# Linux 2.6 内核启动分析

华清远见深圳中心 学员 李枝果 Li zhi quo0532@163. com 2010-6-04

-----

#### 特别说明:

此文章的原始出处 http://www.docin.com/sz\_farsight

此文章是 **华清远见 深圳培训中心** 学员在学习期间的兴趣专题文章,文章的部分内容可能参考了网络资源而没有标明出处,敬请见谅。此文档欢迎转载,但转载请<mark>保留学员个人信息</mark>及**此段说明**,谢谢!

·-----

## 特别说明之后的一般说明 ^\_^:

本文的分析是基于之前分析 SD 卡驱动的内核。

本文由 深圳 华清远见 学员李枝果总结,易松华老师批注

# 一、Makefile 是如何生成 ulmage 的

在研究内核之前,首先我们要认识内核镜像的结构。顶层的 Makefile 生成 ul mage 是按照下面的几个步骤生成的(这里省略了各子目录下\*.o 文件的生成)

1. arm-linux-gnu-ld 用 arch/arm/kernel/vmlinux.lds 、 arch/arm/kernel/head.o、arch/arm/kernel/init\_task.o、各子目录下的 built-in.o、lib/lib.a 、arch/arm/lib/lib.a 生成项层目录下的 vmlinux (根据 arch/arm/kernel/vmlinux、lds 来链接 0xc0008000)

vml i nux. I ds 开始处

批注 [s1]: ELF 格式文件,这 个是原汁原味,不添加任何防 腐剂的真正 kernel,是 linux 杀手级的东东^\_^

**批注** [s2]: 让初学者伤透脑筋的东东,如知详情,先插播广告:来"深圳华清远见"学习,先,哈



```
stem.map vmlinux.lds vmlinux.lds piggy.S
😅 🖀 🗐 🤌 風 🖹 🖺 🔡 🖺 🔞 🐧 🈘 😘 🖺 stack
                                          ▼ | E E E E | AA % at to E
  0 10 10 20 30 40 50 651 OUTPUT ARCH (arm)
  652 ENTRY (stext)
  653
  654 jiffies = jiffies_64;
  655
  656 SECTIONS
  657 {
  658
      .init : \{ /* \text{ Init code and data} \}
  660
          sinittext =
  661
         *(.init.text)
  662
         _einittext = .;
  663
          proc info begin = .;
  664
  665
         *(.proc.info.init)
        __proc_info_end = .;
  667
          arch_info_begin = .;
        *(.arch.info.init)
  668
  669
         _{arch\_info\_end} = .;
  670
          _{tagtable\_begin} = .;
         *(.taglist.init)
  671
         tagtable end = .;
vm/inux.lds 结束处
```

```
em.map vmlinux.lds vmlinux.lds piggy.S
] 😅 👛 🗐 🥥 🚵 🖹 🔡 🔡 🔞 🗞 눸 🖺 stack
                                       _edata =
  778 }
  779
  780
  781
  782
  783
        * (COMMON)
  784
  785
  786
           /* Stabs debugging sections. */
      .stab 0 : { *(.stab) }
      .stabstr 0 : { *(.stabstr) }
  789
      .stab.excl 0 : { *(.stab.excl) }
  790 .stab.exclstr 0 : { *(.stab.exclstr) }
791 .stab.index 0 : { *(.stab.index) }
      .stab.indexstr 0 : { *(.stab.indexstr) }
  792
      .comment 0 : { *(.comment) }
```

2. 生成 system. map, 置于顶层目录之下

3. arm-linux-gnu-obj copy, 去掉顶层 vmlinux 两个段-R. note-R. comment 的调试信息,减小映像文件的大小,此时大概 3. 2M, 生成 arch/arm/boot/I mage

4. gzi p -f -9 < arch/arm/boot/compressed/../Image > arch/arm/boot/compressed/pi ggy. gz, 读入 arch/arm/boot/Image 的内容,以最大压缩比进行压缩,生成 arch/arm/boot/compressed/目录下的pi ggy. gz

5. arm-linux-gnu-gcc,在 arch/arm/boot/compressed/pi ggy. S 文件中是直接<mark>引入 pi ggy. gz 的内容(pi ggy. gz 其实已经是二进制数据了),然后</mark>生成 arch/arm/boot/compressed/pi ggy. o 文件。下面是 pi ggy. S 的内容

**批注** [s3]: 内核符号表,如果 你想调试内核或者查找出错 倪端,多记得关注这个"美女" ^\_^

批注 [s4]: 此时转换成了 BIN 格式文件

批注 [s5]: 通过 incbin 关键字引入的,sorry,你后面有说啊, 重复了哈

```
tem.map vmlinux.lds vmlinux.lds piggy.s

| Section piggydata, #alloc | 2 .globl input_data | 3 input_data:
| 1 .incbin "arch/arm/boot/compressed/piggy.gz" | 5 .globl input_data_end | 6 input_data_end | 7
```

其中所选择的行就是加入了 pi ggy. gz 的内容,通过编译生成 pi ggy. o 文件,以备后面接下来的 I d 链接。

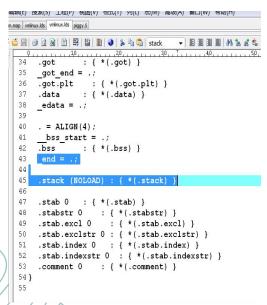
arm-linux-gnu-ld,在 arch/arm/boot/compressed/piggy. o 的基础上,加入重定位地址和参数地址的同时,加入内核解压缩的代码。 arch/arm/boot/compressed/head. o 、 misc. o ),最后生成arch/arm/boot/compressed 目录的 vmlinux,此时在解压缩代码中还含有调试信息(根据 arch/arm/boot/compressed/vmlinux.lds 来链接 0x0) vmlinux.lds 开始处

¥(E) 搜索(S) 工程(P) 视图(V) 格式(T) 列(L) 宏(M) 高级(A) 窗 p vmlinux.lds vmlinux.lds piggy.S 📙 | 🥝 💁 🔊 | 🖹 | 🔡 | 🔡 | 🔌 | 🐎 陷 👛 | stack 10 OUTPUT\_ARCH(arm) 11 ENTRY (\_start) 12 SECTIONS 13 { 14 . = 0; \_text = .; 15 16 .text : { \*(.start) 20 \*(.text) 21 \*(.fixup) 22 \*(.gnu.warning) 23 \*(.rodata) \*(.rodata.\*) 24 25 \*(.glue\_7) \*(.glue\_7t) \*(.piggydata) = ALIGN(4);29

