

# Linux 内核学习起步

albcamus <[albcamus@gmail.com](mailto:albcamus@gmail.com)>

Last Update : 2007年12月06日

## 1, 编译内核

---

### 1.0 获取内核压缩

---

在 <http://www.kernel.org/pub/linux/kernel/v2.6/> 能看到一系列的文件, 如 ChangeLog-2.6.xx、linux-2.6.xx.tar.{gz|bz2} 和 patch-2.6.xx.{gz|bz2}, 你应该下载 linux-2.6.23.tar.bz2 这样的压缩包。

### 1.1 找一个旧的.config 作参考

---

安装内核开发包(以 FC8 为例, 是 kernel-devel-2.6.23.1-42.fc8.i686.rpm), 则 /lib/modules/<version>/build 目录下会有 .config 文件。

```
# cd linux-2.6.24-rc3
# cp /lib/modules/2.6.23.1-42.fc8/build/.config .
```

### 1.2 配置、编译和安装

---

配置:

```
# make menuconfig/gconfig/xconfig/oldconfig/defconfig/allyesconfig/allmodconfig
```

FYI: 一般用 menuconfig 比较多, 如果你喜欢 GUI 风格, 也可以选择基于 Qt 的 xconfig 或

基于 Gtk+ 的 gconfig。

FYI: 可以用 `O=</path/to/build>` 指定编译生成的文件放在哪个目录。举例来说, 如果我的内核源代码目录是 `/usr/src/linux-2.6.23`, 编译内核时指定了:

```
make O=/home/arc/build/linux-2.6.23
```

那么 `/lib/modules/2.6.23` 下的 `source` 和 `build` 这两个符号链接就分别指向了源代码和目标代码目录:

```
# ls -l /lib/modules/2.6.23/source
lrwxrwxrwx 1 root root 21 11-22 13:35 /lib/modules/2.6.23/source -
> /usr/src/linux-2.6.23

# ls -l /lib/modules/2.6.23/build
lrwxrwxrwx 1 root root 28 11-22 13:35 /lib/modules/2.6.23/build ->
/home/arc/build/linux-2.6.23
```

编译:

```
# make
```

```
FYI:  V=1    /*verbose 模式, 把每一部执行的命令都打印出来。
      *一个小技巧就是, 把标准输出重定向到一个文件中,
      *这样以后查找模块之间的依赖关系就很方便了
      */

      C=1    /*需要安装 sparse 程序, 进行严格的静态 C 语法检查。
      *一般开发者会通过它来预防 BUG
      */
```

安装:

```
# make modules_install
# make install
```

FYI: 正常情况下 `make install` 会根据你机器的配置为你定制 `initrd` 文件, 并更新 `grub.conf` 文件中的内容。但如果重新引导时发现无法启动, 注意手工 `mkinitrd`。例如:

```
# mkinitrd -v -preload libata.ko -with=ext3 /boot/initrd-2.6.24-rc3.img
2.6.24-rc3
```

其中 `preload` 指定的模块会在 `/etc/modprobe.conf` 之前加载, 而 `--with` 指定的模块会在这之后加载。

## 1.3 文档

---

```
# make htmldocs //你也可以不用 htmldocs 指令 HTML 格式，而指定 pdfdocs 或
//psdocs

# make mandocs //为 kernel API 生成 man 手册
# make installmandocs /*将 kernel API 的手册页安装到 man 程序能找到的
*目录中，这样就可以 man copy_from_user 了
*/
```

FYI: 执行 make htmldocs/pdfdocs/psdocs 之后，在 O=指定的目录(如果没使用 O=则是源代码目录)的 Documentation/DocBook/下，会生成几份很重要的文档：

```
kernel api    : 内核开发的 API 手册
usb           : USB host 端的 API 手册
gadget       : USB device 端的 API 手册
kernel locking : 内核加锁的 HOWTO 文档
kernel hacking : 内核开发的一些注意事项
```

FYI: 内核源代码目录的 Documentation 目录：

```
kernel-parameter.txt : 内核参数，加在一个 grub entry 的 kernel 指令后面
filesystems/vfs.txt : Linux 虚拟文件系统的深入介绍
memory-barriers      : 关于 barriers 的文档
CodingStyle          : 内核编码的规范
等等.
```

