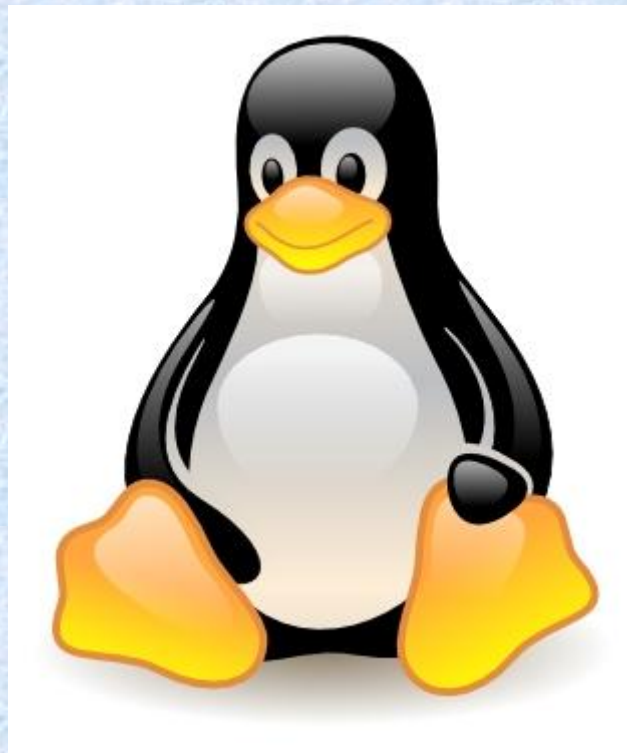


《Linux 内核修炼之道》精华版



前言	5
<i>本书的组织形式</i>	6
精华版方法论部分导读	8
<i>Linux大史记</i>	9
<i>内核学习的方法论</i>	9
<i>驱动开发的方法论</i>	12
<i>Linux内核问题门</i>	13
缅怀已逝的十八年（1991～1998）	14
<i>Linux诞生记</i>	15
<i>泰坦尼克的狂潮</i>	16
缅怀已逝的十八年（1999～2002）	17
<i>提前发生的革命</i>	17
<i>和平、爱情和Linux</i>	19
缅怀已逝的十八年（2003～2006）	20
<i>Ubuntu 4.10</i>	20
<i>Richard Stallman的征婚启事</i>	22
缅怀已逝的十八年（2007～2009）	24
<i>来自微软的指控</i>	24
<i>首款Android手机</i>	28
<i>Linux信用卡</i>	29
KERNEL地图：KCONFIG与MAKEFILE	31

<i>Makefile</i> 不是 <i>Make Love</i>	31
利用 <i>Kconfig</i> 和 <i>Makefile</i> 寻找目标代码.....	32
分析内核源码如何入手？（上）	34
分析 <i>README</i>	35
分析 <i>Kconfig</i> 和 <i>Makefile</i>	37
分析内核源码如何入手？（下）	41
态度决定一切：从初始化函数开始.....	41
内核学习的心理问题	49
内核学习的相关资源	50
内核文档.....	51
经典书籍.....	52
内核社区.....	54
其他网络资源.....	55
模块机制与“HELLO WORLD!”	56
设备模型（上）	59
设备模型（下）	65
内核中 <i>USB</i> 子系统的结构.....	65
<i>USB</i> 子系统与设备模型.....	67
驱动开发三件宝：SPEC、DATASHEET与内核源码	69
LINUX内核问题门——学习问题、经验集锦	70
Linux内核学习常见问题.....	71
Linux内核学习经验.....	74

《Linux内核修炼之道》精华分享与讨论（14）——内核中的链表	76
《Linux内核修炼之道》精华分享与讨论（15）——子系统的初始化：内核选项解析.....	83
《Linux内核修炼之道》精华分享与讨论（16）——子系统的初始化：那些入口函数.....	90
《Linux内核修炼之道》精华分享与讨论（17）——子系统的初始化：以PCI子系统为例	96
《Linux内核修炼之道》精华分享与讨论（18）——选择发行版	101
《Linux内核修炼之道》精华分享与讨论（19）——不稳定的内核API	104
《Linux内核修炼之道》精华分享与讨论（20）——学会使用GIT	106
《Linux内核修炼之道》精华分享与讨论（21）——二分法与PRINTK()	110
二分查找法的基本原理.....	110
<i>printk()</i>	111

前言

至此落笔之际，恰至 Linux 问世 18 周年，18 年的成长，如梦似幻，风雨颇多，感慨颇多。

犹自忆起多年以前一位前辈训导时的箴言：今天的必然正是由之前一系列的偶然所决定的。过去的某年某月，我偶然初识 Linux 就身陷其中，至今仍找不到出去的路，而正是这次乃至之后的多次偶然相联合，从而决定了今日的我要在此写下这些话。那么，当您偶然地拿起这本书，偶然地看到这段话，您是否会问自己：这样的偶然又会导致什么样的必然？

如果您依然决定继续这次的偶然之旅，那么首先请认识两个人，准确的说是一个人和一只企鹅。这个人自然就是 Linus Torvalds，我们也可称他为 Linus 或李纳斯，正是这位来自芬兰的天才，在 1991 年 1 月 2 日，攥着在圣诞节和生日得到的钱，偶然地做出了一个重大的财政决定，分期三年买一台价格 3500 美元得相貌平平得计算机，从而 Linux 开始了。