批注 [s6]: 还能说啥呢,从事解压缩这么辛苦又不受内核特见的工作就是他们了(因为进真正内核之后就不用他们了,为他人做嫁衣裳啊,其人生何其惨烈!,但我们还是很想听他悲惨的人生历程来乐呵乐呵,歹毒吧^^)

注意到了 27 行的吗? \*(.pi ggydata)就表示需要将 pi ggydata 这个段放在这个位置,而 pi ggydata 这个段放的是什么呢?往后翻翻,看看第五步的图片,呵呵,其实就是将按最大压缩比压缩之后的 I mage,压缩之后叫 pi gqy. qz 中的二进制数据。

再看看 vml i nux. l ds 的结束处



注意到了。stack 放在什么位置了吗?放在了arch/arm/boot/compressed/vmlinux的bss\_end之后后,也就是说<mark>放在了piggy.gz代码段的后面</mark>。清楚这一点很重要,这关系到后边理解如何分布堆栈和堆,以及解压的地址空间分配。

- 7. arm-linux-gnu-objcopy, 去掉解压缩代码中的调试信息段,最后生成arch/arm/boot/目录下的zlmage
- 8. /bi n/sh

/home/farsight/Resources/kernel/linux-2.6.14/scripts/mkuboot.s h -A arm -0 linux -T kernel -C none -a 0x30008000 -e 0x30008000 -n 'Linux-2.6.14' -d arch/arm/boot/zlmage arch/arm/boot/ulmage 调用 mki mage 在 arch/arm/boot/zlmage 的基础上加入 64 字节的 ulmage 头,和内核的入口地址,装载地址, 最终生成 arch/arm/boot/目录下的 ulmage 文件。

由上可知,zl mage 中包含了内核的自解压代码(head o mi sc. o),所以在 uboot 将 64 字节的 ul mage 头去掉了之后,内核实际上就变成了一个 zl mage,所以接下来的工作首先会进行内核的自解压。

需要注意的是,内核代码(顶层 vmlinux-→arch/arm/boot/Image-→arch/arm/boot/compressed/piggy.gz-→arch/arm/boot/compressed/piggy.o)在编译链接时的地址是 0xc0008000,而自解压代码(arch/arm/boot/compressed/vmlinux-→ arch/arm/boot/zlmage)的链接地址却是0x0,所以在接下来的代码分析中所有的地址修正都是从0x00000000(自解压代码的链接地址0x0)开始的地址修正到0x30008000处,直到自解压代码执行完成,开始进入真正内核(由 Image 解压成 vmlinux 之后的链接地址是在0xc0008000开始处)

批注 [s7]: 这个说法欠妥,虽然快毕业了,但你还得跟我学哈,应该是放在了.BSS 段后面,由于链接的时候 BSS 实际数据为 0,所以严格来说,应该是链接的时候放在了整个内核压缩包(包括解压缩代码和内核本身压缩镜像后面)

**批注** [s8]: 此时是 bin 格式文件

批注 [s9]: 嗯,如果你们班还 有谁不懂这个,叫他来找我, 我得严肃地,语重心长地,苦 口婆心地和他谈心了,嘿嘿



# 二、内核自解压代码分析

下面分析的是自解压代码, 其链接地址是 0x0

- 一、 arch/arm/boot/compressed/head.s 和 mi sc.c 文件分析
  - 1. 首先来一小段延时,接着保存由 uboot 的 thekernel 传递进来的三个参数中的前两个,r0—>r8, r1-→r7.
  - 2. 导入 LCO 数据对象的数据到一系列的寄存器中,同时计算出这一系列链接地址值(0x0) 到当前空间 0x30008000 的偏移量,如下

- 3. 接下来对导入到各个寄存器中的值(链接时确定,基于链接地址 0x0)修正到当前运行空间 0x30008000,因为 CONFIG\_ZB00T\_ROM 没定义,所以修正了 r2, r3, r5, r6, ip, sp 到当前运行空间,另外也对 r6, ip 两寄存器值所代表的 got(全局偏移表)区间中的所有函数入口 进行地址修正。
- 4. clear bss
- 5. 为调用解压函数而准备 C 语言运行环境,包括打开 cache, 堆空间的分配 堆栈空间是利用下面的语句来分配了一个 4K 的空间 . align

```
.section ".stack", "w" user_stack: .space 4096
```

```
.type LCO, #object @定义了一个数据对象,有点类似于c中的结构体
LCO: .word LCO @ r1
.word _bss_start @ r2
.word _end @ r3
.word zreladdr @ r4 zreladdr 根据注释是virt_to_phys(TEXTADDR), 一般来说zreladdr = 0x300080
.word _start @ r5
.word _got_start @ r6
.word _got_end @ ip
.word _user__stack+4096 @ sp
LC1: .word reloc_end - reloc_start
.size LCO, . - LCO
```

bl cache\_on

/\* 开启指令和数据Cache \*/

mov r1 , sp @ malloc space above stack add r2 , sp , #0x10000 @ 64k max //在堆栈之后分配64K的临时堆空间

综上,可以看出,stack 的空间是直接安排在了 zl mage 代码的\_end 的后面,大小为 4K (具体的地址大小还需要看内核是多大),接着在 stack 的后面紧接着分配了 64K 的临时堆空间,用来进行<mark>自解压的实现</mark>。

这里需要特别说明一下的是 cache\_on 函数中,会打开 mmu,并且会做一些页表初始化的工作,然后在该 head. S 结束的时候会关闭 mmu 和 cache。在需要注意的是在这里初始化的页表并没有真正的使用价值,只是为了打开解压缩代码使用 D-cache 而打开的,所以这里的<mark>页表中映射的虚拟地址和物理地址的关系式等价映射的,也就是 1:1 映射</mark>,也就是说将物理地址映射到等值的虚拟地址上去,这里的页表就是这样映射出来的。详细的注释可以参考

)我的另一篇文章 ...compressed.head.S-v2-cache\_on 注释

根据 zImage 解压之后的 Image 入口地址 0x30008000 (final kernel address),和 zImage 目前的运行地址 (zImage 链接地址是 0x0,但是其前面的自解压代码是地址无关的,所以理论上可以放在 SDRAM 的任意空间,我们这里也是属于通常的情况,在 uboot 中将 uboot 头去掉后的 zImage 放在了 0x30008000 地址上),分三种情况来判断是否留有足够的存放 Image 的空间(4M 足够了,原则是尽可能得直接将 Image 解压到 0x30008000 地址上,避免了解压到其他地址之后再将代码搬运到 0x30008000 地址上来运行的繁琐和时间浪费)。不过通常情况下,kernel (Image)的入口地址和 uImage 的 load 地址或入口地址都设置成了相同的值 0x30008000。

那么关于 I mage 的入口地址 zrel addr 和 ul mage 的 I oad 地址、入口地址 如何设置呢?