图 1.1 menuconfig

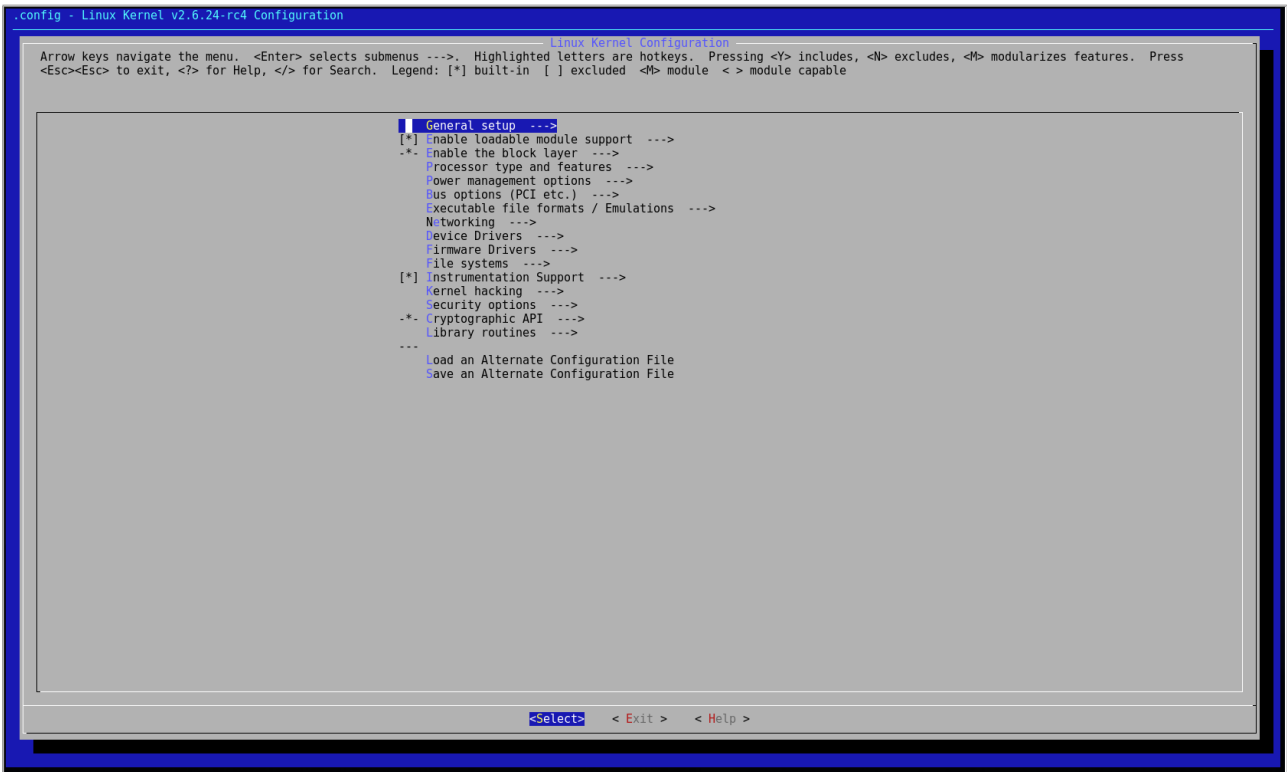


图 1.2 xconfig

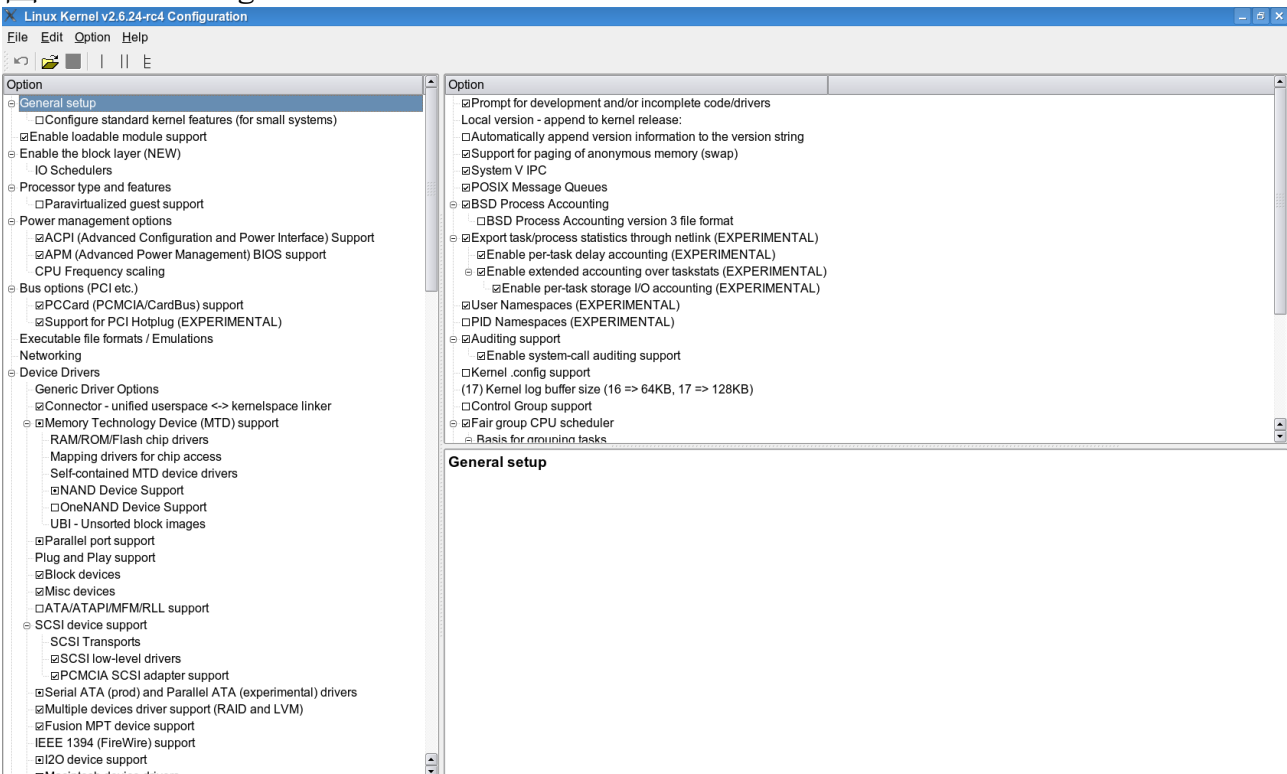
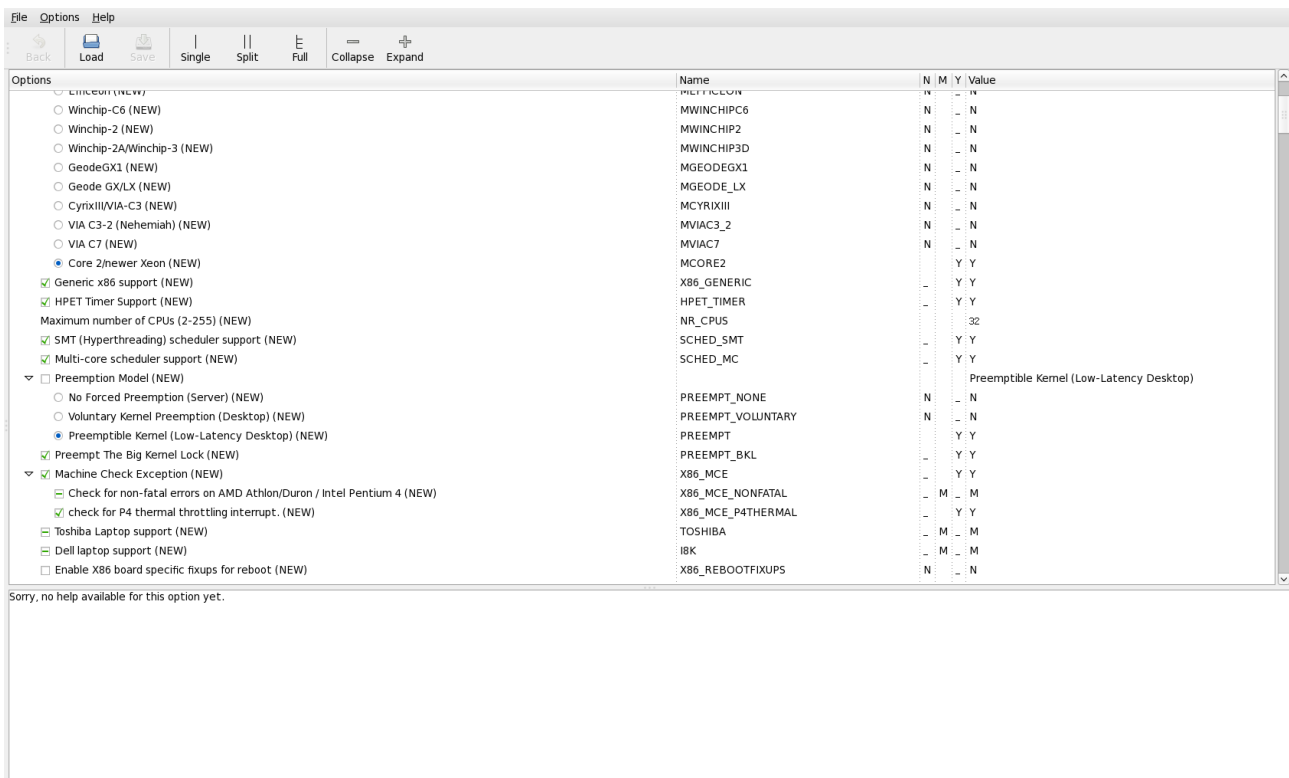


图 1.3 gconfig



## 2, 浏览代码

=====

### 2.1 vim + ctags + cscope

#### 2.1.1 编译 ctags 需要的 tags 文件:

```
# make ARCH=x86 tags
```

FYI: ARCH=<arch>的意思是, 除了指定的体系结构外, 不索引其他体系结构相关的代码.

FYI: 2.6.24 之前, 用 i386 指代 32 位的 x86 处理器, 用 x86\_64 指代 64 位的 x86 处理器。但在 2.6.24 中这 2 者即将合并, 统称为 x86.

#### 2.1.2 编译 cscope 程序需要的 cscope.out 文件:

```
# make ARCH=x86 cscope
```

### 2.1.3 在 Vim 中用 ctags 浏览

```
:ts do_fork //tag select 效果等同于光标停在 do_fork 上, 按 Ctrl-]。
//如果想返回原来的位置, Ctrl-T
:tn //tag next
:tp //tag previous
:tf //tag first
:t1 //tag last
```

FYI: 安装 Vim 的 taglist 插件, 浏览更方便。又, ~/.vimrc 中控制 number:

```
“added by albcamus
function! TlistWrapper2 ()
    Tlist
    if winnr('$') >= 2 " winnr() return the number of current window
        " while winnr('$') returns the number of all windows
        set nonu
    else
        set nu
    endif
endfunction

command! -nargs=0 List call TlistWrapper2()
map <F1> :List<CR>:<CR>
```

则可以按 F1 键来调用 taglist。

### 2.1.4 在 Vim 中用 cscope 浏览

```
:cs add cscope.out #添加一个 connection, 浏览内核时经常发现 cscope 链接断开
                    可以调用:cs a cscope.out
:cs show #显示所有 connections
:cs kill <#> #杀死第#号 connection, 第#号就是:cs show 显示出来的
:cs reset #reset 所有 connections
:cs find c|d|f|g|l|s|t <name>
```

