对应`struct irqaction`）所连接起来的。
- n `status_use_accessors`: 描述中断线当前的状态；由一组位掩码组成，如：`IRQS_ONESHOT`、`IRQS_WAITING`和`IRQS_PENDING`等。
- n `Lock`: 自旋锁。
- n `*name`: `/proc/interrupts` 中显示的中断名称；
- n `depth`: 中断线被激活时，值为0；其值为正数时，表示被禁止的次数；
- n `irq_count`: 记录该中断线发生中断的次数；
- n `irqs_unhandled`: 该IRQ线上未处理中断发生的次数；

## struct irqaction



n 当多个设备共享一条IRQ线时，因为每个设备都要有各自的ISR。为了能够正确处理此条IRQ线上的中断处理程序（也就是区分每个设备），就需要使用irqaction结构体。在这个结构体中，会有专门的handler字段指向该设备的真正的ISR。共享同一条IRQ线上的多个这样的结构体会连接成了一个单链表，即所谓的中断请求队列。中断产生时，该IRQ线的中断请求队列上所有的ISR都会被依次调用。

n irqaction结构体的定义如下：

n 在include/linux/interrupt.h中定义

```
n struct irqaction {  
n     irq_handler_t      handler;  
n     unsigned long      flags;  
n     void               *dev_id;  
n     struct irqaction   *next;  
n     int                irq;  
n     irq_handler_t      thread_fn;  
n     struct task_struct  *thread;  
n     unsigned long      thread_flags;  
n     unsigned long      thread_mask;  
n     const char         *name;  
n     struct proc_dir_entry *dir;  
n } ____cacheline_internodealigned_in_smp;
```

## struct irqaction



- n **handler:** 指向一个具体的硬件设备的中断服务例程，可以从此指针的类型发现与前文我们所定义的中断处理函数声明相同；
- n **flags:** 对应request\_irq函数中所传递的第三个参数，可取IRQF\_DISABLED、IRQF\_SAMPLE\_RANDOM和IRQF\_SHARED其中之一；
- n **name:** 对应于request\_irq函数中所传递的第四个参数，可通过/proc/interrupts文件查看到；
- n **next:** 指向下一个irqaction结构体；
- n **dev\_id:** 对应于request\_irq函数中所传递的第五个参数，可取任意值，但必须唯一能够代表发出中断请求的设备，通常取描述该设备的结构体；
- n **irq:** 中断号
- n 如果一个IRQ线上有中断请求，那么内核将依次调用在该中断线上注册的每一个中断服务程序，但是并不是所有中断服务程序都被执行。一般硬件设备都会提供一个状态寄存器，以便中断服务程序进行检查是否应该为这个硬件服务。也就是说在整个中断请求队列中，最多会有一个ISR被执行，也就是该ISR对应的那个设备产生了中断请求时；不过当该IRQ线上某个设备未找到匹配的ISR时，那这个中断就不会被处理。此时irq\_desc结构中的irqs\_unhandled字段就会加1。

# struct irq\_chip



```
struct irq_chip {
    const char    *name;
    unsigned int  (*startup)(unsigned int irq);
    void          (*shutdown)(unsigned int irq);
    void          (*enable)(unsigned int irq);
    void          (*disable)(unsigned int irq);
    void          (*ack)(unsigned int irq);
    void          (*mask)(unsigned int irq);
    void          (*mask_ack)(unsigned int irq);
    void          (*unmask)(unsigned int irq);
    void          (*eoi)(unsigned int irq);
    void          (*end)(unsigned int irq);
    void          (*set_affinity)(unsigned int irq, const struct cpumask *dest);
    int           (*retrigger)(unsigned int irq);
    int           (*set_type)(unsigned int irq, unsigned int flow_type);
    int           (*set_wake)(unsigned int irq, unsigned int on);

    /* Currently used only by UML, might disappear one day.*/
#ifdef CONFIG_IRQ_RELEASE_METHOD
    void          (*release)(unsigned int irq, void *dev_id);
#endif

    /*
     * For compatibility, ->typename is copied into ->name.
     * Will disappear.
     */
    const char    *typename;
};
```

- n **struct irq\_chip**是一个中断控制器的描述符。通常不同的体系结构就有一套自己的中断处理方式。内核为了统一的处理中断，提供了底层的中断处理抽象接口，对于每个平台都需要实现底层的接口函数。这样对于上层的中断通用处理程序就无需任何改动。
- n 比如s5pv210的外部中断控制器s5p\_irq\_eint。
- n linux/arch/arm/plat-s5p/irq-eint.c

```
static struct irq_chip s5p_irq_eint = {
    .name      = "s5p-eint",
    .mask      = s5p_irq_eint_mask,
    .unmask    = s5p_irq_eint_unmask,
    .mask_ack  = s5p_irq_eint_maskack,
    .ack       = s5p_irq_eint_ack,
    .set_type  = s5p_irq_eint_set_type,
#ifdef CONFIG_PM
    .set_wake  = s3c_irqext_wake,
#endif
};
```