Zrel addr 值在将内核编译链接成 vml i nux 的时候就已经确定了,但是这个值在哪里设定的呢?其实他是在arch/arm/mach-s3c2410/Makefile.boot定义的值,下图所示:

zreladdr-y := 0x30008000 params\_phys-y := 0x30000100

同时也看到了内核中约定的参数存放地址在哪里定义了,这里定义了后,这个变量在哪里引用的呢? arch/arm/boot/Makefile

**批注** [s10]: SP 指针指向了这个 4K 字节的顶端

批注 [s11]: 后面 C 代码 misc.c 的里面的 malloc()函数就是从这里分配内存的,类似我们平常说的堆栈中的堆,从这里可以看到,从内存布局来说,内核解压缩时候,"堆"在"栈"的上面,是不是有种"一切尽在掌握",自己成了上帝的感觉 (我们日常生活中"顾客就是上帝"一般都是喊的口号,但你跟代码打交道嘛,我们还是可以真正做回"上帝"的,嘿嘿)?

批注 [s12]: 心真细啊,明显感觉"长江后浪推前浪,前浪死在沙滩上"啊,这个我可是没提到啊,怕是那天我要在沙滩上啊,好怕怕^\_^



```
# Note: the following conditions must always be true:
# ZRELADDR == virt_to_phys(TEXTADDR)
# PARAMS_PHYS must be within 4MB of ZRELADDR
# INITRD_PHYS must be in RAM
ZRELADDR := $(zreladdr-y)
PARAMS_PHYS := $(params_phys-y)
INITRD_PHYS := $(initrd_phys-y)
```

#### export ZRELADDR INITRD\_PHYS PARAMS\_PHYS

在 arch/arm/boot/compressed/Makefile 文件中有下面片段,看英文注 释的意思就是,将 ZRELADDR 通过链接符号提供给解压程序使用。

# Supply ZRELADDR, INITRD\_PHYS and PARAMS\_PHYS to the decompressor via
# linker symbols. We only define initrd\_phys and params\_phys if the
# machine class defined the corresponding makefile variable.

LDFLAGS\_vmlinux := --defsym zreladdr=\$(ZRELADDR)
ifneq (\$(INITRD\_PHYS),)
LDFLAGS\_vmlinux += --defsym initrd\_phys=\$(INITRD\_PHYS)
endif
ifneq (\$(RAPAMS\_PHYS),)

ifneq (\$(PARAMS\_PHYS).)
LDFLAGS\_vmlinux += --defsym params\_phys=\$(PARAMS\_PHYS)
endif
LDFLAGS\_vmlinux += -p --no-undefined -X \
 \$(shell \$(CC) \$(CFLAGS) --print-libgcc-file-name) -T

最后将顶层的 ymlinux 在最大压缩比压缩之后得到了 I mage。

ul mage 的 l oad 地址和入口地址在哪里设置的呢? 其实应该是 zl mage 的 l oad 地址和入口地址的设定。因为 ul mage 只不过是在 zl mage 的前面加了一个 uboot 头,这个头里面保存了 zl mage 的 l oad 地址和 entry 地址。

怎么设定的是需要追寻到 ul mage 的生成过程了,在得到 zl mage 之后,通过工具 mki mage 才生成 ul mage 的,但是是否还记得当时是传递了很多参数的呢?下图

/bin/sh /home/farsight/Resources/kernel/linux-2.6.14/scripts/mkuboot.sh -A arm -0 linux -T kernel -C none a 0x30008000 o 1 linux-2.6.14 o darch/arm/boot/zImage arch/arm/boot/uImage 最终生成arch/arm/boot/目录下的uImage文件

光标选中的部分是不是传递了两个参数-a 表示 data I oad 地址, -e 表示的就是 entry 地址。但是这个值是在哪里设定的呢? arch/arm/boot/Makefile中,下图选中部分就是

好,现在知道了这两种地址是什么意思了吧!可以看出,通常这三个地址是一样的,都是 0x30008000,那接着往下分析,先贴出代码:



```
* Check to see if we will overwrite ourselves.
   r4 = final kernel address
   r5 = start of this image
   r2 = end of malloc space (and therefore this image)
* We basically want:
   r4 >= r2 -> 0K
   r4 + image length <= r5 -> 0K
        cmp r4 , r2
                         /× r4 = 0x30008000
                    r2 = (zImage+bss size)+stack size 4K + malloc size 64K ×/
        bhs wont_overwrite
        add r0 , r4 , #4096*1024
                                      @ 4MB largest kernel size
        cmp r0 , r5
                         /* r0= 0x30400000 , r5 = 0x30008000 */
        bls wont_overwrite
        mov r5, r2
                              @ decompress after malloc space
       mov r0 , r5 mov r3 , r7
                              \theta r\theta = malloc end or decompress space,
                              @ r1 = sp end or malloc begin
                              @ r2 = malloc end , r3 = architecture ID
           decompress_kernel
        add r\theta , r\theta , #127
                                @ align the kernel length
        bic r\theta , r\theta , #127
```

上图中 r4 代表 I mage 内核的真正入口地址 0x30008000, r5 代表 zI mage 的入口地址, r2 代表 zI mage 的堆空间结束地址。

下面分三种情况来看这段代码是如何判定空间之间是否有覆盖区域的。 第一种情况: r4>=r2, 假如 r4=0x30800000 (也就是 I mage 的入口地址), r2=0x30008000 (z1mage), si zo 4K + malloc si zo 4K

r2=0x30008000+(zImage+bss size)+stack size 4K + malloc size 64K 这中情况 zImage 区域没有覆盖掉解压后 Image 需要的空间,所以可以直接解压到 0x30800000 中去,也省去了解压后代码搬移的过程。