解释:

```
c 查找 name 被哪些函数调用
```

d	查找 name 调用了哪些函数
e	egrep 句型
f	查找名为 name 的文件
g	查找 name 的定义
i	查找#include 本 name 的文件
s	查找 name 这个 C 符号
t	查找何处对 name 赋值

对 Linux 内核来说，如果既有 tags 文件，又有 cscope.out 文件，则 Vim 中:set cst 意味着使用 cscope 风格的^]，也就是:tag 命令。于是：

```

如果想用 ctags 找到 do_IRQ:
:ts do_IRQ
如果想用 cscope 找到 do_IRQ:
:ta do_IRQ

```

FYI: 可以在 ~/.vimrc 中自定义命令，这样就不必每次都输入符号名：

```

nmap <C-\>g :cs find g <C-R>=expand("<word>")<CR><CR>
nmap <C-\>c :cs find c <C-R>=expand("<word>")<CR><CR>
nmap <C-\>s :cs find s <C-R>=expand("<word>")<CR><CR>
nmap <C-\>t :cs find t <C-R>=expand("<word>")<CR><CR>
nmap <C-\>e :cs find e <C-R>=expand("<word>")<CR><CR>
nmap <C-\>f :cs find f <C-R>=expand("<file>")<CR><CR>
nmap <C-\>i :cs find i ^<C-R>=expand("<file>")<CR>$<CR>
nmap <C-\>d :cs find d <C-R>=expand("<word>")<CR><CR>

```

这样，当你在 Vim 中把光标挪到字符串 do\_fork 上，按 Ctrl-\c(先按 Ctrl-\然后按 c)，cscope 就会列出所有调用了 do\_fork 函数的地方。

```

clear local APIC
disable local APIC
lpic_shutdown
verify local APIC
sync Arb IDs
init bsp APIC
setup local APIC
detect_init APIC
init apic mappings
APIC_init_uniprocessor
parse_lpic
parse_nolpic
parse_disable_lpic_timer
parse_lpic_timer_c2_ok
apic_set_verbosity
smp_spurious_interrupt
smp_error_interrupt
apic_intr_init
connect_bsp APIC
disconnect_bsp APIC
lpic_suspend
lpic_resume
apic_pm_activate
init_lpic_sysfs
apic_pm_activate

genapic_32.h (/root/Sources
macro
ASM GENAPIC_H
APICFUNC
IPIFUNC
IPIFUNC
APIC_INIT
struct
genapic

if (!lapic_id_registered())
BUG();

/*
 * Intel recommends to set DFR, LDR and TPR before enabling
 * an APIC. See e.g. "AP-388 82489DX User's Manual" (Intel
 * document number 292116). So here it goes...
 */
init_apic_ldr();

/*
 * Set Task Priority to 'accept all'. We never change this
 * later on.
 */
value = apic_read(APIC_TASKPRI);
value |= ~APIC_TPRI_MASK;
apic_write_around(APIC_TASKPRI, value);

/*
 * After a crash, we no longer service the interrupts and a pending
 * interrupt from previous kernel might still have ISR bit set.
 */
/*
 * Most probably by now CPU has serviced that pending interrupt and
 * it might not have done the ack APIC_irq() because it thought,
 * interrupt came from i8259 as ExtInt. LAPIC did not get EOI so it
 * does not clear the ISR bit and cpu thinks it has already serviced
 * the interrupt. Hence a vector might get locked. It was noticed
 * for timer irq (vector 0x31). Issue an extra EOI to clear ISR.
 */
for (i = APIC_ISR_NR - 1; i >= 0; i--) {
    value = apic_read(APIC_ISR + i*0x10);
    for (j = 31; j >= 0; j--) {
        if (value & (1<<j))
            ack APIC_irq();
    }
}

Tag List 158,5 Bot ~/Sources/linux-2.6/arch/x86/kernel/apic_32.c 852,20 564
kscope tag: setup_local APIC
# line filename / context / line
1 1200 /root/Sources/linux-2.6/arch/x86/kernel/apic_32.c <<APIC_init_uniprocessor>>
    setup_local APIC();
2 1216 /root/Sources/linux-2.6/arch/x86/kernel/apic_64.c <<APIC_init_uniprocessor>>
    setup_local APIC();
3 281 /root/Sources/linux-2.6/arch/x86/kernel/smpboot_32.c <<smp_callin>>
    setup_local APIC();
4 1038 /root/Sources/linux-2.6/arch/x86/kernel/smpboot_32.c <<smp_boot_cpus>>
    setup_local APIC();
5 1049 /root/Sources/linux-2.6/arch/x86/kernel/smpboot_32.c <<smp_boot_cpus>>
    setup_local APIC();
6 214 /root/Sources/linux-2.6/arch/x86/kernel/smpboot_64.c <<smp_callin>>
    setup_local APIC();
7 886 /root/Sources/linux-2.6/arch/x86/kernel/smpboot_64.c <<smp_prepare_cpus>>
    setup_local APIC();
8 48 /root/Sources/linux-2.6/arch/x86/mach-visws/traps.c <<cobalt_init>>
    setup_local APIC();
Choice number (<Enter> cancels): █

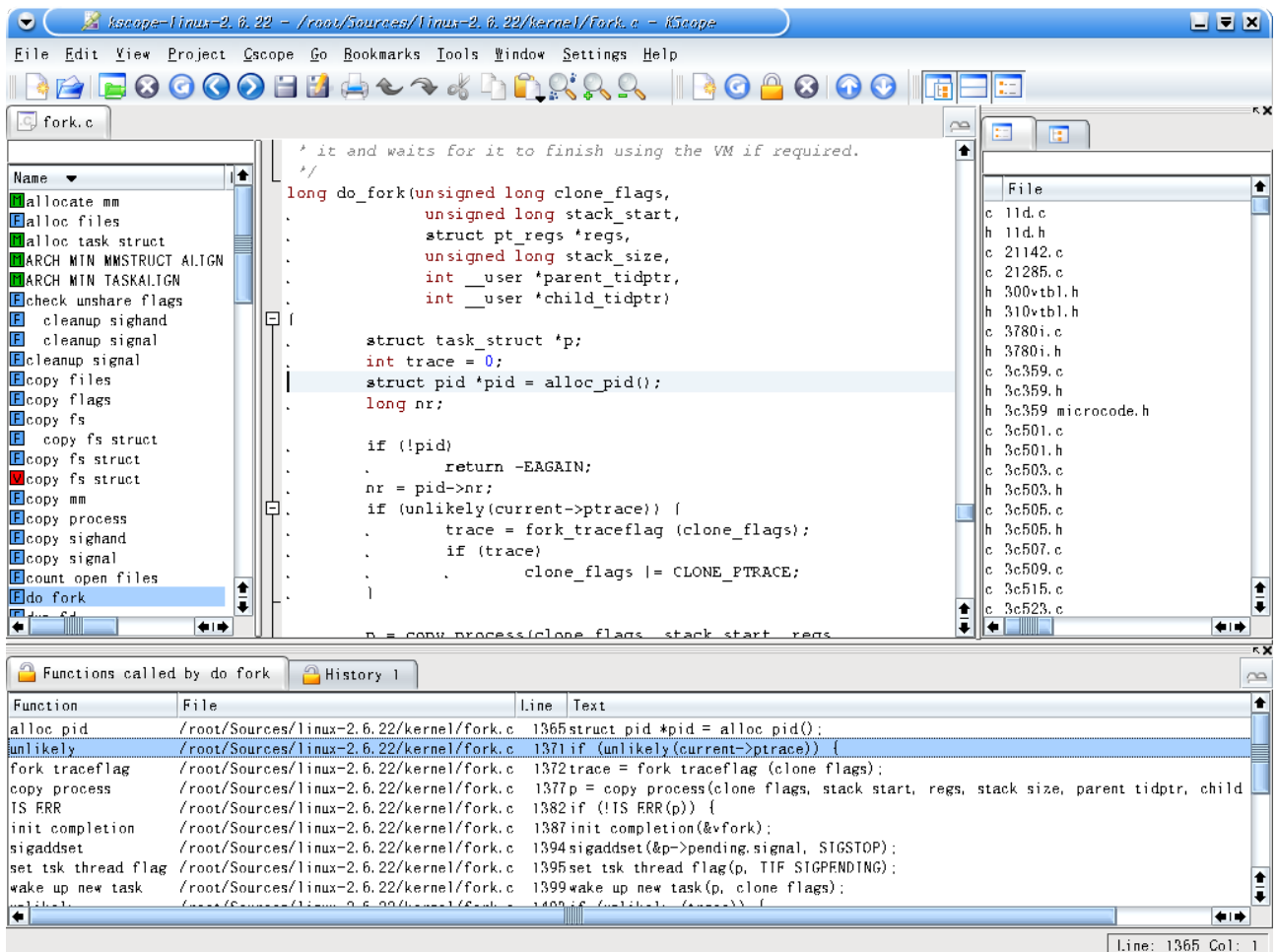
```

