

pkmap_countabs() and abs64()

```
1.  /*
2.   * abs() handles unsigned and signed longs, ints, shorts and chars.  For all
3.   * input types abs() returns a signed long.
4.   * abs() should not be used for 64-bit types (s64, u64, long long) - use abs64()
5.   * for those.
6.   */
7.  #define abs(x) ({                                \
8.          long ret;                                \
9.          if (sizeof(x) == sizeof(long)) {        \ 是否可抢占 (pre
emptable)
10.             long __x = (x);                        \
11.             ret = (__x < 0) ? -__x : __x;          \
12.         } else {                                    \
13.             int __x = (x);                          \
14.             ret = (__x < 0) ? -__x : __x;          \
15.         }                                           \
16.         ret;                                       \
17.     })
18.
19.  #define abs64(x) ({                                \
20.          s64 __x = (x);                            \
21.          (__x < 0) ? -__x : __x;                    \
22.     })
```

判断是否可抢占 (preemptable)

每个process的thread_info都有一个preempt_count。

```

1. struct thread_info {
2.     unsigned long         flags;           /* low level flags */
3.     int                    preempt_count; /* 0 => preemptable, <0 => bug */
4.     mm_segment_t          addr_limit;      /* address limit */
5.     struct task_struct     *task;          /* main task structure */
6.     struct exec_domain     *exec_domain;   /* execution domain */
7.     __u32                  cpu;            /* cpu */
8.     __u32                  cpu_domain;     /* cpu domain */
9.     struct cpu_context_save cpu_context;    /* cpu context */
10.    __u32                   syscall;        /* syscall number */
11.    __u8                    used_cp[16];    /* thread used copro */
12.    unsigned long           tp_value[2];    /* TLS registers */
13. #ifdef CONFIG_CRUNCH
14.     struct crunch_state    crunchstate;
15. #endif
16.     union fp_state         fpstate __attribute__((aligned(8)));
17.     union vfp_state        vfpstate;
18. #ifdef CONFIG_ARM_THUMBEE
19.     unsigned long          thumbee_state; /* ThumbEE Handler Base register
20.     */
21. #endif
22.     struct restart_block   restart_block;
23. };

```

```

1. #define preempt_enable() \
2. do { \
3.     barrier(); \
4.     if (unlikely(preempt_count_dec_and_test())) \
5.         __preempt_schedule(); \
6. } while (0)

```

preempt_count_dec_and_test () check preempt_count - 1后是否为0 , 若是则

```
1.  asmlinkage __visible void __sched notrace preempt_schedule(void)
2.  {
3.      /*
4.       * If there is a non-zero preempt_count or interrupts are disabled,
5.       * we do not want to preempt the current task. Just return..
6.       */
7.      if (likely(!preemptible()))
8.          return;
9.
10.     do {
11.         __preempt_count_add(PREEMPT_ACTIVE);
12.         __schedule();
13.         __preempt_count_sub(PREEMPT_ACTIVE);
14.
15.         /*
16.          * Check again in case we missed a preemption opportunity
17.          * between schedule and now.
18.          */
19.         barrier();
20.     } while (need_resched());
21. }
22. NOKPROBE_SYMBOL(preempt_schedule);
```

define preemptible() (preempt_count() == 0 && !irqs_disabled())

可抢占的判断标准

1. 当前运行process的preempt_count为0 (原来为1 , 减去1后)
2. disable irq

preempt_count值的含义比较复杂

```
1.  /*
2.   * We put the hardirq and softirq counter into the preemption
3.   * counter. The bitmask has the following meaning:
4.   *
5.   * - bits 0-7 are the preemption count (max preemption depth: 256)
6.   * - bits 8-15 are the softirq count (max # of softirqs: 256)
7.   *
8.   * The hardirq count could in theory be the same as the number of
9.   * interrupts in the system, but we run all interrupt handlers with
10.  * interrupts disabled, so we cannot have nesting interrupts. Though
11.  * there are a few palaeontologic drivers which reenable interrupts in
12.  * the handler, so we need more than one bit here.
13.  *
14.  * PREEMPT_MASK:      0x000000ff
15.  * SOFTIRQ_MASK:      0x0000ff00
16.  * HARDIRQ_MASK:      0x000f0000
17.  *      NMI_MASK:      0x00100000
18.  * PREEMPT_ACTIVE:    0x00200000
19.  */
20. #define PREEMPT_BITS    8
21. #define SOFTIRQ_BITS    8
22. #define HARDIRQ_BITS    4
23. #define NMI_BITS        1
```