第二种情况: r4+4M<r5, 假如 r5=0x30800000, r4=30008000。这种情况下也是可以直接将 zI mage 解压到 r4 空间的。同第一种情况,这也不用解压后搬移代码。

第三种情况: 也是我们普遍使用的情况。r4+4M>r5,就拿我们正常的地址值来说。此时 r4=r5=0x30008000,这样就得需要另一种处理方式,先将内核解压到其他空闲的地址,然后将解压后的代码搬运到 0x30008000 处。

- @ r0 = malloc end or decompress space,
- @ r1 = sp end or malloc begin,
- @ r2 = malloc end .
- @ r3 = architecture ID

上面几个寄存器的值都是需要传递给 decompress\_kernel 函数的,该函数将内核解压到 malloc 空间之后去了。之后可定要将解压后的代码搬移会到 0x30008000 处运行。

7. 调用 decompress\_kernel 函数, in arch/arm/boot/compressed/misc.c



Gunzip()函数在 lib/inflate.c 中,在 gunzip 函数中调用的重要函数是inflate()函数,也在 inflate.c 文件中。

这里<mark>不对解压算法做详细分析</mark>,只要知道解压完的代码是放在 mal oc 空 间之后就行了。

在汇编中调用 c 函数返回的值保存在 r0 中的,即现在 r0 中保存的就是解压后代码的长度,但是需要进行 128 字节对齐,所以采用下面的两句代码

add r0, r0, #127

bic r0, r0, #127 @ align the kernel length

9. 下面的图片展示的就是将 head. S 本文件中的用来搬运解压后的内核的代码部分搬运到解压后内核的未尾处,以保证在运行该部分代码搬移解压后内核到 0x30008000 处的时候,自己不会被覆盖掉

```
add r1 , r5, r0
                                  @ end of decompressed kernel
        adr r2, reloc_start
        ldr r3 , LC1
        add r3 , r2 ,
                       \Gamma 3
1:
        ldmia
                r2!, { r8 - r13}
                                          @ copy relocation code
        stmia
                r1!, { r8 - r13}
        ldmia
                 \Gamma 2!, { \Gamma 8 - \Gamma 13}
        stmia
                r1!, { r8 - r13}
        cmp r2 , r3
        blo 1b
        bl cache_clean_flush
                                 @ call relocation code
        add pc, r5, r\theta
```

R1 为解压后内核的结束地址(已经 128 字节对齐了), r2 为本文件中 rel oc\_start 函数的当前运行空间地址, r3 经过修正之后变成了 rel oc\_end 的相 对当前运行空间的地址。这两个标号都是在本 head. S 文件中

搬运结束之后,调用 cache\_cl ean\_fl ush 刷新 cache,因为代码地址变化了不能让 cache 再命中即将被内核覆盖的老地址。接下来就跳到新地址继续运行,也就是到新地址处去执行 reloc start 函数。

另外需要注意的是,zImage 在解压后的搬移和跳转会给 gdb 调试内核带来麻烦。因为用来调试的符号表是在编译是生成的,并不知道以后会被搬移到何处去,只有在内核解压缩完成之后,根据计算出来的参数"告诉"调试器这个变化。

9. 下面是 reloc\_start 函数的代码:

批注 [s13]: 嗯, 其实这个解压缩也比较有技巧的哈



```
* r0
          = decompressed kernel length
 * r1-r3 = unused
* r4
           = kernel execution address
* r5
           = decompressed kernel start
 * r6
           = processor ID
 * r7
           = architecture ID
  r8-r14 = unused
           .align 5
                  add r8, r5, r0
 reloc_start:
          debuq_reloc_start
          mov r1 , r4
           rept
                  r5!, { r\theta , r2 , r3 , r9 - r13 } r1!, { r\theta , r2 , r3 , r9 - r13 }
          ldmia
                                                          @ relocate kernel
          stmia
          cmp r5, r8
          blo 1b
          debug_reloc_end
 call kernel:
                  bl cache_clean_flush
          bl cache_off mov r0, #0
          mov r1 , r7
                                    @ restore architecture number
          mov pc , r4
                                    @ call kernel
```

首先将解压后的内核代码搬运到 0x30008000 地址处,然后清除并关闭 cache,将 r0 和 r1 设置成对应的之后,最后将 pc 指针指向了 0x30008000 处。http://forum.linuxbj.com/node/3

# 三、VMLINUX 内核启动代码分析

从下面开始就进入了真正的内核了,首先执行的代码是在 arch/arm/kernel/head.S 中的汇编代码,然后就进入了内核的 c 代码部分 (init/Main.c),注意,此刻内核的编译链接地址不再是 0x0 了,而是 0x00008000 了。

内核汇编部分: (关于该部分参考我的另一篇详细的分析) 刚刚进入内核汇编时的一些系统状态和寄存器值如下

Mmu 关闭,I Cache 和 D Cache 清除并关闭,r0=0,r1=archi tecture ID。 从 arch/arm/kernel /vml i nux. I ds 文件中可以看出,解压后内核是从 stext 段开始执行的。

- 1. 确保系统处于 SVC 模式并且 FIR、IRQ 都已经关闭、
- 2. 调用\_lookup\_processor\_type 函数,通过协处理器 cp15 读取出系统 cpu 的 id, 然后通过查找内核映像中.proc.info.init 段的所有 proc\_info\_list 数据结构,以判断该内核是否支持该款 cpu。
- 3. 调用\_\_lookup\_machine\_type 函数,检查由 uboot 传递进来的 machine architecture number 是 否 被 内 核 支 持 , 也 就 是 在 内 核 的 . arch.info.init 段中查找看是否有对应 machine number 的 machine\_desc 数据结构体存在。

关于该部分汇编代码的分析由于涉及到很多内容,所以将该部分的分析转移到一个单独的文件中去分析。arch-arm-kernel-head. S 文件中

arch-arm-kernel-head

参考网址:

http://forum.linuxbj.com/node/4

下面进入 start\_kernel 函数去执行, in init/Main.c 文件

- . printk(linux\_banner)打印内核的一些信息,版本,作者,编译器版本,日期等信息
- \_\_setup\_arch(&command\_line); /\* in arm/kernel/setup.c \*/ 函数原型: void \_\_init <mark>setup\_arch</mark>(char \*\*cmdline\_p)