FYI: GNU Global 程序也不错。对喜爱 Emacs 编辑器的朋友，还可以用 etags 和 xscope 来替代。喜欢 Source Insight 的朋友也可以常识一下 kscope，它的界面和 SI 很类似。

kscope 依赖于 KDE 的库，如果你用 gnome 作窗口管理器，也得安装了 KDE 才能安装 kscope。使用 kscope 之前需要安装 ctags、cscope 和 dot 程序(dot 在 graphviz 包中)。

kscope 抓图:





## 2.2 1xr - Linux Cross Reference

我一般都只在网上浏览:


<http://1xr.linux.no/>

<http://users.sosdg.org/~qiyong/1xr/source/>

文件 (E) 编辑 (E) 查看 (V) 历史 (S) 书签 (B) 工具 (I) 帮助 (H)

新建标签页 后退 前进 重新载入 停止 主页 http://kr.linux.no/ident?a=x86\_64&i=do\_fork Google DownThemAll! 下载

Gmail - 收件... 利用异常表处... PBUeng > Beiji... PBUeng > Chin... AP-388 82489... http://....patch Linux identf...



## Cross-Referencing Linux

[\[ source navigation \]](#)  
[\[ identifier search \]](#)  
[\[ freetext search \]](#)  
[\[ file search \]](#)

Version: [\[ 1.0.9 \]](#) [\[ 1.2.13 \]](#) [\[ 2.0.40 \]](#) [\[ 2.2.26 \]](#) [\[ 2.4.18 \]](#) [\[ 2.4.20 \]](#) [\[ 2.4.28 \]](#) [\[ 2.6.10 \]](#) [\[ 2.6.11 \]](#) [\[ 2.6.17.13 \]](#) [\[ 2.6.18 \]](#) [\[ 2.6.20.1 \]](#) [\[ 2.6.22 \]](#) [\[ 2.6.22.6 \]](#)

Architecture: [\[ i386 \]](#) [\[ alpha \]](#) [\[ arm \]](#) [\[ ia64 \]](#) [\[ m68k \]](#) [\[ mips \]](#) [\[ mips64 \]](#) [\[ ppc \]](#) [\[ s390 \]](#) [\[ sh \]](#) [\[ sparc \]](#) [\[ sparc64 \]](#) [\[ x86\\_64 \]](#)

---

Identifier:

### do\_fork

Defined as a function in:

- [kernel/fork.c, line 1356](#)

Defined as a function prototype in:

- [include/linux/sched.h, line 1436](#)

Referenced (in 36 files total) in:

- arch/alpha/kernel/process.c:
  - [line 252](#)
  - [line 258](#)
- arch/arm/kernel/process.c, line 426
- arch/arm/kernel/sys\_arm.c:
  - [line 238](#)
  - [line 255](#)
  - [line 260](#)
- arch/arm26/kernel/process.c, line 368
- arch/arm26/kernel/sys\_arm.c:
  - [line 242](#)
  - [line 259](#)
  - [line 264](#)
- arch/avr32/kernel/process.c:
  - [line 96](#)
  - [line 353](#)
  - [line 362](#)
  - [line 369](#)
- arch/blackfin/kernel/process.c:
  - [line 216](#)
  - [line 226](#)
  - [line 242](#)
- arch/cris/arch-v10/kernel/process.c:
  - [line 106](#)

完成

### 3. 内核相关的图书

=====

## 入门推荐： LKD2 - 《Linux 内核设计与实现-第 2 版》

从入门开始，介绍了诸如中断、系统调用、虚拟文件系统、同步与互斥、内存管理、进程控制等方面，内容比较浅显易懂，是入门的好书。

优点： 适合入门（个人认为，没有比 LKD2 更优秀的内核入门图书）

缺点： 内容不够深入，覆盖面不广。（对高手来说估计就像休闲材料）

### 3.1 ULK3 - 深入理解 Linux 内核-3rd

一本很全面的 Linux 内核原理书。书中以 2.6.11 为示例版本，着重讲述了 x86 平台的 Linux 内核实现。

优点： 深入，全面

缺点： N/A

我觉得看完 ULK3 就是高手了:) 而且由于书中着重介绍了 X86 体系结构，也称得上半个 x86 专家了。

### 3.2 《Linux 内核源代码情景分析》

以 2.4.0 为例讲解，注重代码级别的剖析，对中断、内存管理、文件系统、SMP、PCI 和 USB、IPC 的讲解都是代码级别的深入细致。

优点： 深入

缺点： 针对的内核版本较旧，且每个「情景」都很长，不容易坚持读完。

FYI : 新手不要从《情景分析》开始学习内核，这样只会增长你的学习周期。

### 3.3 LDD3 - Linux 设备驱动程序-3rd

LDD3 写的很精彩。但如果要学会写具体的驱动程序，还是得参照硬件的 datasheet，读一个内核中现成的驱动程序。

FYI : 内核中自带的驱动程序 skeleton:

drivers/net/pci-skeleton.c 和 drivers/usb/usb-skeleton.c, 分别是为 PCI/USB 驱动程序提供的参考代码。

### 3.4 现代体系结构上的 Unix 系统 - 内核程序员的 SMP 和 Caching 技术

这本书着重讲解各种体系结构上的 Unix 实现注意事项, e.g. SMP 的同步与互斥、Cache 一致性问题。

优点: 作者知识面非常广, 原理讲得很清楚。

缺点: 94 年的书, 比较旧

### 3.5 Intel & AMD CPU 参考手册

最好带着问题有针对性的去读 Intel & AMD 的手册。

## 4, 第一个模块: Hello Kernel

=====

### 4.1 代码

-----

```
# cat Makefile
    ifneq ($(KERNELRELEASE),)
        obj-m := hello.o
    else
```

```

KBUILD := /lib/modules/`uname -r`/build
modules:
    make -C $(KBUILD) M=$(shell pwd) modules
clean:
    rm -f *.o *~ core .depend *.cmd *.ko *.mod.c
    rm -rf .tmp_versions
endif

```

```
# cat hello.c
```

```

#include <linux/kernel.h>
#include <linux/module.h>
#include <linux/init.h>

MODULE_LICENSE("GPL");
MODULE_DESCRIPTION("The hello world demo module");
MODULE_AUTHOR("albcamus<albcamus@gmail.com>");

static int __init hello_init(void)
{
    printk("Hello kernel!\n");

    return 0;
}

static void __exit hello_exit(void)
{
    printk("Bye!\n");
}

module_init(hello_init);
module_exit(hello_exit);

```

## 4.2 编译、加载和查看

---

编译:

```
# make
```

加载

```
# insmod ./hello.ko
```

查看:

```
# dmesg |tail -1
```

```
Hello Kernel!
```

卸载:

```
# rmmod hello
```

FYI : 接下来做什么? 在 Hello world 之后, 有很多可以练习的小模块可以做, 例如截获系统调用表(这个最好只拿来玩), 写一个 dummy 字符设备驱动