企鹅则是 Linux 的标志，很多人可能不知道 Linus，但是却可能知道这只企鹅，这是一个奇怪的现象，就像很多人知道微软，却不知道比尔盖茨。不管怎么说，是 Linus 塑造了这只企鹅，并让它有一副爽透了的样子，就像刚刚吞下一扎啤酒。除此之外，这只企鹅还要很特别，其他的企鹅都是黑嘴巴黑脚蹼，但它却是黄嘴巴黄脚蹼，这使它看上去好像是鸭子与企鹅的杂交品种，也许它是唐老鸭在南极之旅中与一只当地企鹅一夜倾情的结晶。

其次，在您继续之前，我还想请您问自己一个问题：我在强迫自己学习内核么？我很希望您能回答不是，但希望与现实往往都有段不小的距离，因为很多时候，我都发现身边的人是因为觉得内核很高深而强迫自己喜欢的。强迫自己去喜欢一个人是多么痛苦的事情。或许，针对这个问题，最让人愉悦的回答是“说实话，我学习的热情从来都没有低落过。”正如 Linus

在自己的自传《Just for Fun》中希望的那样。

本书的组织形式

本书将 Linux 内核的学习分为四个层次：全面了解，掌握基本功；兴趣导向，选择重点深度钻研；融入社区，参与开发做贡献；坚持，坚持，再坚持。总结起来，就是“全面了解抓基本，兴趣导向深钻研；融入社区做贡献，坚持坚持再坚持。”（如果您是一个修真小说爱好者，尽可以将其与炼气、筑基、结丹和元婴等层次相对应。）

第一层次修炼的内容包括了前三章，目的是希望您能够对 Linux 以及内核有个全面的认识 and 了解，掌握分析 Linux 内核源代码的分析方法。

第 1 章主要介绍了 Linux 的 18 年成长史，或许您会乐意陪我一起缅怀下这过去的十八年。

第 2 章介绍内核的配置和编译过程，和任何大型软件源码的学习一样，学会编译和配置是第一步。

第 3 章介绍学习内核需要的基础，内核的体系结构、目录结构、代码特点，浏览内核代码的工具，最后，突出强调了内核源码分析过程中极为重要的两个角色——Kconfig 和 Makefile，并以 USB 子系统为例，演示了如何利用这两个角色进行代码分析。

第二层次的修炼包括了第 4 ~ 11 章的内容，对内核多数部分的工作原理进行介绍。按照认识的发展规律，在第一层次修炼中已经对内核有个全局的认识和了解之后，接下来就应该以兴趣为导向，寻找一个子系统或模块，对其代码深入钻研和分析，不懂的地方就通过社区、邮件列表或者直接发 Email 给 maintainer 请教等途径弄懂，切勿得过且过，这样分析下来，对同步、中断等等内核的很多机制也同样会非常了解，俗话说一通则百通就是这个道

理。

因此第二层次的各个章节里，只是阐释重点的概念和工作原理，帮助您在分析该部分代码时进行理解，并不求详尽。

第 4 章讨论系统的初始化，万事开头难，系统的初始化是一个很复杂的过程，不过对于内核源码的学习来说，以这个部分开始应该是个不错的选择。特别是子系统初始化的讨论，应该是您选择任何内核子系统开始分析时都需要了解的内容。

第 5 章讨论系统调用，它是应用程序和内核间的桥梁，学习并理解它是我们走向内核的一个很好的过渡。

第 6 章讨论内核的中断处理机制，包括几乎任何一本内核书籍都没有涉及的通用 IRQ 层。

第 7 章讨论进程的内存抽象，以及进程如何被创建和销毁。如果我们将计算机上运行的操作系统以及各种各样的软件看作一系列有机的生命体，而不是死的指令集合，那么这就是一个进程的世界，只不过与我们人类世界不同的是，进程世界里的个体是一个一个鲜活的进程，而不是人。人的世界有道德与法律去制约管理，进程的世界同样也有自己的管理机制，这就是第 7 章所要展示的内容——进程管理。

第 8 章讨论进程的调度，重点讨论了在内核历史上具有重要地位的 O(1)调度器和最新的 CFS 调度器。

第 9 章讨论内存管理，内存就是进程的家，这里讨论内核如何为每个进程都分配一个家，并尽量的去做到“居者有其屋”，以及保证每个家的安全。

第 10 章讨论文件系统，主要是虚拟文件系统 (VFS)，它通过在各种具体的文件系统之

上建立一个抽象层，屏蔽了不同文件系统间的差异。

第 11 章讨论设备驱动，对于驱动开发来说，设备模型的理解是根本，spec、datasheet 与内核源代码的利用是关键。

通过第二层次的修炼，您应该对至少一到两个部分有了很深入的理解，对内核代码采用的通用手法也已经很拈熟，那么您应该开始进入第三层次，努力融入到内核的开发社区，此时的您已经不会再是社区中潜水小白的角色，而是会针对某个问题发表自己的见解。您已经可以尝试参与到内核的开发中去，即使仅仅修改了内核中的一个错误单词，翻译了一份大家需要的文档，也是做出了自己的贡献，会得到大家的认可。

本书中第三层次只包括了两章的内容，这是因为内核的修炼之道越往后便越依赖于自己，任何参考书都替代不了自己不断的反思与总结。

第 12 章讨论参与内核开发需要了解的一些基础信息。

第 13 章讨论了内核的调试技术，与第 12 章一样，您可以仅仅将这些内容看成内核修炼中的一些 tips。

至于最后的第四层次，更是仅有两个字——坚持。能够在内核的修炼之道上走多远，都取决于我们能够坚持多久，或许一个用一个公式概括更为合适：心态+兴趣+激情+时间+X=Y。