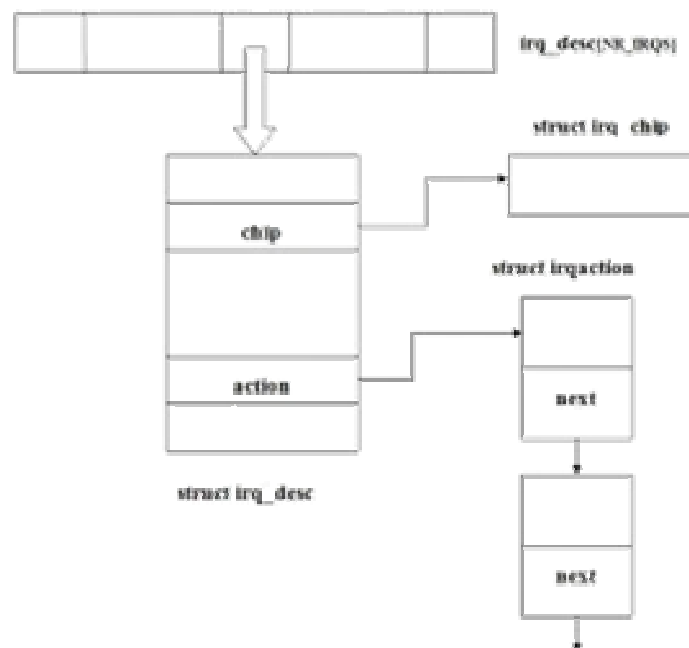
## struct irq\_chip



```
00031: static inline void s5p_irq_eint_mask(struct irq_data *data)
00032: {
00033:     u32 mask;
00034:
00035:     mask = __raw_readl(SSP_EINT_MASK(EINT_REG_NR(data->irq)));
00036:     mask |= eint_irq_to_bit(data->irq);
00037:     __raw_writel(mask, SSP_EINT_MASK(EINT_REG_NR(data->irq)));
00038: }
```

```
00040: static void s5p_irq_eint_unmask(struct irq_data *data)
00041: {
00042:     u32 mask;
00043:
00044:     mask = __raw_readl(SSP_EINT_MASK(EINT_REG_NR(data->irq)));
00045:     mask &= ~(eint_irq_to_bit(data->irq));
00046:     __raw_writel(mask, SSP_EINT_MASK(EINT_REG_NR(data->irq)));
00047: }
-----
```

## 三个数据结构的关系



中断处理程序的调用过程：在asm\_do\_IRQ函数中，通过对irq\_desc结构体中handler\_irq字段的引用，调用handler\_irq所指向的公共服务程序；在这个公共服务程序中会调用hand\_IRQ\_event函数；在hand\_IRQ\_event函数中，通过对irqaction结构体中handler字段的引用最终调用我们所写的中断处理程序。

struct irq\_chip描述了中断最底层的部分；而struct irqaction则描述最上层具体的中断处理函数；而与中断向量所对应的struct irq\_desc则类似一个中间层，将中断中的硬件相关的部分和软件相关的部分连接起来。

# IRQ中断申请



```
#include <linux/interrupt.h>
static inline int __must_check
request_irq(unsigned int irq, irq_handler_t handler, unsigned long flags,
            const char *name, void * dev_id)
```

- n irq----要申请的硬件中断号
- n handle----r是向系统登记的中断处理程序(上半部)，是一个回调函数，中断发生时，系统调用它，将dev\_id参数传递给它。
- n irqflags----中断处理的属性,可以指定中断的触发方式和处理方式。
  - q 触发方式:  
IRQF\_TRIGGER\_NONE 、 IRQF\_TRIGGER\_RISING、 IRQF\_TRIGGER\_FALLING  
、  
IRQF\_TRIGGER\_HIGH、 IRQF\_TRIGGER\_LOW、 IRQ\_TYPE\_EDGE\_BOTH
  - q 处理方式:  
IRQF\_DISABLED - keep irqs disabled when calling the action handler.  
IRQF\_SAMPLE\_RANDOM - irq is used to feed the random generator  
IRQF\_SHARED - allow sharing the irq among several devices  
IRQF\_TIMER - Flag to mark this interrupt as timer interrupt  
IRQF\_NO\_SUSPEND - Do not disable this IRQ during suspend  
IRQF\_NO\_THREAD - Interrupt cannot be threaded  
IRQF\_ONESHOT - Interrupt is not reenabled after the hardirq handler finished.  
Used by threaded interrupts which need to keep the irq line disabled until the threaded handler has been run

- n name----设置中断名称，在cat /proc/interrupts中可以看到此名称。
- n void \* dev\_id ----用作共享中断线的指针。一般设置为这个设备的设备结构体或者NULL。它是一个独特的标识,用在当释放中断线时以及可能还被驱动用来指向它自己的私有数据区，来标识哪个设备在中断。这个参数在真正的驱动程序中一般是指向设备数据结构的指针。在调用中断处理程序的时候它就会传递给中断处理程序的void \* dev\_id。如果中断没有被共享， dev\_id可以设置为 NULL。
- n request\_irq的返回值：
  - q 为0表示成功，
  - q 返回-EINVAL表示中断号无效，
  - q 返回-EBUSY表示中断已经被占用，且不能共享