程序，添加一个系统调用，等等。

(TBD: 如果有时间, 就讲 *intercept-execve* 模块)

## 5, 内核相关的网络资源

---

<http://www.kernel.org>

- git trees
- bugzilla
- LKML (和其他的 mailing list)

FYI: 两个新闻组 `comp.os.linux.development.system` (讨论 linux 内核程序设计), 和 `fa.linux.kernel`(LKML 的 USENET 版)  
新闻组可以用 `evolution/thunderbird` 等订阅, 也可以在 `http://groups.google.com` 访问。

<http://kerneltrap.org>

Linux 和 BSD 内核的技术新闻。 如果没时间跟踪 LKML, 那么经常浏览 `kerneltrap` 是个好主意。

<http://lwn.net>

Linux weekly news

<http://www.linuxforum.net>

从1999年开始，积累了大量的高质量内核技术文档资源。

<http://linux.chinaunix.net>

积累壮大中。有很多高手活跃，如zx\_wing、puppylove、daameon等。

[http://www.ibm.com/developerworks/cn/views/linux/articles.jsp?view\\_by=search&search\\_by=%E5%86%85%E6%A0%B8](http://www.ibm.com/developerworks/cn/views/linux/articles.jsp?view_by=search&search_by=%E5%86%85%E6%A0%B8)

ibm developers works 网站有很多关于内核的文档，大都写得通俗易懂。

<http://zh-kernel.org/mailman/listinfo/linux-kernel>

由Li Yang @ freescale 维护，活跃着很多在Kernel 开发领域活跃的华人，如Herbert Xu, Mingming Cao, Bryan Wu, Shaohua Li, Yanmin Zhang 等都有参与，此外还有各个社区的高手如CLF的wheelz、daameon、zh11g等。

该社区有翻译的内核文档，还有Intel OTC的Linux ACPI项目。

## 6，其他的知识

=====

在学习Linux内核时，我们经常发现很多内核之外的知识要掌握，否则就会成为

学习之路上的绊脚石。大体上来说，各个领域都有它的特定知识需要掌握，例如要学习 Linux 内核的网络协议栈实现，就不可能不学习 TCP/IP；要学习 Linux 是如何管理物理内存，就不可能不学习某个 CPU 的内存管理机制。等等。

这里只说一点：GCC 的 C 扩展 和 inline ASM。

GCC 的 C 扩展：

```
$ info gcc 'c e'
```

INLINE ASM(针对 i386)：

《情景分析》第一章讲解 inline ASM 比较透彻。另外有两篇很好的文档：

<http://www.ibiblio.org/gferg/ldp/GCC-Inline-Assembly-HOWTO.html>

[http://groups.google.com/group/muc.lists.linux-kernel/browse\\_thread/thread/ca8860dd41ddd42b](http://groups.google.com/group/muc.lists.linux-kernel/browse_thread/thread/ca8860dd41ddd42b)

针对 x86-64 的 inline 汇编和 i386 区别不大。

## 7, 内核调试

=====

### 7.1 *printk*

朴素的 `printk` 永远是最受欢迎的。

缺点：有时 BUG 导致直接 reboot 了，`printk` 的内容甚至都看不到。不过没关系，我们可以用 `console=/dev/ttyS0` 把 `console` 的输出定向到串口，在另一端接收；也可以用 `netconsole`(详见 `Documentation/networking/netconsole.txt`)



优点： 随时随地可用。

FYI： 使用 GCC 或 C99 的可变参数宏包裹 kprintk

C99 方式:

```
=====  
#define MYDEBUG  
#ifdef MYDEBUG  
#define dbg_print(fmt, ...) \  
printk("%s-%d: "fmt, __FILE__, __LINE__,  
##__VA_ARGS__);  
#else  
#define dbg_print(fmt, ...) do {} while (0)  
#endif
```

GCC 方式:

```
=====  
#define MYDEBUG  
#ifdef MYDEBUG  
#define dbg_print(fmt, args...) \  
printk("%s-%d: "fmt, __FILE__, __LINE__,  
##args);  
#else  
#define dbg_print(fmt, args...) do {} while(0)  
#endif
```

然后就可以使用:

```
dbg_print("This is a debug messge, i = %d, and string is %d\n", i, str);
```

## 7.2 kdb

由 SGI 公司开发的一系列补丁。

<http://oss.sgi.com/projects/kdb/>

优点： 更新很快，基本能和官方内核同步。 稳定。

缺点： 只能做到汇编级别的源代码单步，不能做到 C 代码单步。

## 7.3 kgdb

linsys 公司开发的一系列补丁。 Linux 内核维护者 Andrew Morton 就使用它。

优点： C 源代码级别的单步。

缺点： BUG 较多，需要串口，两台机器。 免费版本支持的内核版本很少。

FYI： 这个问题可以规避:

a), VMWare 中可以指定虚拟机的串口为 FIFO 文件，从而两台虚

虚拟机实现"串口"互联:

b), 较新的 kgdb 已经支持 kgdboe - Over Ethernet. 只要以太网卡驱动程序支持 NETPOLL 即可。

FYI : kgdb 即将被合并到官方内核中, 不过最早也得是 2.6.25.

FYI : 比 linsys 和 sourceforge 要全的一个地方:

<http://www.kernel.org/pub/linux/kernel/people/agk/patches/2.6/>

另 Andrew Morton 的-mm tree 中自带 kgdb:

<http://www.eu.kernel.org/pub/linux/kernel/people/akpm/patches/2.6/>

## 7.4 kprobes/systemtap

kprobes 自 2.6.9 加入到官方内核中。

systemtap 是由 redhat 和 Intel 等公司开发的调试工具, 类似于 Solaris 的 DTrace。它基于 kprobes。

kprobes

优点: 几乎内核中的任何地方都可以插入监测点, 而且可以使用它来截获那些不导出的符号(只要在 kallsyms 中有)

缺点: N/A

systemtap:

下载:

```
git clone git://sources.redhat.com/git/systemtap.git
```

优点: 用户态, 功能强大

缺点: 需要学习一门新的语言(stap)

文档:

主页:

<http://sourceware.org/systemtap/documentation.html>

使用 systemtap 调试内核:

<http://www.ibm.com/developerworks/cn/linux/1-systemtap/>

Linux 下的一个全新的性能测量和调试诊断工具 systemtap, 第3部分:

<http://www.ibm.com/developerworks/cn/linux/1-cn-systemtap3/>

## 7.5 VMWare workstation 6.x

VMWare workstation 在 VCPU 中实现了 gdb stub

优点: 简单(不需要串口、不需要多台机器), C 源代码单步。

缺点: 并不针对 Linux 系统, 不能针对线程栈做分析。学习意义大于工程意义。

文档:

<http://blog.chinaunix.net/u/548/showart.php?id=282575>

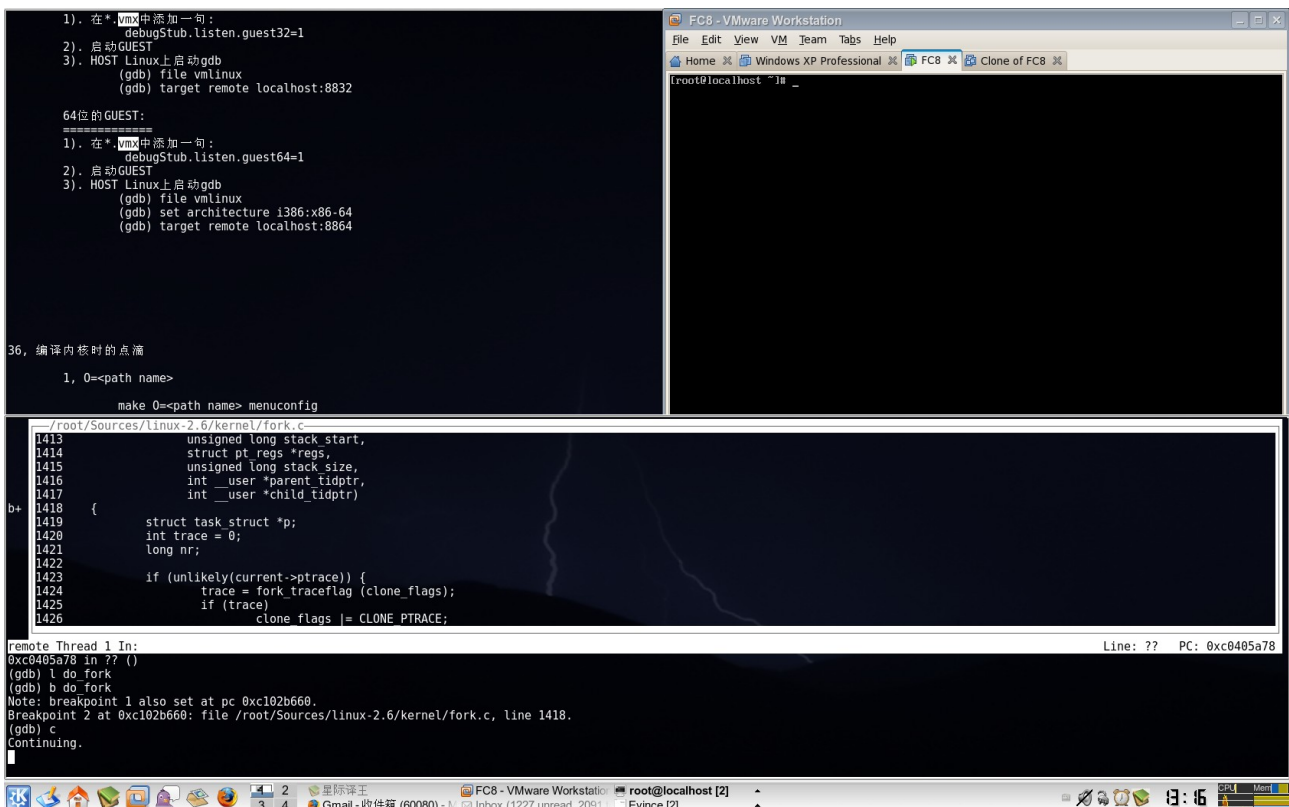


图: WS6.0

# 8, 现代硬件: ARCH & DRIVER

---

## 8.1 现代 X86 CPU

Intel 64 and IA 32 CPU Manuals:

<http://www.intel.com/products/processor/manuals/index.htm>

AMD CPU Manuals:

<http://developer.amd.com/devguides.jsp>

一些 Linux 支持(或需要更好地支持)的现代 X86 CPU 的特性:

多处理器: SMP, 同步和互斥, APIC

内存管理: PAE, PSE, MTRR/PAT, Cache/TLB, x86-64, NUMA

系统调用: `int 0x80`, `sysenter/sysexit`, `syscall/sysret`

电源管理: ACPI

-> 电源管理

-> MPS

-> PCI IRQ Routing

-> PNP

-> E820

扩展指令: MMX, SSE/SSE2/SSE3/SSE4, 3DNow!

安 全 : NX

虚拟机 : VMX(Intel VT), SVM(AMD-V)

## 8.2 外围设备

外围设备资料，建议下载 Intel ICH8 南桥芯片的 datasheet:

<http://www.intel.com/design/chipsets/datashts/313056.htm>

南桥芯片的 datasheet 会囊括 host controllers 手册。当然，具体到某种设备，例如 USB 或 SATA，还需要看相关的规范。

# 9, 80386 的分段机制

=====

## 1. 80386 为段映射提供的硬件

- o GDT 表, LDT 表, TSS, 都在内存中。(Linux 的实现是: 这 3 个东西都是每个 CPU 有一个)
- o 每个 CPU 中都有一个 gdtr 指向 GDT 表, 一个 ldtr 指向 LDT 表, 一个 tr 指向 TSS。
- o 每个 CPU 有 6 个段寄存器: CS, DS, SS, ES, FS, GS。它们都是 16-bit 的。  
其中 CS、DS 分别指向代码和数据段, SS 指向堆栈段, 都是特殊的; 而 ES、FS、GS 可以用作普通目的, 指向哪个 data segment 都成。

这些段寄存器, 叫作"Segment Selector"。亦即: 由它的值决定一个地址(16 位的 segment selector 加上 32 位的线性地址)采用 GDT 或 LDT 中的哪个 entry — 每个 entry 都代表一个"Segment", 所以一个 entry 就是一个"Segment Descriptor"。

- o 一个 segment selector 是这样的:

```
15 _____ 3 | 2 | 1 | 0 |  
| _____ INDEX _____ | TI | RPL |
```

- RPL: Request Privilege Level。  
CS 寄存器的[1:0]就叫 CPL, current privilege level。
- TI: Table Indicator。TI==0 时, 表示使用 GDT 表; TI==1 时, 使用 LDT 表。
- INDEX: 在 GDT 或 LDT 表中的位置。

例如, TI==0,使用 GDT 表。INDEX 为 2, gdt\_r==0x00020000, 而 GDT 表每个 entry 的大小为 8 字节(见下条), 于是会找到:

$$0x00020000 + (2 * 8) == 0x00020010$$

这个地址就是 GDT 表中相应的 segment descriptor 在内存中的地址。

- o GDT 表的一个 entry, 代表一个 segment descriptor, 长度为 8 字节(64-bit), 格式为:

(参考 ULK3 P39, Intel 手册 V3 Figure 4-1 和 Linux 内核的 struct desc\_struct 定义)

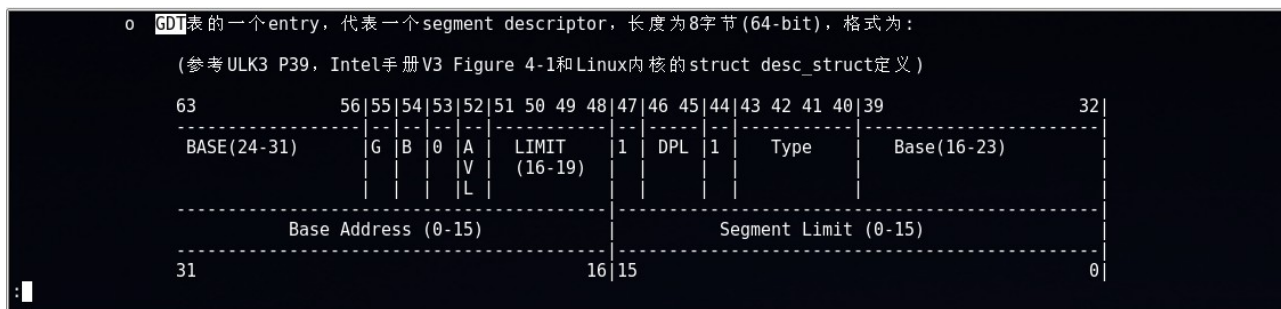


Table 48-1

A	Accessed
AVL	Available to System programmer
B	Big
C	Conforming
D	Default
DPL	Descriptor Privilege Level
E	Expansion Direction
G	Granularity
R	Readable
LIMIT	Segment Limit
W	Writable
P	Present

Linux 内核对这个 entry 格式的定义:

```
struct desc_struct {
    ul6 limit0;
    ul6 base0;
    unsigned base1 : 8, type : 4, s : 1, dp1 : 2, p : 1;
```

```

        unsigned limit : 4, avl : 1, l : 1, d : 1, g : 1, base2 : 8;
    } __attribute__((packed));

```

注意对 Code 段, 来说, limit 由 Limit 和 G 决定; 对 Data 段和 Stack 段来说, Limit 还依赖于 B 和 E 位。但 Linux 根本不管 B 和 E 位。