革命尚未成功，我等仍需努力。——与君共勉之。

精华版方法论部分导读

到目前为之，博客上分享的精华篇都可以归为方法论的范畴，在很多时候，都是方法论

要比细节紧要得多。而这些精华篇又可细分为三个专题：Linux 大史记；内核学习的方法论；驱动开发的方法论。

Linux大史记

除去那些精彩的“门”，我们生活中乏味的事情太多了，所以不希望再去按惯例花个一二页的篇幅乏味的写个“Linux 简介”，就将几天中出去溜弯的时间贡献了出来，逐年逐月的搜集整理了一些 Linux 成长过程中所发生的重要的事情，抑或一些非常有趣儿的事情。

开始时本以为这是一件很轻易的事，起码应该比统计公布房价上涨多少的事情轻易的多，利用 google，完成这么一件事情又有何难？但是意外的是，貌似很难找到类似的归纳整理，或许能够看到某个时间段内的所谓的 top10 之类的字眼，但里面的罗列似乎大都满足不了有趣儿的要求。所以里面有些月份是个空白，不管如何，大家可以了解了解，看看是否有很多自己不知道的有趣闻轶事？

[缅怀已逝的十八年 \(1991 ~ 1998 \)](#)

[缅怀已逝的十八年 \(1999 ~ 2002 \)](#)

[缅怀已逝的十八年 \(2003 ~ 2006 \)](#)

[缅怀已逝的十八年 \(2007 ~ 2009 \)](#)

内核学习的方法论

透过现象看本质，兽兽门无非就是一些人体艺术展示。同样往本质里看过去，学习内核，就是学习内核的源代码，任何内核有关的书籍都是基于内核，而又不高于内核的。

所以这个专题的前三个精华篇就是专注于介绍如何入手分析内核源代码的，这里前无来者的突出强调了“Kernel 地图”的概念，虽然 Goggle 带着 Goggle 地图远去了，可 Kernel

地图仍然在继续。

Kernel地图：Kconfig与Makefile

毫不夸张地说 ,Kconfig 和 Makefile 是我们浏览内核代码时最为依仗的两个文件。基本上 ,Linux 内核中每一个目录下边都会有一个 Kconfig 文件和一个 Makefile 文件。对于一个希望能够在 Linux 内核的汪洋代码里看到一丝曙光的人来说 ,将它们放在怎么重要的地位都不过分。

我们去香港 ,通过海关的时候 ,总会有免费的地图和各种指南拿 ,有了它们在手里我们才不至于无头苍蝇般迷惘的行走在陌生的街道上。即使在内地出去旅游的时候一般来说也总是会首先找份地图 ,当然了 ,这时就是要去买了 ,拿是拿不到的 ,不同的地方有不同的特色 , 只不过有的特色是服务 , 有的特色是索取。

Kconfig 和 Makefile 就是 Linux Kernel 迷宫里的地图。地图引导我们去认识一个城市 , 而 Kconfig 和 Makefile 则可以让我们了解一个 Kernel 目录下面的结构。我们每次浏览 kernel 寻找属于自己的那一段代码时 ,都应该首先看看目录下的这两个文件。

分析内核源码如何入手？（上）

既然要学习内核源码 ,就要经常对内核代码进行分析 ,而内核代码千千万 ,还前仆后继的不断往里加 ,这就让大部分人都有种雾里看花花不见的无助感。不过不要怕 ,孔老夫子早就留给我们了应对之策 :敏于事而慎于言 ,就有道而正焉 ,可谓好学也已。这就是说 ,做事要踏实才是好学生好同志 ,要遵循严谨的态度 ,去理解每一段代码的实现 ,多问多想多记。如果抱着走马观花 ,得过且过的态度 ,结果极有可能就是一边看一边丢 ,没有多大的收获。

[分析内核源码如何入手？（下）](#)

下面的分析，米卢教练说了，内容不重要，重要的是态度。就像韩局长对待日记的态度那样，严谨而细致。

只要你使用这样的态度开始分析内核，那么无论你选择内核的哪个部分作为切入点，比如 USB，比如进程管理，在花费相对不算很多的时间之后，你就会发现你对内核的理解会上升到另外一个高度，一个抱着情景分析，抱着 0.1 内核完全注释，抱着各种各样的内核书籍翻来覆去的看很多遍又忘很多遍都无法达到的高度。请相信我！

让我们在 Linux 社区里发出号召：学习内核源码，从学习韩局长开始！

对于学习来说，无论是在学校的课堂学习，还是这里说的内核学习，效果好或者坏，最主要取决于两个方面——方法论和心理。注意，我无视了智商的差异，这玩意儿玄之又玄，岔开了说，属于迷信的范畴。

因此继介绍分析内核源码方法的三个精华篇之后，又针对内核学习过程中最为常见的两个心理误区做了阐述。

[内核学习的心理问题](#)

而心理上的问题主要有两个，一个是盲目，就是在能够熟练适用 Linux 之前，对 Linux 为何物还说不出来个道道来，就迫不及待的盲目的去研究内核的源代码。这一部分人会觉得既然是学习内核，那么耗费时间在熟悉 Linux 的基本操作上纯粹是浪费宝贵的时间和感情。不过这样虽然很有韩峰同志的热情和干劲儿，但明显走入了一种心理误区。重述 Linus 的那句话：要先会使用它。