  - b. mdesc = setup\_machi ne(machi ne\_arch\_type)函数 取出当前系统的在内核. arch. i nfo. i ni t 段内对应的 machi ne\_desc 数据 结构体的地址 l i st,并打印出 l i st-name 内容
  - c. machi ne\_name = mdesc->name 设置全局变量 machi ne\_name

数,在 proc\_arm920. S 文件中定义的。

- d. tags = phys\_to\_virt(mdesc->boot\_params), 修改默认的参数地址, 使用 uboot 传递进来的参数的真实地址 0x30000100 并将其转换成虚拟地址 0xc0000100
- e. 参数解析

```
if (tags->hdr.tag == ATAG_CORE) {
   if (meminfo.nr_banks!= 0) /* meminfo defined in setup.c */
        squash_mem_tags(tags);
   parse_tags(tags);
}
static struct meminfo meminfo __initdata = { 0, };in steup.c

所以这里会调用 parse_tags(tags) (in steup.c) 函数
```

批注 [s14]: 嗯, 重中之重啊

Steup\_arch()-→parse\_tags()-→parse\_tag(), all function in steup.c

\_\_tagtable\_begin, \_\_tagtable\_end 这两个参数是在arch/arm/kernel/vmlinux.lds文件中确定的,这个区间存放了各种参数的解析函数,可以直接引用

parse\_tag()函数就是查找这个区间的各种参数头,来和传递寄来的参数头进行匹配,找到了之后就利用其中的函数指针 parse 来进行参数解析,具体怎么解析需要分析各种参数的作用了。如果匹配不到的话,会打印信息说 I gnori ng unrecogni sed tag 0x%08x\n,表示 uboot 传递进来了一个内核识别不了的参数。

f. struct mm\_struct init\_mm = INIT\_MM(init\_mm);

init\_mm.start\_code = (unsigned long) &\_text;
init\_mm.end\_code = (unsigned long) &\_etext;
init\_mm.end\_data = (unsigned long) &\_edata;
init\_mm.brk = (unsigned long) &\_end;
\_text, \_etext, \_edata, \_end 都是在 arch/arm/kernel/vmlinux.lds
文件中确定的

- g. parse\_cmdline(cmdline\_p, from)解析 uboot 传递寄来的 commond\_line 字符串,通过 from 指针内的参数解析到 command\_line 这个全局数组里,然后将该数组的地址赋值给上上级传进来的参数 command\_line (start\_kernerl)。这里只分析了当中的 mem、initrd 部分,另外一处是 start\_kernel()->parse\_option()用于解析其余部分。将 mem 和 initrd 参数从原始 command\_line 中扣出,留下的部分放在\*cmdline\_p 中(start\_kernel 传进来的参数,第二次分析会用到)。因为接下来是建立内存的映射页表,需要用到 mem、initrd 两个参数的内容,所以在这里提前解析出来。
- h. paging\_init(&meminfo, mdesc); /\* 这部分的主要工作建立页表,初始 化内存 \*/ in arch/arm/mm/init.c 文件中

memtable\_init(mi)为系统内存创建页表,是一个比较重要的函数,详细的参考网址

http://bbs.sjtu.edu.cn/bbstcon,board,Embedded,reid,1165977462,

的内容,需要注意的是这个函数中将中断向量表映射到了 0xffff0000 开始处

调 用 mdesc->map\_io() 函 数 , 在 smdk2410 中 位 于 arch/arm/mach-s3c2410/mach-smdk2410.c文件中的 smdk2410\_map\_io() 函数,关于这部分的外设的静态映射分析请参考易松华老师的文档: ARM Linux 静态映射分析, smdk2410\_map\_io()函数主要做了下面几件事情:

1)、iotable\_init(s3c\_iodesc, ARRAY\_SIZE(s3c\_iodesc)) 建立GPIO, IRQ, MEMCTRL, UART 的静态映射表

**批注** [s15]: 别参考了, 到你们 这个学习阶段, 应该当我已经 死在沙滩上了, 自己分析哈 ^\_^ 2)、(cpu->map\_io)(mach\_desc, size) 建立其他外设的静态映射表,包括网卡,LCD,看门狗等

这里需要明确一点的是在 map. h 文件中的

 $S3C2410\_ADDR(x)$  ((void \_\_i omem \*)0xF0000000 + (x))

表明我们处理的 I0 和外设的寄存器被静态映射到了 0xF0000000 开始地方,其中每一种外设的空间占 1M 大小。

这里必须得提到的是内存的三种映射和使用方式: 1. 处理器的 GPI 0,各种外设(IRQ,UART,MEMCRTL,WATCHDOG,USB等)的寄存器的静态映射,也就是上面红色字体体现出来的。这是映射到内核空间内,所以用户程序是不能访问的,只能由内核来访问。2. 内核空间访问的虚拟地址(比如是代码的地址或者是 kmal I oc 空间的地址)都是在 0x30008000 开始后的空间,也就是说这中虚拟地址是通过静态映射得到的,这部分也只能由内核访问。phys\_to\_virt 和 virt\_to\_phys ,3. 除此之外的物理地址都是通过 mmu 的规则映射出来的,TTB 已经确定在 0x30004000 开始的 16K空间就是一级页表的空间。通过这中方式映射的虚拟地址用户程序和内核程序都可以访问,不过解析式通过 mmu 来完成的而已。参考arch/arm/kernel/head. S 中建立的 4M空间手工页表就是按照 mmu 的规则建立的(只不过是按照 sector 的方式建立)

3)、对时钟, Uart 的简单初始化。

h. request\_standard\_resources(&meminfo, mdesc) 为 memory , kernel\_text,kernel\_data,video\_ram 在资源树种申请标准资源。它完成 实际的资源分配工作,如果参数 new 所描述的资源中的一部分或全部已经被 其它节点所占用,则函数返回与 new 相冲突的 resource 结构的指针。否则 就返回 NULL

i.cpu\_init()函数分析in arch/arm/kernel/setup.c

打印一些关于 cpu 的信息,比如 cpu i d, cache 大小等。另外重要的是设置了 I RQ、ABT、UND 三种模式的 stack 空间,分别都是 12 个字节。最后将系统切换到 svc 模式。