其中代码段描述符、数据段描述符都可以放在 GDT 或 LDT 中作为 entries, 但 TSS 段描述符只能放在 GDT 表中作一个 entry, LDT 描述符只能放在 GDT 表中作一个 entry。

- o 对应 cs/ds/ss/es/fs/gs 这 6 个 segment selectors, 可以当前指定 6 个 segment descriptors(在 gdt 表或 ldt 表中)。为此, 80x86 CPU 提供了 6 个软件不可见的 64-bit 寄存器, 存储这 6 个当前的 segment descriptors。这样, 例如 cs 和 ds 已确定, 那么 CPU 就把相应的 descriptor 加载到不可见寄存器中, 从此不再访问内存中的 GDT/LDT 表读 descriptors, 直到相应的 cs/ds 等 selector 被改变。

## 2. Linux 对段映射的实现。

ULK3 说, 2.6 Linux 只有在 x86 处理器需要时, 才在段映射一事上敷衍敷衍它。

- o 作为 GDT 表中的一个 entry, \_\_KERNEL\_CS 的定义:

```

#define GDT_ENTRY_KERNEL_CS (GDT_ENTRY_KERNEL_BASE + 0)
#define __KERNEL_CS (GDT_ENTRY_KERNEL_CS * 8)

```

解释:

Linux 把内核代码段作为 GDT 表的第 12 项, 所以其 INDEX 就是 12。INDEX 在 CS 的高 13 位, 低 3 位的 TI==0, 因为用的是 GDT; RPL==0, 因为是 RING0 权限。这样, \_\_KERNEL\_CS 就是  $12 \ll 3$ , 亦即  $12 * 8$ 。

- o 而 \_\_USER\_CS:

```

#define GDT_ENTRY_DEFAULT_USER_CS      14
#define __USER_CS (GDT_ENTRY_DEFAULT_USER_CS * 8 + 3)

```

解释:

Linux 把用户代码段作为 GDT 表的第 14 项, 所以其 INDEX 是 14。INDEX 在 CS 的高 13 位, 所以左移 3 位, 亦即乘上 8。至于低 3 位, 因为用的是 GDT, 所以 TI==0; 又因为用户态, 所以 RPL==3。所以 \_\_USER\_CS 就是  $14 \ll 3 + 3$ , 亦即  $14 * 8 + 3$ 。

- o \_\_KERNEL\_DS 是 GDT 表中第 13 项, \_\_USER\_DS 是 GDT 表中第 15 项。

好吧，我们看看 Linux 定义的这 4 个 segments，其 base address 和 limit 的设置。

```
GDT[12] == kernel cs
GDT[13] == kernel ds
GDT[14] == user cs
GDT[15] == user ds
```

在 GDT 表中这 4 个 entries 的 base address 都是 0x00000000，其 limit 都是 0xffff。每个 entry 的 G 标志被设为 1。（注意，如果一个 entry 的 G 设为 0，表示 limit 的单位是 byte，那么 limit == 0xffff 就指定了 1MB 的空间；如果 G 被设置为 1，则 limit 的单位是 4096 字节，于是 limit==0xffff 就表示  $(2^{20}) \times (2^{12}) = 4GB$  的空间）

Linux 就是这么敷衍 80386 的 segmenting unit 的：4 个段，每个都是 flat 的 4GB 大小。分段，什么都不做。逻辑地址.offset == 线性地址。（注意，intel 的“逻辑地址”定义是 48bit 的，其中 16bit 是 segment selector）

arch/x86/kernel/cpu/common.c :

```
DEFINE_PER_CPU(struct gdt_page, gdt_page) = { .gdt = {
    [GDT_ENTRY_KERNEL_CS] = { 0x0000ffff, 0x00cf9a00 },
    [GDT_ENTRY_KERNEL_DS] = { 0x0000ffff, 0x00cf9200 },
    [GDT_ENTRY_DEFAULT_USER_CS] = { 0x0000ffff, 0x00cffa00 },
    [GDT_ENTRY_DEFAULT_USER_DS] = { 0x0000ffff, 0x00cff200 },

    {其他的忽略}
} };
EXPORT_PER_CPU_SYMBOL_GPL(gdt_page);
```

我们来看看，\_\_KERNEL\_CS 的 GDT entry 赋值为 { 0x0000ffff, 0x00cf9a00 } 是什么意思。

\_\_KERNEL\_CS 在 GDT 表中的 entry 值:

对照着 Table 48-1，知道 base address 为 0x00000000，limit 为 0xffff，Granularity 为 1 - 亦即 4KBytes 为单位。

0000 0000 1100 1111 1001 1010 0000 0000
0000 0000 0000 0000 1111 1111 1111 1111

- o UP Linux 中只有一个 GDT 表，但 SMP 中，每个 CPU 都有自己的 GDT 表。  
(include/asm-x86/desc\_32.h)



```

struct gdt_page
{
    struct desc_struct gdt[GDT_ENTRIES];
} __attribute__((aligned(PAGE_SIZE)));

DECLARE_PER_CPU(struct gdt_page, gdt_page);

```

每个CPU都有一个GDT表，变量名就叫做gdt\_page，是个struct gdt\_page结构，其中域gdt[GDT\_ENTRIES]就是GDT表，每个entry都是一个struct desc\_struct结构体变量。

- o Linux下基本不用LDT。  
所以只定义了一个默认的LDT，放在default\_ldt数组中。它包含5项，但只有两个常用：一个兼容iBCS操作系统的可执行文件，一个兼容Solaris/x86的可执行文件。

Wine需要修改ldt，于是Linux提供了modify\_ldt(2)系统调用。

## 10. IA32 Linux 的 2-level 页面映射

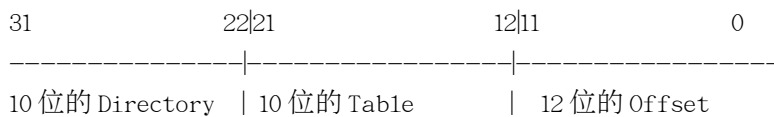
=====

参考：INTEL 手册 V3: Ch3, Section 3.7: Page Translation using 32-bit Physical Addressing

### *PART I: 线性地址的划分*

-----

2-level: 32 位的线性地址被划分为 3 部分:



- 高 10 位       : 又叫 Page Directory, PD
- 中间 10 位   : Page Directory, PT
- 低 12 位      : 页内 offset