第二个就是恐惧。人类进化这么多年，面对复杂的物体和事情还是总会有天生的惧

怕感，体现在内核学习上面就是：那么庞大复杂的内核代码，让人面对起来该情何以堪啊！

即使有好的方法和坚强的心理，我们在内核学习过程中仍需要利用很多好的资源。其实，韩峰同志已经在日记里告诉了我们资源的重要性，因此我们在学习韩峰同志严谨细致的态度同时，还要领悟他对资源的灵活运用。只有在以内核源码为中心，坚持各种学习资源的长期建设不动摇，才能达到韩局长那样的高度，俯视 Linux 内核世界里的人生百态。

[内核学习的相关资源](#)

待到山花烂漫时，还是那些经典在微笑。

驱动开发的方法论

因为至少在国内大部分内核相关的开发都是驱动的开发，所以在内核学习的方法论之后，专门用一个专题，从模块机制、设备模型、驱动三件宝三个层次介绍了驱动开发的方法论。

[模块机制与 “Hello World!”](#)

有一种感动，叫泪流满面，有一种机制，叫模块机制。显然，这种模块机制给那些 Linux 的发烧友们带来了方便，因为模块机制意味着人们可以把庞大的 Linux 内核划分为许许多多小的模块。对于编写设备驱动程序的开发者来说，从此以后他们可以编写设备驱动程序却不需要把她编译进内核，不用 reboot 机器，她只是一个模块，当你需要她的时候，你可以把她抱入怀中（insmod），当你不再需要她的时候，你可以把她一脚踢开（rmmod）。

[设备模型（上）](#)

设备模型（下）

对于驱动开发来说，设备模型的理解是根本，毫不夸张得说，理解了设备模型，再去那些五花八门的驱动程序，你会发现自己站在了另一个高度，从而有了一种俯视的感觉，就像凤姐俯视知音和故事会，韩峰同志俯视女下属。

顾名思义就知道设备模型是关于设备的模型，既不是任小强们的房模，也不是张导的炮模。对咱们写驱动的和写不写驱动的人来说，设备的概念就是总线 and 与其相连的各种设备了。电脑城的 IT 工作者都会知道设备是通过总线连到计算机上的，而且还需要对应的驱动才能用，可是总线是如何发现设备的，设备又是如何和驱动对应起来的，它们经过怎样的艰辛才找到命里注定的那个他，它们的关系如何，白头偕老型的还是朝三暮四型的，这些问题就不是他们关心的了，而是咱们需要关心的。在房市股市千锤百炼的咱们还能够惊喜的发现，这些疑问的中心思想中心词汇就是总线、设备和驱动，没错，它们就是咱们这里要聊的 Linux 设备模型的名角。

驱动开发三件宝：spec、datasheet 与内核源码

设备模型之外，对于驱动程序的开发者来说，有三样东西是不可缺少的：第一是协议或标准的 spec，也就是规范，比如 usb 协议规范；第二是硬件的 datasheet，即你的驱动要支持的硬件的手册；第三就是内核里类似驱动的源代码，比如你要写触摸屏驱动的话，就可以参考内核里已经有的一些触摸屏驱动。

Linux内核问题门

继前面三个专题之后，为了感谢精华篇发布过程中很多朋友的关心与支持，便以“问题门”为题为拙作《Linux 内核修炼之道》制作了一个小插曲，希望通过对大家内核学习过程中遇到的问题与经验心得做一番展示，来帮助还在门外的朋友寻找到这扇门的钥匙。

陈宪章说：“学贵有疑，小疑则小进，大疑则大进。疑者，觉悟之机也，一番觉悟一番长进。”

培根说：“多问的人将多得。”

还在学校的时候导师在激情讲演之后对着会议室里形态各异但均静默不语的我们痛心疾首的说：“会提问题很重要啊，同志们！不会提问题怎么有资格做研究！”

这样铿锵有力的训诫今日想起仍觉深受刺激，于是就要不可避免得要做出一些反应来。不过一是因为咱这年代还没有非主流的说法，二是因为也没有冯仰妍同学的性别优势，不可能受到刺激就整出个门来。咱能够做到的最大反应也就是在这里开贴专门探讨探讨内核学习的相关问题，为了稍微增加那么一些广告效应，就称为“问题门”吧。

使用“问题门”的称呼，一是内心里潜藏的那点低级趣味想去沾点近些年层出不穷各种各样的“门”的仙气，二是在内核的学习过程中的确实实实在在的存在着这样的一个“门”，横亘在我们的面前，跨过去便海阔天空是另一番世界，但却是让无数人亮折腰，百思不得其钥匙。

缅怀已逝的十八年（1991～1998）

至此落笔之际，恰至 Linux 问世 18 周年，18 年的成长，风雨颇多，感慨颇多，谨以这些许年来的点滴之事为 Linux 的成人礼添彩。

如果你尚未与 Linux 亲密接触过，那么希望这里的内容可以成为你初识 Linux 的见证。

如果你已经是个 Linux 达人，那么就选个安静的早晨，抑或下午，陪我一起缅怀下这过去的

十八年吧。

Linux诞生记

1987 年

MINIX 诞生，而我也已端坐于学堂之中，隐去一身的稚气，能够摇头晃脑的吟诵几句诗赋了。若真是冥冥中自有定数的话，或许这时就暗定了 4 年后 Linux 的诞生。