释放IRQ

```
void free_irq(unsigned int irq,void *dev_id);
```

## 1、中断处理程序的限制：

在中断时间内运行，不能向用户空间发送或者接收数据。

不能做任何导致休眠的操作。

不能调用**schedule**函数。

无论快速还是慢速中断处理例程，都应该设计成执行时间尽可能短。

## 2、上半部的中断处理函数的参数和返回值

**irqreturn\_t (\*handler)(int irq, void \*dev\_id)**

irq----中断号

dev\_id----驱动程序可用的数据区，通常可传递指向描述设备的数据结构指针。

**typedef irqreturn\_t (\*irq\_handler\_t)(int, void \*);**

```
/**
 * enum irqreturn
 * @IRQ_NONE          interrupt was not from this device
 * @IRQ_HANDLED       interrupt was handled by this device
 * @IRQ_WAKE_THREAD    handler requests to wake the handler thread
 */
enum irqreturn {
    IRQ_NONE           = (0 << 0),
    IRQ_HANDLED        = (1 << 0),
    IRQ_WAKE_THREAD    = (1 << 1),};
```

- n `/proc/interrupts`反映系统的中断信息
  - q 第一列是 IRQ 号
  - q 给出每个中断线发生中断的次数。
  - q 给出处理中断的可编程中断控制器。
  - q 给出在该中断号上注册中断处理例程的设备名称。
  
- n `proc/stat`记录了几个关于系统活动的底层统计信息, 包括(但不仅限于)自系统启动以来收到的中断数。

## n 驱动禁止和使能特定的中断

q `#include <asm/irq.h>.`

q `void disable_irq(unsigned int irq);`

q `irq`为中断号，使所选择的中断线无效。使一个中断栈无效。这个函数要等待任何挂起的处理程序在退出之前已经完成。如果你在使用这个函数，同时还持有IRQ处理程序可能需要的一个资源，那么，你就可能死锁。要小心地从IRQ的上下文中调用这个函数。

q `void enable_irq(unsigned int irq);`

q `irq`为中断号，重新启用这条IRQ线上的中断处理。可以在IRQ的上下文中调用这个函数。

q `void disable_irq_nosync(unsigned int irq);`

q `irq`为中断号，使所选择的中断线无效。使一个中断栈无效。与`disable_irq`不同，这个函数并不确保IRQ处理程序的现有实例在退出前已经完成，将立即关闭中断。可以从IRQ的上下文中调用该函数。

- q `#include <asm/irq.h>.`

- n 禁止所有中断

- q `void local_irq_save(unsigned long flags);`

- 将把当前中断状态保存到flags中，然后禁用当前处理器上的中断发送。注意, flags 被直接传递, 而不是通过指针来传递。

- q `void local_irq_disable(void);`

- 不保存状态而关闭本地处理器上的中断发送; 只有我们知道中断并未在其他地方被禁用的情况下，才能使用这个版本。

- n 打开所有中断

- q `void local_irq_restore(unsigned long flags);`

- 恢复由 local\_irq\_save 存储于 flags 的状态, 而 local\_irq\_enable 无条件打开中断.

- q `void local_irq_enable(void);`



## n **S5PV210**外部中断中断号

n 在头文件: linux/arch/arm/plat-s5p/include/plat/irqs.h

n `#define IRQ_EINT(x) ((x) < 16 ? ((x) + S5P_EINT_BASE1) \`  
`: ((x) - 16 + S5P_EINT_BASE2))`

n x----0~31。

q	<code>#define S5P_EINT_BASE1</code>	<code>(S5P_IRQ_VIC0(0))</code>
q	<code>#define S5P_IRQ_VIC0(x)</code>	<code>(S5P_VIC0_BASE + (x))</code>
q	<code>#define S5P_VIC0_BASE</code>	<code>S5P_IRQ(0)</code>
q	<code>#define S5P_IRQ(x)</code>	<code>((x) + S5P_IRQ_OFFSET)</code>
q	<code>#define S5P_IRQ_OFFSET</code>	<code>(32)</code>
q	<code>#define S5P_EINT_BASE2</code>	<code>(IRQ_VIC_END + 1)</code>
q	<code>#define IRQ_VIC_END</code>	<code>S5P_IRQ_VIC3(31)</code>
q	<code>#define S5P_IRQ_VIC3(x)</code>	<code>(S5P_VIC3_BASE + (x))</code>
q	<code>#define S5P_VIC3_BASE</code>	<code>S5P_IRQ(96)</code>

## n **S5PV210 timer**中断号

n 在头文件: linux/arch/arm/plat-s5p/include/plat/irqs.h

n #define IRQ\_TIMER0 S5P\_TIMER\_IRQ(0)

n #define IRQ\_TIMER1 S5P\_TIMER\_IRQ(1)

n #define IRQ\_TIMER2 S5P\_TIMER\_IRQ(2)

n #define IRQ\_TIMER3 S5P\_TIMER\_IRQ(3)

n #define IRQ\_TIMER4 S5P\_TIMER\_IRQ(4)

q S5P\_TIMER\_IRQ(x) (11 + (x))

## n S5PV210 UART号

n 在头文件: linux/arch/arm/plat-s5p/include/plat/irqs.h

n #define IRQ\_S3CUART\_RX0                      IRQ\_S5P\_UART\_RX0

n #define IRQ\_S3CUART\_RX1                      IRQ\_S5P\_UART\_RX1

n #define IRQ\_S3CUART\_RX2                      IRQ\_S5P\_UART\_RX2

n #define IRQ\_S3CUART\_RX3                      IRQ\_S5P\_UART\_RX3

q #define IRQ\_S5P\_UART\_RX0                      (IRQ\_S5P\_UART\_BASE0 +UART\_IRQ\_RXD)

q #define IRQ\_S5P\_UART\_BASE0                      (16)

q #define UART\_IRQ\_RXD                              (0)