每个进程都有自己的页表，用以找到其 PD 的物理地址，称为 PGD: Page Global Directory。

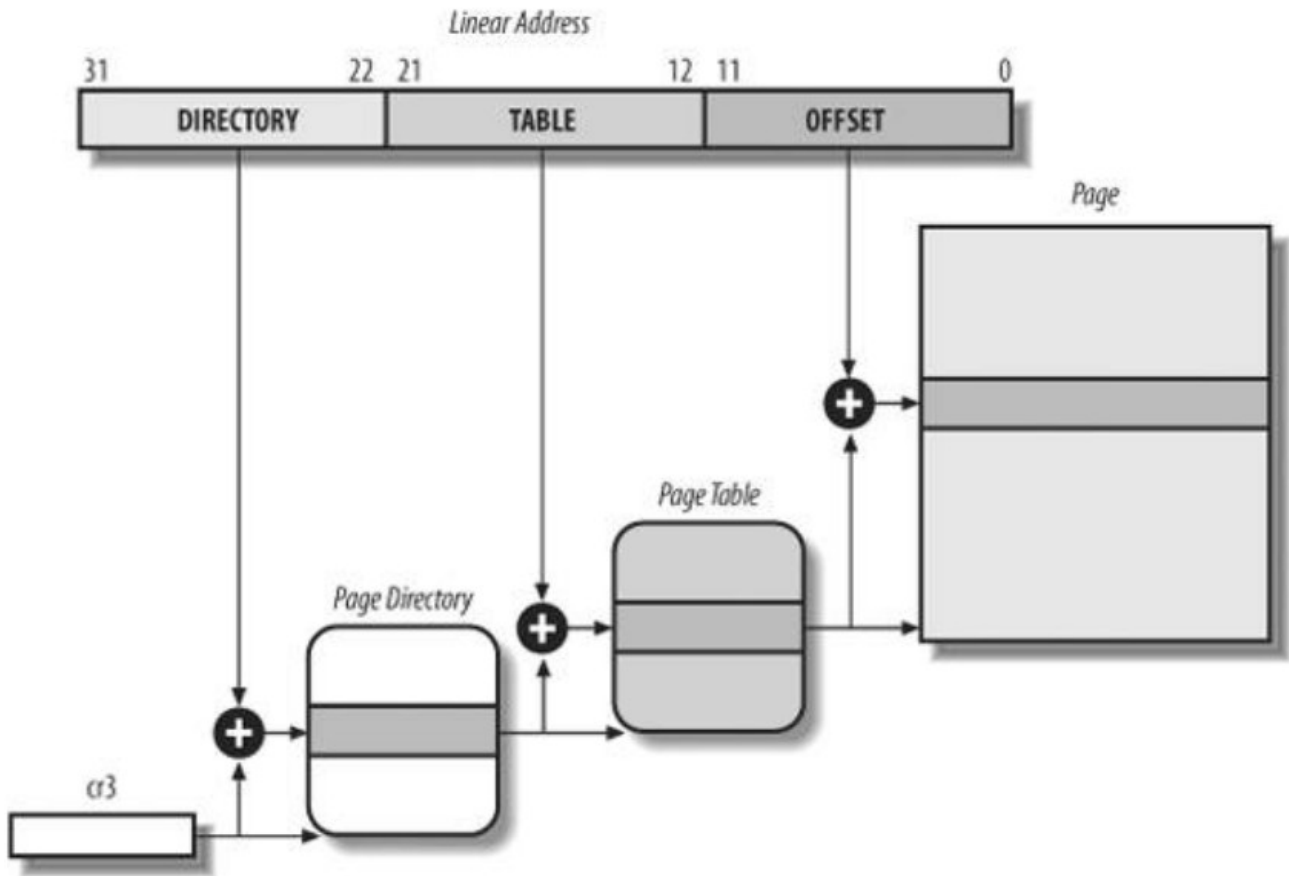
每个进程都有自己的页表，其 PGD 的物理地址放在 CR3 寄存器中，见 `schedule > context_switch > switch_mm > load_cr3(next->pgd)`。

## *PART II: MMU 分页单元寻找物理地址的过程*

---

- 1) MMU 根据 CR3 中的 `pgd` 值，和要访问的那个线性地址的高 10bit(就是在 PD 中的 `offset`)，在 PD 中找到相应的 `entry` —— 一个 PD 的 `entry`，是一个 32 位的物理地址。其中高 20bit 是 Page Table 的 Base Address。
- 2) MMU 根据 Page Table 的 Base Address，和要访问的那个线性地址的中间 10bit(`[12:21]`)，在 PT 中找到相应的 `entry` —— 一个 PT 的 `entry`(就是传说中的 PTE)，是一个 32 位的物理地址。其中高 20bit 是 Page 的 Base Address -- 找到 Page 的 base address 就找到了页面的物理地址。
- 3) MMU 根据 Page 的物理地址(which 低 12 位都是零，因为页首地址是 4096 对齐的)，和要访问的那个线性地址的低 12bit，找到该地址在该页面中的位置 — 这就是传说中的「物理地址」，和那个要访问的线性地址对应的。

线性地址=>物理地址(来自 ULK3):



### PART III: 标志

---

- \* PD和PT中的entries, 其格式是一样的。只不过:
- \* PD的entry约束的是一个包含PT的页面, 而PT的entry约束的是一个包含有效数据的页面。

BIT	Flag	说明
0	P, Present	<p>为1, 则Page在RAM中; 为0, 则Page不在RAM中。</p> <p>如果CPU试图访问某线性地址, 结果在中途发现PDE或PTE的P标志为0, 则会把该线性地址写入CR2寄存器, 然后触发Page Fault。CPU并不改写PDE或PTE的flags, 而是让OS来做这些:</p> <ol style="list-style-type: none"> <li>1. 把页面从交换设备上swap到RAM中</li> <li>2. 为该物理页面建立起映射, 也就是修改相应的PDE和PTE, 并且将P设置为1。其他的flags如Dirty和Accessed, 页可以在此时设置。</li> <li>3. Invalidate这个PTE在TLB Cache中的条目</li> <li>4. 从Page-Fault handler中返回</li> </ol> <p>我觉得, CPU遭遇PDE.P=0和PTE.P=0的情况还不一样。后者说明页面不存在, 而前者说明连页表都不存在。</p>
1	R/W, Read/Write	<p>为0, 则只读; 为1, 则可读写。</p> <p>指定页面, 或一组页面的读写属性。(啊, 为什么会是「一组页面」? Intel手册说, PTE约束一个页面, 而PDE约束了一个Page Table, 从而可以约束更多页面 — Bull shit!)</p> <p>当CRO.WP = 0 (CPU复位时默认就这样), 且CPU处于Supervisor模式, 则任何页面都是可以读/写的。当CPU处于User模式, 则Supervisor模式的页面既不可读也不可写, 即使是User模式的页面, 也要受这里的读写属性来约束。</p> <p>从PPro开始, CPU支持copy-on-write, 亦即:</p> <ol style="list-style-type: none"> <li>a) CRO.WP = 1时, 如果某User模式的页面只读, 那么Supervisor模式的CPU去尝试写它, 会发生page fault。</li> <li>b) CRO.WP = 1时, 如果某Supervisor模式的页面只读, 那么不管CPU处于什么模式, 都不能写它</li> </ol>
2	U/S, User/Supervisor	<p>为0, 表示页面是Supervisor态的; 为1, 表示页面是User态的。</p>
3	PWT, Write-Through	<p>为0, 则write-back; 为1, 则write-through。</p> <p>如果CRO.CD = 1, 则CPU会忽略该bit。</p>
4	PCD, Cache Disabled	<p>为1, 则该页面不被Cache; 为0, 则页面被Cache</p> <p>当为MMIO的页面建立映射时, 应该置该bit以禁止CPU Cache它。</p>
5	A, Accessed	<p>为1, 表示页面被访问过了; 为0, 表示没被访问过。</p> <p>一般来说, 一个页面刚被load进内存时, 该位置0。什么时间它被访问了, CPU就会把该位置1。</p>
6	D, Dirty	<p>***Note, PDE中这一bit不起作用。只有PTE中的该bit起作用***</p> <p>为1, 表示页面Dirty; 为0, 表示页面Clean。</p>
7 @ PDE (4K pages) PS, Page Size		<p>为0表示4K大小; 为1表示4M大小。(如果为1, 且PAE被激活, 则页面大小为2M)</p>
7 @ PTE (4K Pages) PAT, Page Attribute Table 和12 @ PDE (4M Pages) Index		<p>PAT在Pentium III中引入。与PCD和PWT标志一起使用, 来在PAT表中索引到某个entry。</p> <p>对于不支持PAT的CPU, 该位应当永远置0。</p> <p>FIXME PAT有空研究一下</p>
8	G, Global Page	<p>该位在PPro CPU中引入。为1, 则这是一个Global Page; 为0则不是GP。</p> <p>注意, 当CR4.PGE=1时才有效。</p> <p>Global Page的意思是: 它是全局的, 因此当进程切换时(或CR3被写时), 不需要invalidate该页的映射在TLB中的entry。</p> <p>FIXME: 我觉得Linux的内核态内存就是这样。各个进程的CR3所指向的页表都有一部分是指向了内核内存的页表, 但是所有进程都使用想通的内核内存页表, 因此不需要在进程切换时改变什么。</p>
[9:11]	Avail, Available for system programmer's use	N/A

Thank you!

Any Questions?