1991 年

Linus Torvalds，一个芬兰的大学生，对于他不能按照意愿访问大学 UNIX 服务器而感到很愤怒，于是开始为一个以后被称为“Linux”的内核而工作，并于这一年的 10 月 5 日发布了 Linux 0.01。

1992 年

4 月，第一个 Linux 新闻组“comp.os.linux”建立。10 月，第一个可以安装的 Linux 版本 SLS 发布。同年，我拿到了平生的第一个毕业证。

1993 年

8 月，第一本关于 Linux 的著作《Linux Installation and Getting Started Version 1》出版。而这一年，我最敬佩的语文老师患病离去了，从此，我知道了生活中不仅仅只有欢聚，还有伤别。

1994 年

Linux 1.0 发布，并采用 GPL (GNU General Public License , 通用公共许可证) 协议。
大家要 Linus Torvalds 想一只吉祥物，Linus 突然想到小时候去动物园被一只企鹅追着满地

打滚，还被咬了一口！既然想不到其它的吉祥物了，干脆就以这支企鹅来当吉祥物算了！

泰坦尼克的狂潮

1995 年

4 月，召开首届 Linux 博览会，一个以 Linux 为特征的商业展览博览会。几个月后，我迎来了第二个中学阶段。

1996 年

Linux 2.0 发布，它第一个支持了 SMP (对称多处理器) 架构。此时 Linux 的全球用户已经达到了 350 万左右。

1997 年

首例 Linux 病毒 “Bliss” 被发现。电影《泰坦尼克号》所用的 160 台 Alpha 图形工作站中，有 105 台采用了 Linux。

1998 年

1 月，第一份 Linux 新闻周刊出版，同时，Netscape 宣布他们将在自由软件许可协议下发布浏览器的源代码，这为 Linux 和自由软件的发展提供了广阔空间。

2 月，Eric Raymond 和他的朋友门提出了 “open source” 的概念，申请了该商标特权并且组建了 opensource.org 网站，从而开始推动 Linux 的商业化发展。

4 月，Linux 广泛被美国国家公共新闻广播报道，标志 Linux 在主流、非技术性的媒体界首次出现。

5 月，Google 搜索引擎开始流行，不仅仅是因为它是最好的搜索引擎，而且还因为它

是基于 Linux 和具有 Linux 特色的搜索网页。

6 月，“从来没有一个用户向我提起 Linux，Linux 就像众多的免费产品一样，虽然它是很小的，却得到了一群忠诚的拥护者。” 比尔盖茨在 6 月 25 日的《PC 周刊》上说。

7 月，KDE 和 GNOME 的桌面之争在其拥护者之间愈演愈烈，Linux 以实际行动表明 KDE 非常好用，在这种情况下，KDE1.0 诞生了。Oracle、Informix、Sybase 都宣布将积极支持 Linux。Linux 开始成为一个家喻户晓的词。

9 月，Dave Whiting 和 Dwight Johnson 创建了 LinuxToday.com，该网站后来被 Internet.com 收购，不过它一直是访问量最高和最容易阅读的 Linux 入门网站。

12 月，一篇来自 IDC 的报导说 Linux 的发行量在 1998 年涨了 200%以上，它的市场占有率也增加了 150%以上。Linux 拥有 17%的市场占有率并且增长率超过了市场上其它任何一个系统。

同年，我迎来了人生中一个非常重要的时刻：我上大学了！

缅怀已逝的十八年（1999 ~ 2002）

提前发生的革命

1999 年

1 月，“Linux 2.2 已经发布，我终于可以松口气了” 创造者 Linus Torvalds 说。

3 月，首届 LinuxWorld 讨论会和博览会在加州的圣何塞举行，作为 Linux 第一个大的商业化的贸易展示活动，它无疑向世界昭示了 Linux 的到来。

8 月 ,SG 宣布了与 Red Hat 的合作关系 ,并且开始大规模的为内核的发展做贡献。Red Hat 进行了首次公开募股 ,股价马上涨到了 50 美元 ,在那个时候这个价似乎很高。摩托罗拉公司与 Lineo 建立了合作关系 ,进入 Linux 领域并提供嵌入式系统产品 ,支持和培训服务。Sun 宣布了 Sun 公共源许可 (Sun Community Source License) 下发行 StarOffice 和开发一个网络版本的办公套件。

9 月 , Red Hat 的股票达到了 135 美元 , 这个价格在那个时候似乎是难以置信的高。

10 月 , Sun 宣布它将在 Sun 公共源许可下公布 Solaris 的源代码。

12 月 , VA Linux Systems 的首次公开募股价格是 30 美元/股 , 这个价格很快涨到了 300 美元 , 它在 NASDAQ 历史上创造了最高的首次公开募股价格。

这一年 , 网络进入了宿舍 , QQ、mud 等也进入了我们的生活。

2000 年

1 月 , VA Linux Systems 宣布创建我们非常熟悉的 SourceForge , 到去年底 , SourceForge 已经接到了超过 12000 个项目 , 拥有 92000 个注册的开发者。

2 月 , 最近的 IDC 报告显示 Linux 现在成为 “服务器电脑上第二个最受欢迎的操作系统” , 在 1999 年占了 25% 的服务器操作系统销售额 , Windows NT 为 38% , 占第一位 , NetWare 为 19% , 排名第三 , IDC 以前曾预测过 Linux 将在 2002 或 2003 年到达第 2 位 , 这场革命提前发生了。