即只有bit0 - bit7是对preempt的counter , bit8 - bit15 is for softirq

bit16-bit19 is for hardirq, bit20 is for NMI

临时禁止preemptable

因为现在kernel也是部分可抢占的，所以如果kernel中有code希望临时禁止preemptable，可如下

```
preempt_disable();
```

```
do something
```

```
preempt_enable();
```

check kernel是否运行于“特殊”状态

```
1.  #define in_irq()                (hardirq_count())
2.  #define in_softirq()           (softirq_count())
3.  #define in_interrupt()         (irq_count())
4.  #define in_serving_softirq()    (softirq_count() & SOFTIRQ_OFFSET)
```

in_irq() --- hardware interrupt running 判断

```
1.  #define irq_count()             (preempt_count() & (HARDIRQ_MASK | SOFTIRQ_MASK \
2.                                     | NMI_MASK))
```

check preempt_count的bit16 - bit19是否为0来判断是否运行在hardware interrupt中。

in_softirq() --- software irq running 判断

in_interrupt()

```
1.  #define irq_count()      (preempt_count() & (HARDIRQ_MASK | SOFTIRQ_MASK \
2.                                  | NMI_MASK))
```

"interrupt"是 **hardware interrupt, software interrupt and NMI** running的合称。

check 当前是否可抢占

```
1.  #ifdef CONFIG_PREEMPT_COUNT
2.  # define preemptible()  (preempt_count() == 0 && !irqs_disabled())
3.  #else
4.  # define preemptible()  0
5.  #endif
```

enable / disable pagefault

include/linux/uaccess.h

```

1.  /*
2.   * These routines enable/disable the pagefault handler in that
3.   * it will not take any locks and go straight to the fixup table.
4.   *
5.   * They have great resemblance to the preempt_disable/enable calls
6.   * and in fact they are identical; this is because currently there is
7.   * no other way to make the pagefault handlers do this. So we do
8.   * disable preemption but we don't necessarily care about that.
9.   */
10. static inline void pagefault_disable(void)
11. {
12.     preempt_count_inc();
13.     /*
14.      * make sure to have issued the store before a pagefault
15.      * can hit.
16.      */
17.     barrier();
18. }
19.
20. static inline void pagefault_enable(void)
21. {
22. #ifndef CONFIG_PREEMPT
23.     /*
24.      * make sure to issue those last loads/stores before enabling
25.      * the pagefault handler again.
26.      */
27.     barrier();
28.     preempt_count_dec();
29. #else
30.     preempt_enable();
31. #endif
32. }

```

从code上看，更准确的function name应该是current_process_pagefault_disable()

该函数只是disable当前process产生pagefault。因为当产生pagefault时必然伴随着current process要抢占(它要被schedule出去sleep)。

但问题是，pagefault产生了，但又不能被抢占,那current process该怎么办？有点糊涂了？！

要看page fault handler才能明白怎么办。

access kernel memory safely

mm/mmaccess.c

```

1.  /**
2.   * probe_kernel_read(): safely attempt to read from a location
3.   * @dst: pointer to the buffer that shall take the data
4.   * @src: address to read from
5.   * @size: size of the data chunk
6.   *
7.   * Safely read from address @src to the buffer at @dst. If a kernel fault
8.   * happens, handle that and return -EFAULT.
9.   */
10.
11. long __weak probe_kernel_read(void *dst, const void *src, size_t size)
12.     __attribute__((alias("__probe_kernel_read"))));
13.
14. long __probe_kernel_read(void *dst, const void *src, size_t size)
15. {
16.     long ret;
17.     mm_segment_t old_fs = get_fs();
18.
19.     set_fs(KERNEL_DS);
20.     pagefault_disable();
21.     ret = __copy_from_user_inatomic(dst,
22.                                     (__force const void __user *)src, size);
23.     pagefault_enable();
24.     set_fs(old_fs);
25.
26.     return ret ? -EFAULT : 0;
27. }
28. EXPORT_SYMBOL_GPL(probe_kernel_read);
29.
30. /**
31.   * probe_kernel_write(): safely attempt to write to a location
32.   * @dst: address to write to
33.   * @src: pointer to the data that shall be written
34.   * @size: size of the data chunk
35.   *
36.   * Safely write to address @dst from the buffer at @src. If a kernel fault
37.   * happens, handle that and return -EFAULT.
38.   */
39. long __weak probe_kernel_write(void *dst, const void *src, size_t size)
40.     __attribute__((alias("__probe_kernel_write"))));
41.
42. long __probe_kernel_write(void *dst, const void *src, size_t size)
43. {
44.     long ret;
45.     mm_segment_t old_fs = get_fs();
46.
47.     set_fs(KERNEL_DS);
48.     pagefault_disable();
49.     ret = __copy_to_user_inatomic((__force void __user *)dst, src, size);
50.     pagefault_enable();
51.     set_fs(old_fs);
52.
53.     return ret ? -EFAULT : 0;