- n **S5PV210 MMC**中断号
- n `/* linux/arch/arm/mach-s5pv210/include/mach/irqs.h`
- n `#define IRQ_HSMMMC0` `S5P_IRQ_VIC1(26)`
- n `#define IRQ_HSMMMC1` `S5P_IRQ_VIC1(27)`
- n `#define IRQ_HSMMMC2` `S5P_IRQ_VIC1(28)`
  
- q `#define S5P_IRQ_VIC1(x)` `(S5P_VIC1_BASE + (x))`
- q `#define S5P_VIC1_BASE` `S5P_IRQ(32)`
- q `#define S5P_IRQ(x)` `((x) + S5P_IRQ_OFFSET)`
- q `#define S5P_IRQ_OFFSET` `(32)`

1. 发生中断时，ARM执行异常向量`vector_irq`的代码。
2. 在`vector_irq`内，最终会调用中断处理的总入口函数`asm_do_IRQ`。
3. `asm_do_IRQ`根据中断号调用`irq_desc`数组项的`handle_irq`。
4. `handle_irq`会使用`chip`成员中的函数来设置硬件，如清中断、禁止中断、重新使能中断等。
5. `handle_irq`调用`irq_action`中注册的处理函数`handler`。
  -

# 中断处理程序的局限性



- n 负责对硬件做出迅速响应并完成时间要求很严格的操作
  - q 异步方式执行，可能会打断其他重要代码的执行
  - q 中断处理程序执行过程中
    - n 最好情形下，同级中断被屏蔽
    - n 最坏情形下，当前处理器上所有其他中断都会被屏蔽
  - q 中断处理程序不在进程上下文中运行，不能被阻塞
- n 中断处理程序
  - q 上半部：中断处理程序
    - n 简单快速，执行时禁止部分或全部中断
  - q 下半部：推后工作的机制
    - n 执行与中断处理密切相关但中断处理程序本身不执行的工作
  - q 执行期间可以响应中断

-- Linux把中断处理程序分成了两部分：

- n 上半部：实际响应中断的处理程序。
- n 下半部：被上半部调用，通过开中断的方式进行。

q 下半部两种机制实现：

- n Tasklet
  - n 工作队列work queue
- 
- n 上半部的功能是“登记中断”，当一个中断发生时，它进行相应地硬件读写后就把中断例程的下半部挂到该设备的下半部执行队列中去。因此，上半部执行的速度就会很快，可以服务更多的中断请求。
  - n 但是，仅有“登记中断”是远远不够的，因为中断的事件可能很复杂。因此，Linux引入了一个下半部，来完成中断事件的绝大多数使命。
  - n 下半部和上半部最大的不同是下半部是可中断的，而上半部是不可中断的，下半部几乎做了中断处理程序所有的事情，而且可以被新的中断打断！下半部则相对来说并不是非常紧急的，通常还是比较耗时的，因此由系统自行安排运行时机，不在中断服务上下文中执行。

# tasklet的实现



tasklet数据结构为struct tasklet\_struct，每一个结构体代表一个独立的小任务，在<linux/interrupt.h>定义如下

n

```
struct tasklet_struct
{
    struct tasklet_struct *next; /*指向下一个链表结构*/
    unsigned long state; /*小任务状态*/
    atomic_t count; /*引用计数器*/
    void (*func)(unsigned long); /*小任务的处理函数*/
    unsigned long data; /*传递小任务函数的参数*/
};
```