3 月 , 嵌入式 Linux 协会 (Embedded Linux Consortium) 成立。

8 月 , HP、Intel、IBM 以及 NEC 宣布开放源代码发展实验室 (OSDL , Open Source Development Lab) 成立。

9 月， Trolltech 发布了 GPL 下的 Qt 库。

11 月， IBM 宣布将在 2001 年投资 10 亿美元在 Linux。首部基于 Linux 的手机 IMT-2000 在韩国发布。

这一年的某一天，和同学坐在学校四大发明广场上观看同一首歌演出，困意盎然，期间那个粗犷的名歌星的一句话却惊醒了我：“希望你们交通大学为中国的交通事业做出更大的贡献”，大意如此，我顿时无语，他的语言竟然和他的外表一样粗犷。

这一年的暑假，我第一次来到江南，在西湖断桥对面的饭馆里，透过落地窗恰恰看到湖里荷花的位置，要了份西湖醋鱼和一瓶啤酒，坐到下午四点钟，然后顺着苏堤白堤静静的走下去，直到绕湖一周再次回到断桥，已是晚上八点，坐在湖边的长凳上，一夜无语。

和平、爱情和Linux

2001 年

1 月，期待已久的 Linux 2.4 发布。

3 月，Linux2.5 内核高级会议在加州圣何塞举行，它或许是历史上 Linux 内核 hacker 最完整的一次聚会。

4 月，IBM 在几个城市鼓吹“和平、爱情和 Linux”（Peace, Love and Linux）时遇到了麻烦。

6 月，Sharp 宣布基于 Lineo 嵌入式系统的 Linux PDA 即将上市。

这一年底，找工作的季节，我深刻认识了 IT 泡沫和 9.11，找所谓的好工作无门和出国

无门，我无奈选择考研。

2002 年

Linus Torvalds 将 Linux 2.4 交由巴西 18 岁的内核开发人员 Marcelo Tosatti 维护，自己则带领 Linux 2.5 的开发工作。

这一年，我从一个交大到了另一个交大，这个转变似乎很平淡，并不深刻。

缅怀已逝的十八年（2003 ~ 2006）

Ubuntu 4.10

2003 年

1 月，NEC 宣布将在其手机中使用 Linux，代表着 Linux 成功进军手机领域。

6 月，IDC 分析师称，2003 年 Linux 服务器在西欧的销售量将达到 18.2 万台，到 2007 年，销售量将增至这个数字的三倍，销售收入将翻一番，达到 19 亿美元。

8 月，韩国国家航空公司和 IBM 联合发布声明，表示韩国航空公司将把该公司的核心业务移植到 IBM 的 eServer 服务器当中完成，其中操作系统则采用 Linux。

9 月，三星在推出了首款基于 Linux 系统平台的 CDMA 智能手机 SCH-i519。

11 月，Linux 2.6 发布，它被认为是第一款真正意义上的企业级内核，这是 Linux 内核从 2001 年以来第一次的大改动。

这一年，我第一次在电视直播里看着自己喜欢的米兰夺得了冠军杯。

2004 年

1 月，X.Org 基金会成立。

2 月，Linux 标准 2.0 出台，规范了所有能被称为 Linux 操作系统所应该有的特性。

5 月，基于 Linux 的路由系统出现。

10 月 20 日，Ubuntu 首个版本发布，在五年后的今天 Ubuntu 已经是 Linux 桌面发行版的一个成功典范。

11 月，Firefox 1.0 发布，它成为大众关注的焦点，IE 降低了 1 个点的市场份额——像这种事已经多年没有发生过了。Firefox 已经成为了微软 IE 的强有力的对手。

又到了找工作的季节，宣讲会、笔试、面试，我就要离开学校了么？

2005 年

10 月，Firefox 的下载量突破了 1 亿大关，这表明，只要产品好，开放源代码软件也能够获得普通用户的青睐。

11 月，Sun 开放了除 Java 之外的几乎所有软件，这使得它在一夜间成为了最大的开放源码软件厂商之一。

12 月，Red Hat 公布了第三季度业报，销售收入增长了 43.6%，利润增长了 114%。

这一年夏天，遭遇了到目前为止最为严重的一次失窃，除了 IQ 卡，所有的卡都随着钱夹子消失了，到工行补办牡丹卡时，那慵懒的上海女人说，必须要上海土生土长的本地人来担保，仅仅拥有上海户口的人是不行的。

Richard Stallman的征婚启事

2006 年

6 月，自由软件之父Richard Stallman在自己的网站<http://www.stallman.org/> 上发布了一则 “征婚启事”。

I'm a single atheist white man, 52, reputedly intelligent, with unusual interests in politics, science, music and dance.

I'd like to meet a woman with varied interests, curious about the world, comfortable expressing her likes and dislikes (I hate struggling to guess), delighting in her ability to fascinate a man and in being loved tenderly, who values joy, truth, beauty and justice more than "success"--so we can share bouts of intense, passionately kind awareness of each other, alternating with tolerant warmth while we're absorbed in other aspects of life.

My 22-year-old child, the Free Software Movement, occupies most of my life, leaving no room for more children, but I still have room to love a sweetheart. I spend a lot of my time traveling to give speeches, often to Europe, Asia and Latin America; it would be nice if you were free to travel with me some of the time.

If you are interested, write to rms at stallman dot org and we'll see where it leads.