```



```
54.     }
55.     EXPORT_SYMBOL_GPL(probe_kernel_write);
```

分析如下：

```
pagefault_disable();

ret = __copy_from_user_inatomic(dst,
    (__force const void __user *)src, size);
```

在kernel产生page fault，进入do_page)fault() / arch/arm/mm/fault.c

```
1.  static int __kprobes
2.  do_page_fault(unsigned long addr, unsigned int fsr, struct pt_regs *regs)
3.  {
4.      struct task_struct *tsk;
5.      struct mm_struct *mm;
6.      int fault, sig, code;
7.      unsigned int flags = FAULT_FLAG_ALLOW_RETRY | FAULT_FLAG_KILLABLE;
8.
9.      if (notify_page_fault(regs, fsr))
10.         return 0;
11.
12.     tsk = current;
13.     mm = tsk->mm;
14.
15.     /* Enable interrupts if they were enabled in the parent context. */
16.     if (interrupts_enabled(regs))
17.         local_irq_enable();
18.
19.     /*
20.      * If we're in an interrupt or have no user
21.      * context, we must not take the fault..
22.      */
23.     if (in_atomic() || !mm)
24.         goto no_context;
25.
26.     if (user_mode(regs))
27.         flags |= FAULT_FLAG_USER;
28.     if (fsr & FSR_WRITE)
29.         flags |= FAULT_FLAG_WRITE;
30.     .....
```

运行到

```
if (in_atomic() || !mm)

    goto no_context;
```

由于pagefault_disable()的缘故，这里的in_atomic()返回true，执行流运行如下code

```
1. no_context:
2.     __do_kernel_fault(mm, addr, fsr, regs);
3.     return 0;
```

```
1.  /*
2.   * Oops. The kernel tried to access some page that wasn't present.
3.   */
4.  static void
5.  __do_kernel_fault(struct mm_struct *mm, unsigned long addr, unsigned int fsr,
6.                   struct pt_regs *regs)
7.  {
8.      /*
9.       * Are we prepared to handle this kernel fault?
10.     */
11.     if (fixup_exception(regs))
12.         return;
13.
14.     /*
15.      * No handler, we'll have to terminate things with extreme prejudice.
16.     */
17.     bust_spinlocks(1);
18.     printk(KERN_ALERT
19.            "Unable to handle kernel %s at virtual address %08lx\n",
20.            (addr < PAGE_SIZE) ? "NULL pointer dereference" :
21.            "paging request", addr);
22.
23.     show_pte(mm, addr);
24.     die("Oops", regs, fsr);
25.     bust_spinlocks(0);
26.     do_exit(SIGKILL);
27. }
```

由于

__copy_from_user_inatomic(dst, (__force const void __user *)src, size)会安装 exception handler，所以

这里的fixup_exception(regs))会返回true。

```
1. int fixup_exception(struct pt_regs *regs)
2. {
3.     const struct exception_table_entry *fixup;
4.
5.     fixup = search_exception_tables(instruction_pointer(regs));
6.     if (fixup) {
7.         regs->ARM_pc = fixup->fixup;
8. #ifdef CONFIG_THUMB2_KERNEL
9.         /* Clear the IT state to avoid nasty surprises in the fixup */
10.        regs->ARM_cpsr &= ~PSR_IT_MASK;
11. #endif
12.     }
13.
14.     return fixup != NULL;
15. }
```