- i. 用\_\_mach\_desc\_SMDK2410\_type 结构体中特定成员来初始化系统的这么写全局变量 i ni t\_arch\_i rq, i ni t\_machi ne, system\_timer
- 3. sched\_init()函数

初始化每个处理器的可运行队列,设置系统初始化进程即 0 号进程。也就是调用了 i ni t\_i dl e (current, smp\_processor\_id())初始化 idle 当前进程。

- 4. preempt\_disable()禁止抢占
- 5. 建立系统内存页区(zone)链表 build\_all\_zonelists()
- 6. printk(KERN\_NOTICE "Kernel command line: %s\n", saved\_command\_line); 打印出从 uboot 传递过来的 command\_line 字符串, 在 setup\_arch 函数中获得的。
- 7. parse\_early\_param(),这里分析的是系统能够辨别的一些早期参数(这个函

数甚至可以去掉,\_\_setup 的形式的参数),而且在分析的时候并不是以setup\_arch(&command\_line)传出来的 command\_line 为基础,而是以最原生态的 saved command line 为基础的。

对于比较新的版本真正起作用的函数,与 parse\_earl y\_param();相比,此处对解析列表的处理范围加大了,解析列表中除了包括系统以 setup 定义的启动参数,还包括模块中定义的 param 参数以及系统不能辨别的参数。

command\_line是setup\_arch函数传递出来的值;

\_\_start\_\_\_param 是 param 参数的起始地址,在 System. map 文件中能看到

\_\_stop\_\_\_param - \_\_start\_\_\_param是参数个数

unknown\_bootopti on是对应与启动参数不是param的相应处理函数(查看parse\_one()就知道怎么回事)

函数 parse\_one()既可以处理 param 参数,又可以处理\_\_setup 的形式的参数,还可以处理不能识别的参数。当然这些都是依靠传递进来的参数进行分支处理的。

参数的处理详情参考网络上另一篇文章: Linux 启动 bootargs 参数分析

- 9. sort\_main\_extable() 将放在\_\_start\_\_ex\_table 到\_\_stop\_\_\_ex\_table 之间的\*(\_\_ex\_table)区域中的 struct exception\_table\_entry 型全局结构变量按 insn 成员变量值从小到大排序,即将可能导致缺页异常的指令按其指令二进制代码值从小到大排序。
- 10. 在前面的 setup\_arch-→paging\_init-→ memtable\_init 函数中为系统创建页表的时候,中断向量表的虚地址 init\_maps,是用alloc\_bootmem\_low\_pages分配的,ARM 规定中断向量表的地址只能是0或0xFFFF0000,所以该函数里有部分代码的作用就是映射一页到0或0xFFFF0000。

trap\_i ni t 函数做了一下的工作: 把放在. Lcvectors 处的系统 8 个意外入口跳 转 指 令搬 到 高端 中断 向量 0xffff0000 处,再将\_\_stubs\_start 到\_\_stubs\_end 之间的各种意外初始化代码搬到 0xffff0200 处。将系统调用的返回句柄拷贝到 0xffff0500 处。刷新 0xffff0000 处 1 页范围的指令 cache,将 DOMAI N\_USER 的访问权限由 DOMAI N\_MANAGER 改成 DOMAI N\_CLI ENT 权限。

- 11. rcu\_i ni t()函数初始化当前 cpu 的读、复制、更新数据结构(struct rcu\_data) 全局变量 per\_cpu\_rcu\_data 和 per\_cpu\_rcu\_bh\_data.
- 12. i ni t\_I RQ 函数

初始化系统中支持的最大可能中断数的中断描述结构 struct i rqdesc 变量数组 i rq\_desc[NR\_IRQS], 把每个结构变量 i rq\_desc[n]都初始化成预先定义好的坏中断描述结构变量 bad\_i rq\_desc, 并初始化该中断的链表表头成员结构变量 pend。

执行 init\_arch\_irq 函数,该函数是在 setup\_arch 函数最后初始化的一个全局函数指针,指向了 smdk2410\_i nit\_irq 函数(in mach-smdk2410.c),实际上是调用了 s3c24xx\_init\_irq 函数,在该函数中,首先清除所有的中断未决标志,之后就初始化中断的触发方式和屏蔽位,还有中断句柄初始化,这里不是最终用户的中断函数,而是 do\_level\_IRQ 或者 do\_edge\_IRQ 函数,在这两个函数中都使用过\_\_do\_irq 函数来找到真正最终驱动程序注册在系统中的中断处理函数。



```
status = 0;
do {
    ret = action->handler(irq, action->dev_id, regs);
    /* 这里才是用户的真正中断处理程序的执行 */
    if (ret == IRQ_HANDLED)
        status |= action->flags;
    retval |= ret;
    action = action->next;
} while (action);
if (status & SA SAMDIE DANDOM)
```

接着初始化外部中断的一些参数,最后补充初始化 uart 和 ADC 中断。

#### 13. pi dhash\_i ni t()函数

设置系统中每种 pi d\_hash 表中的 hash 链表数的移位值全局变量 pi dhash\_shi ft,将 pi dhash\_shi ft 设置成 mi n(12)。分别为每种 hash 表的 连续 hash 链表表头结构申请内存,把申请到的内存虚拟基址分别传给 pi d\_hash[n] (n=1~3),并将每种 hash 表中的每个 hash 链表表头结构 struct hl i st\_head 中的 fi rst 成员指针设置成 NULL

#### 14. init\_timers()函数

初始化当前处理器的时间向量基本结构 struct tvec\_t\_base\_s 全局变量 per\_cpu\_tvec\_bases, 初始化 per\_cpu\_tvec\_bases 的自选锁成员变量 lock

# 15. softirq\_init函数

内核的软中断机制初始化函数,



```
void __init SOftirq_init(void)
{
/*HI_SOFTIRQ用于实现bottom half, TASKLET_SOFTIRQ用于公共的tasklet*/
    open_softirq(TASKLET_SOFTIRQ, tasklet_action, NULL);
    open_softirq(HI_SOFTIRQ, tasklet_hi_action, NULL);
}
```