我，单身，无神论者，白人，52岁，据说比较聪明，对于政治、科学、音乐和舞蹈有着不同寻常的兴趣。

我想寻找这样一位女士：爱好广泛，对世界充满好奇心，能够清晰表达她的爱憎（我痛恨动脑筋猜测），乐于使男人着迷，渴望被温柔地爱，对于快乐、真理、美和正义的评价高于“成功”。这样的话，我们就能不断对另一方产生热烈而又美好的了解，当我们被生活中其他东西吸引的时候，彼此就能感到宽容的温暖。

我有一个22岁的孩子——自由软件运动——他占据了我大部分的生活，没有精力再抚养更多的孩子了，但是我仍然会投入的爱我的爱人。我有大量时间花在巡回演讲上，经常要去欧洲、亚洲和拉丁美洲。如果你有空在某些时间陪我一起旅行，那就最好了。

如果你有兴趣的话，请写信到 rms@stallman.org，让我们看看会有什么结果。

7月，Ubuntu 被授予 PC World 2006 World Class Award，证明了 Ubuntu 成为 2006 年世界最好的 100 个产品之一。Ubuntu 越来越显示出他的不凡实力，虽说他是免费的，但是后台却是商业公司 Canonical，加上太空人老板的聪明才智，逐渐的开始商业合作，比如和 Sun 合作，对有需要的客户提供 Linux 支持服务。

8月，Linux业界另外一位狂人，Linuspire公司总裁Kevin Carmony宣布推出免费版本的Freespire 1.0，该版本中附带有二进制的商业硬件驱动程序，在Linux社区中引起轩然大波。27日，网站http://linux.inet.hr/poll_filesystem.html上推出 “Your favorite file system?”（你最喜欢的文件系统？）投票活动。

9月，16日是“国际软件自由日”（SFD，Software Freedom Day 2006）。

10月，Oracle Unbreakable Linux 发布，Oracle 成为第一个推出自有 Linux 服务的

非操作系统软件厂商。17 日，FSG（自由标准组，一个非赢利的致力于开发和促进自由开放软件的标准的组织）宣布与 O'Reilly Media 合作，共同为 Linux 应用程序开发人员提供类似 MSDN 的服务，该服务将作为 LSB (Linux Standard Base) Developer Network 的一个组成部分。

11 月，微软和 Novell 达成一揽子协议，号称要改善 Linux 和微软操作系统的兼容问题。

如图 1.1 所示，看着昔日的对手用“+”连起来是否会觉得古怪？

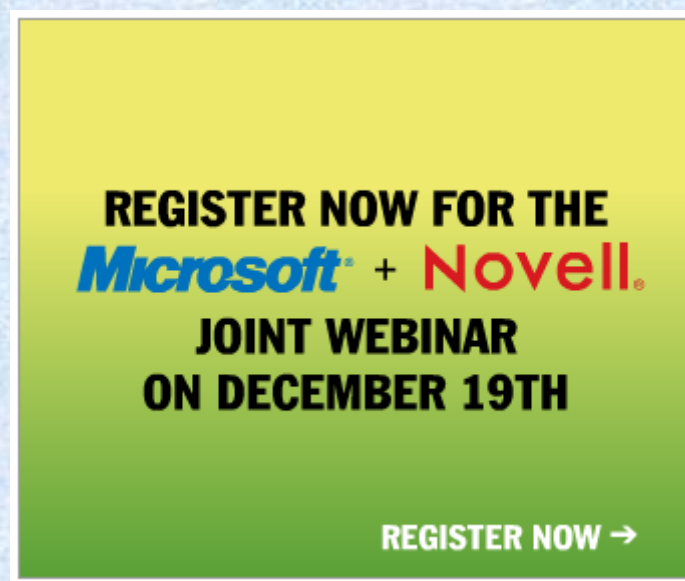


图 1.1 微软 + Novell

这一年，三次去青岛，回来时遭遇三次严重的飞机晚点，让我疑惑这个世界怎么了？

缅怀已逝的十八年（2007～2009）

来自微软的指控

2007 年

1 月，虚拟人生游戏（Second Life）客户端开源。两大 Linux 领导社团 OSDL 和 Free

Standard Group 宣布合并为新的 Linux Foundation (Linux 基金会), 此举将促进社区的资源整合, 也使 Linux 在企业市场能够更加高效地参与竞争。

2 月, Bill Xu 发起了一个 “致招商银行的公开信” 的行动, 希望用这种方式促使招商银行改变在公众服务中使用专属软件的作法, 取消客户端上的 ActiveX 技术, 而转用其他公开的、开放的、不限制用户平台的技术。据说, 浦发银行的网络银行能很好的支持 Firefox。

3 月, Novell 推出模仿苹果的 “Mac vs PC” 广告, 它在广告中插入了第三者 :Linux——一位迷人的年轻女子。Novell 用此来宣传预装 Novell Linux 的 PC, 一共发布了三个视频, 你可以在 www.youtube.com 上看到它们。15 日, Novell 公开表示, 同意从总费用上说 Linux 比 Windows 要昂贵, 这使它在开源社区的名誉进一步恶化。