n state的取值参照下边的枚举型：

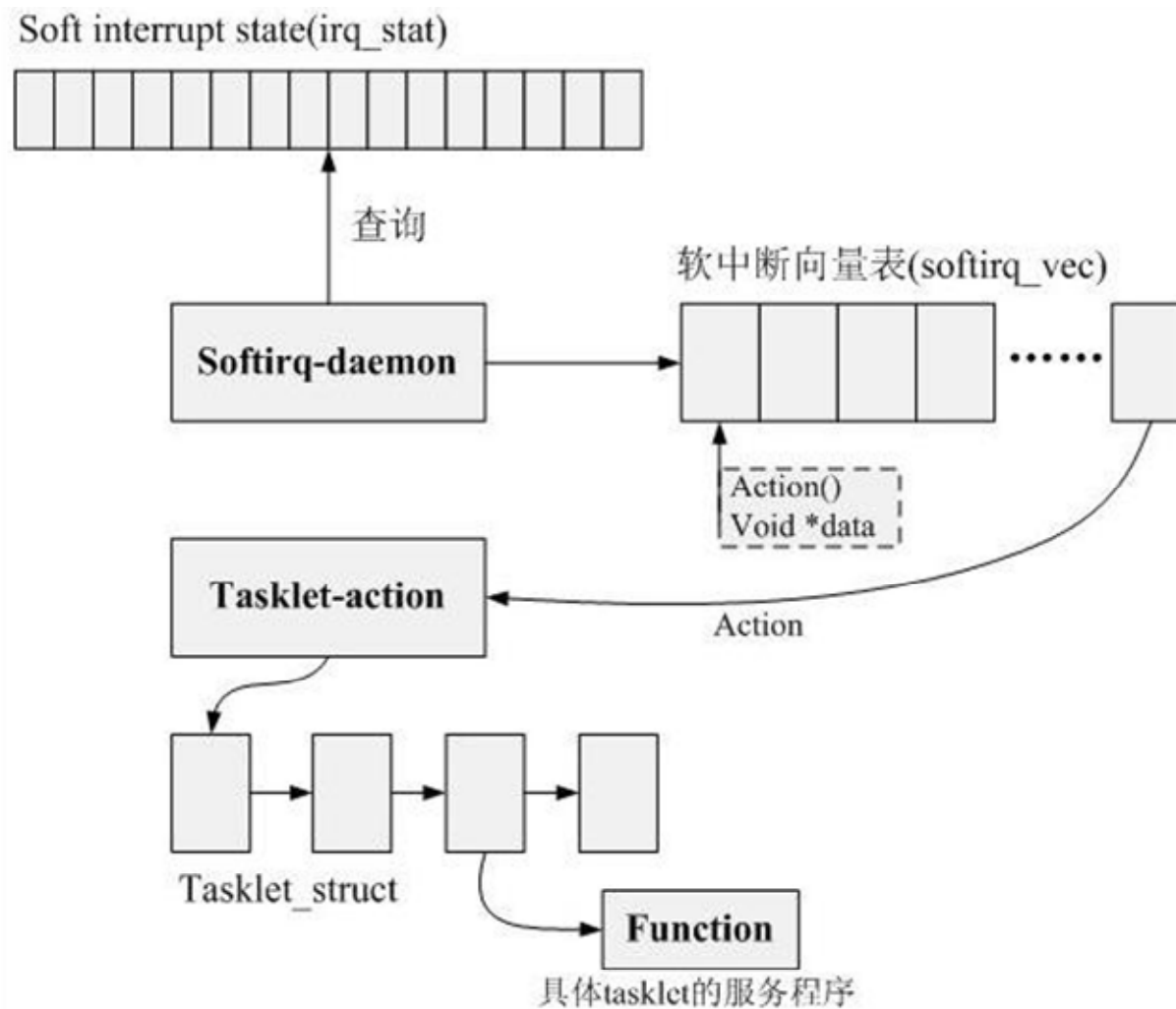
```
enum
{
    TASKLET_STATE_SCHED, /*小任务已被调度，等待运行*/
    TASKLET_STATE_RUN /*正在运行，仅在多处理器上使用*/
};
```

n count域是小任务的引用计数器。只有当它的值为0的时候才能被激活，并其被设置为挂起状态时，才能够被执行，否则为禁止状态。



- n **ksoftirqd**是一个后台运行的内核线程，它会周期的遍历软中断的向量列表，如果发现哪个软中断向量被挂起了，就执行对应的处理函数。
- n **tasklet** 所对应的处理函数就是**tasklet\_action**，这个处理函数在系统启动时初始化软中断时，就在软中断向量表中注册。
- n **tasklet\_action()** 遍历一个全局的 **tasklet\_vec** 链表。链表中的元素为 **tasklet\_struct**结构体。

# tasklet运行



# tasklet程序设计



- `#include <linux/interrupt.h>`
- n 1、声明和使用小任务tasklet
- n 静态的创建一个tasklet:
  - q `#define DECLARE_TASKLET(name, func, data) \`  
`struct tasklet_struct name = { NULL, 0, ATOMIC_INIT(0), func, data }`
  - q `#define DECLARE_TASKLET_DISABLED(name, func, data) \`  
`struct tasklet_struct name = { NULL, 0, ATOMIC_INIT(1), func, data }`
- n name 是 tasklet 的名字,
- n func 是执行 tasklet 的函数;
- n data 是 unsigned long 类型的 func 参数。
- n 这两个宏的区别在于计数器设置的初始值不同, 前者为0, 后者为1。为0的表示激活状态, 为1的表示禁止状态。
- n 动态创建一个tasklet
- n `void tasklet_init(struct tasklet_struct *t,`  
`void (*func)(unsigned long), unsigned long data);`
- n 例:
  - q `static struct tasklet_struct my_tasklet;`
  - q `tasklet_init(&my_tasklet,tasklet_handler,0);`

## n 2、小任务处理函数程序

n `void tasklet_handler(unsigned long data)`

## n tasklet应用注意

- q tasklet不能睡眠，不能在tasklet\_handler内部使用信号量或其他阻塞式函数。
- q 由于tasklet运行时允许响应中断，如果tasklet与中断处理程序之间共享了某些数据，必须做好预防工作（锁机制）。
- q 不允许两个两个相同类型的tasklet同时执行。
- q 一个tasklet不会抢占另外一个tasklet。
- q 索引号小的tasklet在索引号大的tasklet之前执行。
- q tasklet被调度以后，只要有机会它就会尽可能早地运行。在它还没有得到运行机会之前，如果有一个相同的tasklet又被唤醒了，那么它只会运行一次。