16. time\_init()这个函数是用来做体系相关的 timer 的初始化

```
void __init time init(void)
    if ( system_timer ->offset == NULL)
         system_timer ->offset = dummy_gettimeoffset;
    system_timer ->init();
还记得在 setup_arch 最后初始化了三个全局的变量吗?下图所示
  * Set up various architecture-specific pointers
  .
//运行mach-smdk2410.c文件中定义的结构体__mach_desc_SMDK2410中的特定函数和结构
 init_arch_irq = mdesc->init_irq;/* smdk2410_init_irq */
system_timer = mdesc->timer; /* 3c24xx_timer */
 init_machine = mdesc->init_machine;/* smdk2410_init_fs2410 */
static void __init s3c2410_timer_init (void)
     s3c2410_timer_setup();
     setup_irq(<mark>IRQ_TIMER4</mark>, & s3c2410_timer_irq);
struct sys_timer s3c24xx_timer = {
     .init = s3c2410_timer_init,
     offset
                = s3c2410_gettimeoffset,
                = s3c2410_timer_setup
-};
```

很显然, time\_init 函数中的 if 条件不成立, 所以就是直接执行了函数 s3c2410\_timer\_init 函数。可以看出, 系统使用了 timer4 来作为系统的定时器。

### 17. consol e\_i ni t 函数

初始化系统的控制台结构,该函数执行后调用 printk 函数将 log\_buf 中所有符合打印级别的系统信息打印到控制台上

a. tty\_register\_ldisc()注册默认的 TTY 线路规程 大家研究 tty 的驱动就会发现了在用户和硬件之间 tty 的驱动是分了三 层最底层当然是 tty 驱动程序了,主要负责从硬件接受数据,和格式化 上层发下来的数据后给硬件。在驱动程序之上就是线路规程了,他负责 把从 tty 核心层或者 tty 驱动层接受的数据进行特殊的按着某个协议的



格式化,就像是 ppp 或者蓝牙协议,然后在分发出去的。在 tty 线路规 程之上就是 tty 核心层了。 b. 接下来就是执行平台相关的控制台初始化代码 s3c24xx\_serial\_initconsole 在 vml i nux. I ds. S 中连接脚本汇编中有这段代码 \_\_con\_initcall\_start = .; \*(.con\_initcall.init) \_con\_initcall\_end = .; 因此我们再此处调用的就是 con\_i ni tcall. i ni t 段的代码, 将 fn 函数放到. con\_initcall.init 的输入段中,如下: #define console\_initcall(fn) \ static initcall\_t \_\_initcall\_##fn \ \_\_attri bute\_used\_\_ attribute\_((\_\_section\_\_(".con\_initcall.init")))=fn //在我们串口驱动里面有 这么一个注册语句: console\_initcall(s3c24xx\_serial\_initconsole); //因此我们的控制台初始化流 程就是: start\_kernel ->consol e\_i ni t->s3c24xx\_seri al \_i ni tconsol e call = \_\_con\_initcall start; /\* console\_initcall \*/ while (call < \_\_con\_initcall\_end) {</pre> (\*call)(); call++; } 18. profile\_init()函数 /\* 对系统剖析做相关初始化, 系统剖析用于系统调用\*/ //profile 是用来对系统剖析的,在系统调试的时候有用。 //需要打开内核选项,并且在 bootargs 中有 profile 这一项才能开启这个功 能 profile 只是内核的一个调试性能的工具,这个可以通过 menuconfig

/\* 中 profiling support 打开。

1. 如何使用 profile:

首先确认内核支持 profile, 然后在内核启动时加入以下参数: profile=1 或者其它参数, 新的内核支持 profile=schedule 1

2. 内核启动后会创建/proc/profile 文件,这个文件可以通过 readprofile 读取,

如 readprofile /proc/kallsyms -m sort -nr ~/cur\_profile.log,

或者 readprofile -r -m /proc/kallsyms |sort -nr,

或者 readprofile -r && sleep 1 && readprofile -m /proc/kallsyms

|sort -nr >~/cur\_profile.log

3. 读取/proc/profile 可获得哪些内容?

根据启动配置 profile=? 的不同,获取的内容不同:

如果配置成 profile=? 可以获得每个函数执行次数,用来调试函数性能很有用

如果设置成 profile=schedule ?可以获得每个函数调用 schedule 的次数,用来调试 schedule 很有用

19. local\_irq\_enable() 使能 IRQ 中断

20. mem\_i ni t()

最后内存初始化,释放前边标志为保留的所有页面,这个函数结束之后就不能再使用 alloc\_bootmem(), alloc\_bootmem\_low(), alloc\_bootmem\_pages()等申请低端内存的函数来申请内存,也就不能申请大块的连续物理内存了

21. kmem\_cache\_i ni t()

执行高速缓存内存管理即 slab 分配器的初始化

22. numa\_policy\_init()/

if (late\_time\_init)

late\_time\_init();

calibrate\_delay();

//计算机系统的 BogMIPS 数值,即处理器没秒钟执行的指令数

23. pidmap\_init();

pgtabl e\_cache\_i ni t();

prio\_tree\_init();

/\*

初始化无符号长整型全局数组 index\_bits\_to\_maxindex[BITS\_PER\_LONG] 的每个组员

将每个组员 index\_bits\_to\_maxindex[n] 设置成-1,将最后的 index\_bits\_to\_maxindex[BITS\_PER\_LONG-1]设置成~OUL \*/

24. anon\_vma\_i ni t();

/\*

该函数调用 kmem\_cache\_creat() 函数,为匿名逆序内存区域链表结构 struct anon\_vma

创建高速缓存内存描述结构 kmem\_cache\_t 变量,为该变量命名为anon\_vma,其对象的

构造函数指针指向 void anon\_vma\_ctor 函数,析构函数指针为 NULL,将创建的

kmem\_cache\_t 结构变量地址传给全局指针 anon\_vma\_chachep

\*/

- 25. fork\_init(num\_physpages);/\* 进程创建的相关初始化 \*/
- 26. proc\_caches\_init();
   buffer\_init();
   /\*

调用 kmem\_cache\_create("buffer\_head", sizeof(struct buffer\_head), 0, SLAB\_RECLAIM\_ACCOUNT|SLAB\_PANIC, init\_buffer\_head, NULL)函数为缓冲区描述结构 struct buffer\_head 创建高速缓存内存描述结构 kmem\_cache\_t 变量 \*/

27\_\_securi ty\_i ni t(); /\* 打印安全架构方面的信息 \*/

28. vfs\_caches\_i ni t (num\_physpages);

radix\_tree\_i ni t();

signals\_init();

调用 kmem\_cache\_create("si gqueue",

sizeof(struct sigqueue),

\_alignof\_\_(struct siggueue),

SLAB\_PANIC, NULL, NULL) 函数为信号队列结构 struct si gqueue 创建高速缓存内存描述结构 kmem\_cache\_t 变量,名字叫 si gqueue,不要求其对象按照处理器硬件 cache\_li ne 大小对齐,没有定义其对象的构造和析构函数,将创建好的 kmem\_cache\_t 结构变量的地址传给全局指针si hqueue\_cachep.