## n 3、调度tasklet

n 调度小任务时引用调度函数就能使系统在合适的时候进行调度。函数原型为：

```
n static inline void tasklet_schedule(struct tasklet_struct *t)
{
    if (!test_and_set_bit(TASKLET_STATE_SCHED, &t->state))
        __tasklet_schedule(t);
}
```

q 调度执行指定的tasklet。

q 将定义后的 tasklet 挂接到 cpu 的 tasklet\_vec 链表。而且会引起一个软 tasklet 的软中断，既把 tasklet 对应的中断向量挂起 (pend) 。

n 调度函数放在中断处理的上半部处理函数中，这样中断申请的时候调用处理函数（即 irq\_handler\_t irq\_handler）后，转去执行下半部的小任务。

## n 4、删除tasklet

n tasklet\_kill(&tasklet);

从挂起的队列中去掉一个tasklet。在处理一个经常重新调度它自身的tasklet的时候，从挂起的队列中移去已调度的tasklet会很有用。这个函数首先等待该tasklet执行完毕，然后再将它移去。由于该函数可能会引起休眠，所以禁止在中断上下文中使用它。

n 使用tasklet作为下半部的处理中断的设备驱动程序模板如下：

n /\*中断处理下半部\*/

```
n void my_do_tasklet(unsigned long)
{
    ...../*编写处理事件方法*/
}
```

n /\*定义tasklet和下半部函数并关联\*/

```
n DECLARE_TASKLET(my_tasklet, my_do_tasklet, 0);
```

n /\*中断处理上半部\*/

```
n irqreturn_t my_interrupt(unsigned int irq,void *dev_id)
{
    .....
    /*调度my_tasklet函数，根据声明将去执行my_tasklet_func函数*/
    tasklet_schedule(&my_tasklet)
    .....
}
```

```
n  /*设备驱动的加载函数*/
n  int xxx_open(...)
  {
      .....
      /*申请中断, 转去执行my_interrupt函数并传入参数*/
      result=request_irq(my_irq,my_interrupt,IRQF_DISABLED,"xxx",NULL);
      .....
  }

n  /*设备驱动模块的卸载函数*/
n  void xxx_release(...)
  {
      .....
      /*释放中断*/
      free_irq(my_irq,my_interrupt);
      .....
  }
```

- n **void tasklet\_enable(struct tasklet\_struct \*t)**  
调用**tasklet\_enable()**函数可以激活一个**tasklet**，要激活**DECLARE\_TASKLET\_DISABLED()**创建的**tasklet**，也要调用这个函数。
- n **void tasklet\_disable(struct tasklet\_struct \*t)**  
**tasklet\_disable()**函数用来禁止某个指定的**tasklet**。如果该**tasklet**当前正在执行，这个函数会等到它执行完毕再返回。
- n **void tasklet\_disable\_nosync(struct tasklet\_struct \*t)**  
**tasklet\_disable\_nosync()**函数也可以用来禁止指定的**tasklet**，不过它无需在返回前等待**tasklet**执行完毕。这样做往往不太安全，因为我们无法估计该**tasklet**是否仍在执行。
- n **void tasklet\_kill(struct tasklet\_struct \*t)**  
通过调用**tasklet\_kill()**函数从挂起的队列中去掉一个**tasklet**。该函数的参数是一个指向某个**tasklet**的**tasklet\_struct**的长指针。在处理一个经常重新调度它自身的**tasklet**的时候，从挂起的队列中移去已调度的**tasklet**会很有用。这个函数首先等待该**tasklet**执行完毕，然后再将它移去。由于该函数可能会引起休眠，所以禁止在中断上下文中使用它。



- n 工作队列（**work queue**）是另外一种将中断的部分工作推后的一种方式，它可以实现一些**tasklet**不能实现的工作，比如工作队列机制可以睡眠。
- n 工作队列将推后的工作交给一个称之为**工作者线程（worker thread）**的内核线程去完成（单核下一般会交给默认的**线程events/0**）。因此，在该机制中，当内核在执行中断的剩余工作时就处在进程上下文（**process context**）中。
- n 工作队列可以重新调度甚至睡眠。
- n 对于**tasklet**机制（中断处理程序也是如此），内核在执行时处于中断上下文（**interrupt context**）中。而中断上下文与进程毫无瓜葛，所以在中断上下文中不能睡眠。
- n 因此，当推后的那部分中断程序需要睡眠时，工作队列毫无疑问是最佳选择；否则用**tasklet**。