29. page\_writeback\_init()

统计系统中所遇内存节点的通用内存页区中高页数水印数之外的额外内存总页数之和传给 buffer\_pages.

- 30. proc\_root\_init(); /\* 只有在系统支持 proc 文件系统即配置了 CONFIG\_PROC\_FS 选项事才被调用 \*/
- 31. check\_bugs(); /\* 在 arm 中是验证内存一致性 \*/
- 32. rest\_init(), 这是另一个重要的函数

该函数创建 init 内核进程,也就是 1 号进程,然后原来的为 0 号系统启动进程进入空闲状态

**linux** 进程流程: 0 号内核<mark>进程</mark>->1 号内核进程->1 号用户进程 下面分析 rest\_i ni t()函数

- a. kernel\_thread(init, NULL, CLONE\_FS | CLONE\_SIGHAND) in arch/arm/kernel/process.c 文件 在这个函数中创建了1号内核进程,即init
- b. schedule()进入idle之前,主动调用下schedule,看有没有其他进程
- c. cpu\_i dle()0 号进程进入空闲态

**批注 [s16]:** 笔误啊, 应该是内 核线程



init 函数分析 in main.c

- a. lock\_kernel() 这是1号进程的lock
- b. 不需要关注那些和 smp 相关的函数
- c. populate\_rootfs()函数,必须在initcalls之前调用它来安装文件系统, 因为有些驱动可能会访问

```
d. do_basic_setup()驱动加载
  #define module_init(x) __initcall(x);
#define __initcall(fn) device_initcall(fn)
```

上面所涉及的宏都来自 include/linux/init.h 文件

```
#define core_initcall(fn)
                                __define_initcall("1", fn)
#define postcore_initcall(fn)
                                    __defi ne_i ni tcal I ("2", fn)
#define arch_initcall(fn)
                                __define_initcall("3", fn)
#define subsys_initcall(fn)
                                __define_initcall("4", fn)
#define fs_initcall(fn)
                                __define_initcall("5", fn)
                                __define_initcall("6", fn)
#define device initcall(fn)
                                __define_initcall("7", fn)
#define Late_initcall(fn)
#define _define_initcall(level,fn) \
static initcall_t_initcall_##fn __attribute_used__ \
 attribute (( section (".initcall" level ".init"))) = fn
```

```
if (board != NULL) {/* board 在本文件中定义和在本文件中初始化 */
    struct platform_device **ptr = board->devices;
    int i;

for (i = 0; i < board->devices_count; i++, ptr++) {
        ret = platform_device_register(*ptr);

        if (ret) {
            printk(KERN_ERR "s3c24xx: failed to add board devi
        }
    }

    /* mask any error, we may not need all these board
        * devices */
    ret = 0;
}

return ret;
} end s3c_arch_init ?
```

## arch\_initcall(s3c\_arch\_init);

上图中的 arch\_i ni tcall 是属于 i ni tcall 第 3 级的, 驱动的 modul e\_i ni t 属于第 6 级,总之,不管是什么级别,被这些关键字修饰的说有函数都会在 i ni t 1 号内核进程中的 do\_basi c\_setup 函数中的 do\_i ni tcalls()函数来调用执行,当然执行的先后顺序就是按照 1 到 7 排列的。



```
顺序是:
Start_kernel *
   --→rest_i ni t()*
      --→ kernel_thread()* 开启 init 1号内核进程
        --→Ini t* 内核进程
           --→ popul ate_rootfs() 安装文件系统
           --→do_basic_setup()* 系统所有被 initcall 修饰的
函数被调用,7个级别的函数一次执行,其中 modlue_init 修饰的设备驱
动程序是在第6个级别上。
              --→i ni t_workqueues()初始化工作队列
              --→usermodehel per_i ni t()启动用户态的 khel per
进程
              --→driver_init()设备模型中一些基础结构体初
始化和设备的注册
              --→sysctl_init()系统调用初始化
              --→sock_init()网络socket 初始化
              --→do initcalls()*
                --→(*call)()* 具体的 initcall 函数调用
             →prepare_namespace()挂载文件系统主要工作所在
的地方
             ←¬→name_to_dev_t()这一步很重要,解析了参数中
的 root 选项, 获得要挂载的设备号。如果 bootargs 不是以 root=/dev/**
开始,则进入循环,我们也可以通过设备号,例 root=31:03 等价于
/dev/mtdblock3。root 参数的名字不能过长,即/dev/后面的个数不能超
过31。
              --→mount_root();/* 挂载设备 */
           --→free_initmem()将以init标示的函数空间释放,也
就是我们启动的时候经常会冒出来: Freeing init memory: 116K
           --→sys_open() and sys_dup(0)打开标准终端 0 1 2
           --→run_i ni t_process (execute_command) 运行第一个
用户进程,init=/initrd ,command_line 中传递进来的。
      --→ cpu idle() 系统启动 0 号进程进入空闲态
```

在 init 1号进程运行完下面的代码时,整个系统就算起来了

run\_i ni t\_process("/sbi n/i ni t"); run\_i ni t\_process("/etc/i ni t"); run\_i ni t\_process("/bi n/i ni t"); run\_i ni t\_process("/bi n/sh");

3

### 参考资料:

ARM Li nux 静态映射分析--易松华老师,华清远见内核启动过程-PXA255---易松华老师,华清远见Li nux 启动 bootargs 参数分析

■ <a href="http://blog.chinaunix.net/u3/99423/article.html">http://blog.chinaunix.net/u3/99423/article.html</a>

http://bbs.sjtu.edu.cn/bbstcon, board, Embedded, reid, 1165977462.

html

http://forum.linuxbj.com/node/4

全文完

2010/6/4 lizgo 李枝果 lizhiguo0532@163.com

**批注 [s17]**:后面的一段没看的太仔细,主要是困了哈,只能以后再仔细看吓了\*\_\*