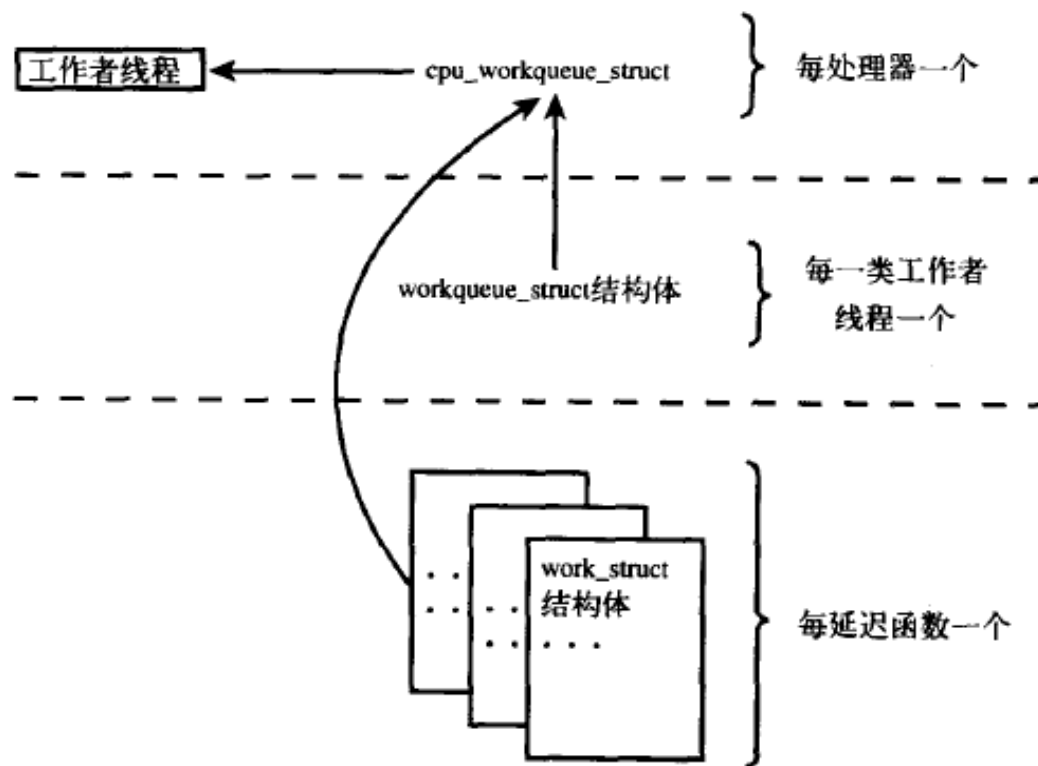
- n 我们把推后执行的任务叫做工作（work），描述它的数据结构为work\_struct。
- n 这些工作以队列结构组织成工作队列（workqueue），其数据结构为workqueue\_struct。
- n 工作者线程就是负责执行工作队列中的工作。系统默认的工作者线程为events，自己也可以创建自己的工作线程。

- n 工作队列子系统是一个用于创建内核线程（工作者线程）的接口，通过它创建的进程负责执行由内核其他部分排队到队列里的任务
- n 工作队列子系统提供一个默认的工作者线程（events/n）来处理需推后的工作
  - q 其中n是处理器的编号，每个处理器对应一个线程
  - q 工作队列最基本的形式
    - n 将需要推后执行的任务交给特定的通用线程
- n 工作队列可以让驱动程序创建一个专门的工作者线程来处理需要推后的工作
  - q 主要适用于处理器密集型和性能要求严格的任务

# 工作队列的核心组成

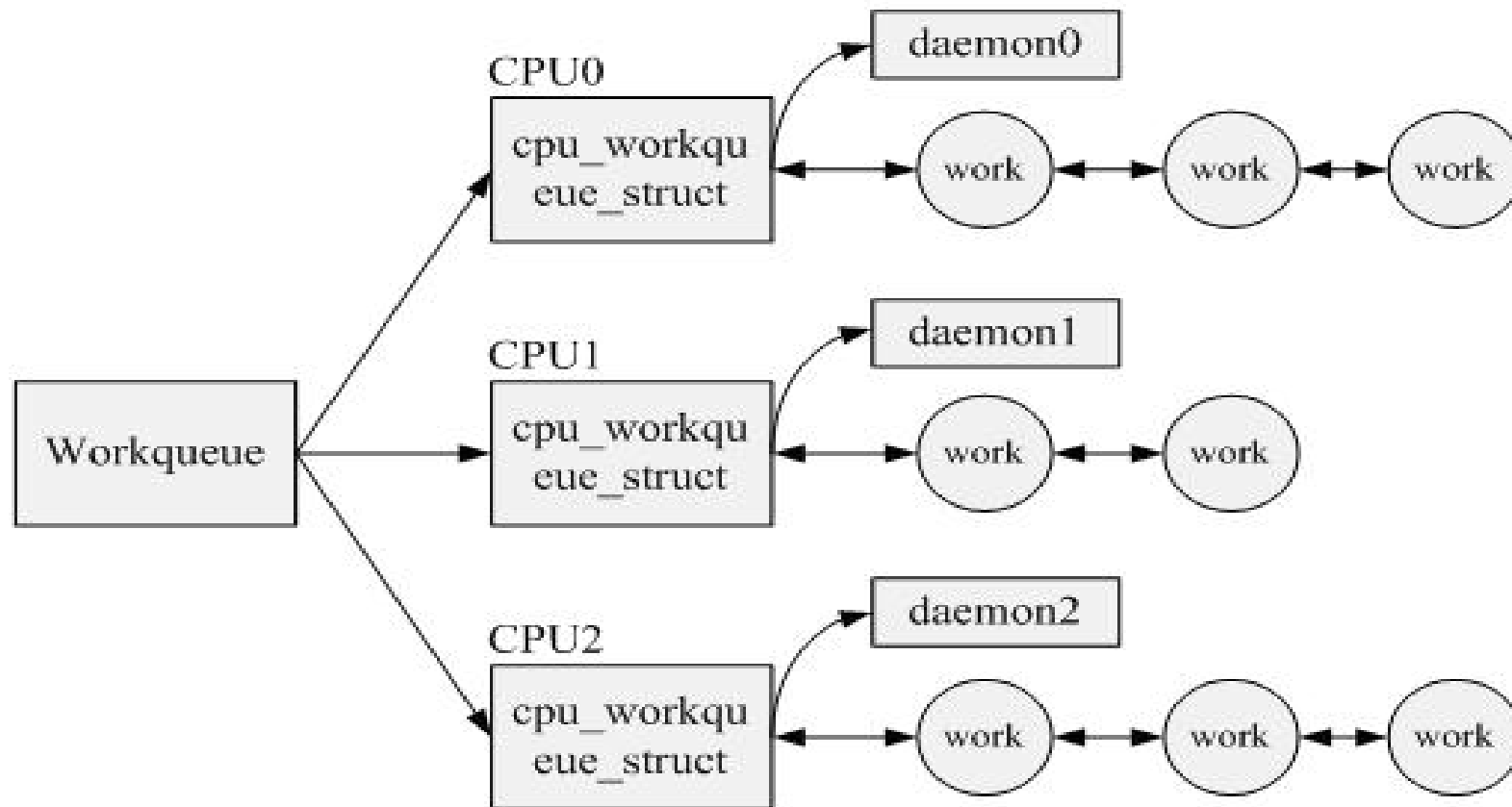


- n 工作者线程
  - q 工作队列: `workqueue_struct`
  - q CPU相关工作队列: `cpu_workqueue_struct`
  - q 任务: `work_struct`



- n 工作队列在进程上下文中执行；
- n 可以阻塞；
- n 可以被重新调度；
- n 使用工作队列的两种形式：
  - q 缺省工作者线程
  - q 自建的工作者线程
- n 在工作队列和内核其他部分之间使用锁机制就像在其他的进程上下文一样；
- n 默认允许响应中断；
- n 默认不持有任何锁。

# workqueue



工作队列work\_struct结构体，位于/include/linux/workqueue.h

```
n typedef void (*work_func_t) (struct work_struct *work);

n struct work_struct {
    atomic_long_t data; /*传递给处理函数的参数*/
    #define WORK_STRUCT_PENDING 0/*工作是否正在等待处理标志*/
    #define WORK_STRUCT_FLAG_MASK (3UL)
    #define WORK_STRUCT_WQ_DATA_MASK (~WORK_STRUCT_FLAG_MASK)
    struct list_head entry; /* 连接所有工作的链表*/
    work_func_t func; /* 要执行的函数*/
    #ifdef CONFIG_LOCKDEP
        struct lockdep_map lockdep_map;
    #endif
};
```

n 这些工作结构被连接成链表，组成工作队列。当一个工作者线程被唤醒时，它会执行它的链表上的所有工作。工作被执行完毕，它就将相应的work\_struct对象从链表上移去。当链表上不再有对象的时候，它就会继续休眠。可以通过DECLARE\_WORK在编译时静态地创建该结构，以完成推后的工作。

- n `#include <linux/workqueue.h>`
- n 在需要调度工作队列的时候，引用类似tasklet\_schedule()函数如：
- n `int schedule_work(struct work_struct *work);` //put work task in global workqueue
- n 如果有时候并不希望工作马上就被执行，而是希望它经过一段延迟以后再执行。在这种情况下，可以调度指定的时间后执行函数：
- n `int schedule_delayed_work(struct delayed_work *work, unsigned long delay);`
- n 其中是以delayed\_work为结构体的指针，而这个结构体的定义是在work\_struct结构体的基础上增加了一项timer\_list结构体。
- n 

```
struct delayed_work {  
    struct work_struct work;  
    struct timer_list timer; /* 延迟的工作队列所用到的定时器，当不需要延迟时初始化为NULL */  
};
```
- n 这样，便使预设的工作队列直到delay指定的时钟节拍用完以后才会执行。
- n 返回值：
- n Returns zero if @work was already on the kernel-global workqueue and \* non-zero otherwise.



- n `#include <linux/workqueue.h>`
- n 初始化一个work
- n `INIT_WORK(struct work_struct *work, work_func_t func);`
- n `DECLARE_WORK (struct work_struct *work, work_func_t func);`
- n `DECLARE_DELAYED_WORK (struct work_struct *work, work_func_t func);`

n 使用缺省工作队列处理中断下半部的设备驱动程序模板如下：

n /\*中断处理下半部\*/

```
n void my_do_work(unsigned long)
{
    ...../*编写自己的处理事件内容*/
}
```

n /\*定义工作队列和下半部函数并关联\*/

```
n static DECLARE_WORK(my_work, my_do_work);
```

n /\*中断处理上半部\*/

```
n irqreturn_t my_interrupt(unsigned int irq,void *dev_id)
{
    .....
    schedule_work(&my_work)/*调度my_work 处理函数
    .....
}
```

```
n  /*设备驱动的加载函数*/
n  int __init xxx_init(void)
  {
      .....
      /*申请中断,转去执行my_interrupt函数并传入参数*/
      result=request_irq(my_irq,my_interrupt,IRQF_DISABLED,"xxx",NULL);
      .....
  }

n  /*设备驱动模块的卸载函数*/
n  void __exit xxx_exit(void)
  {
      .....
      /*释放中断*/
      free_irq(my_irq,my_interrupt);
      .....
  }
```

# Thank You !

[www.gec-edu.org](http://www.gec-edu.org)