4 月, Dell 推出预装 Ubuntu 操作系统笔记本。

5 月, 微软声称 Linux 内核侵犯了微软的 42 项专利, 而用户界面和其它设计方面也有 65 项侵权, OpenOffice.org 也被指控侵犯 45 项专利, 还有 83 项是针对其它免费开源软件。同一个月, 微软加拿大网站推出了一个 “Get the Facts” (了解真相) 页面, 如图 1.2 所示, 赤裸裸地对 Linux 进行了攻击, 有趣的是页面上方放置了一张《The Highly Reliable Times》报纸截图, 标题模仿《纽约时报》风格。“报纸” 中写道: “我们采用 Linux 平台以后每周至少遭遇一次系统崩溃问题。而迁移到微软 Windows Server 2003 后真正消灭了系统崩溃问题, 另外我们还能获得厂商支持。”

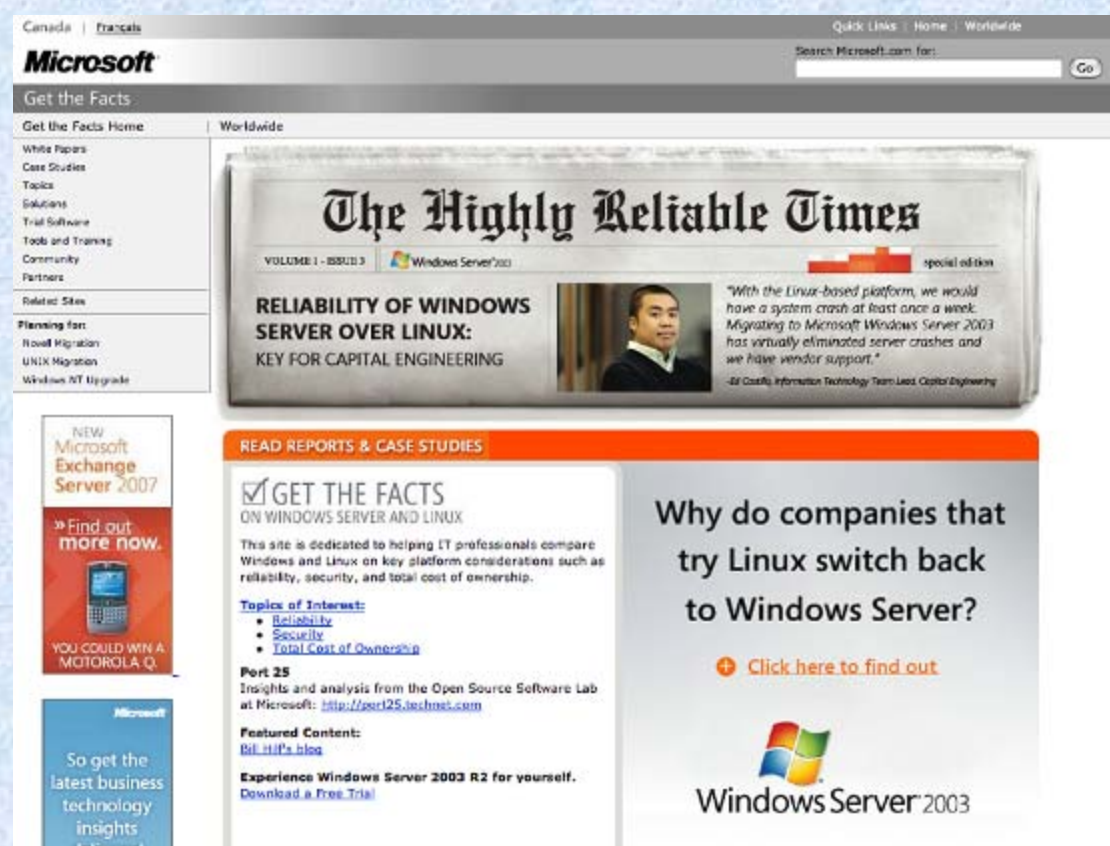


图 1.2 Get the Facts 页面

还是 5 月，Firefox 在 Linux 中显示的表单控件，特别是单选框，比较丑陋问题被修正，如图 1.3 所示。

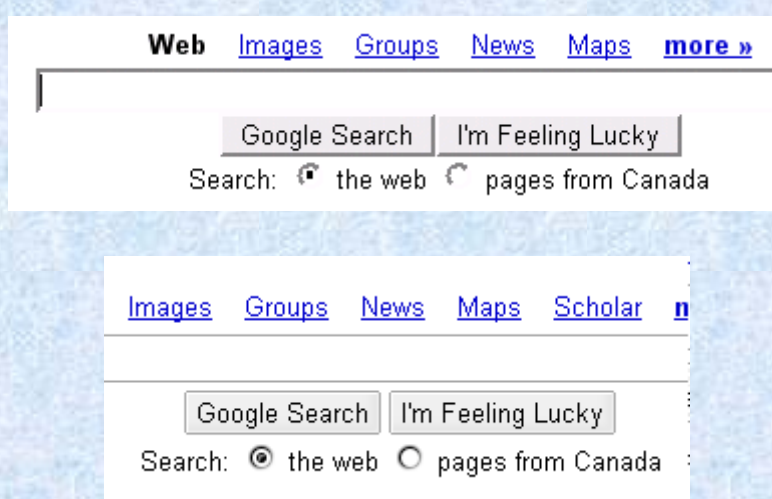


图 1.3 Firefox 表单控件修正前后比较

6月,5日微软和Linux发行商Xandros宣布,双方达成了一个技术和法律上的合作。

Red Hat、Ubuntu与Mandriva拒绝与微软进行专利交易。28日,Google桌面搜索Linux版正式发布,截图如图1.4所示。29日,第三版GNU通用公共许可证GPLv3发布。

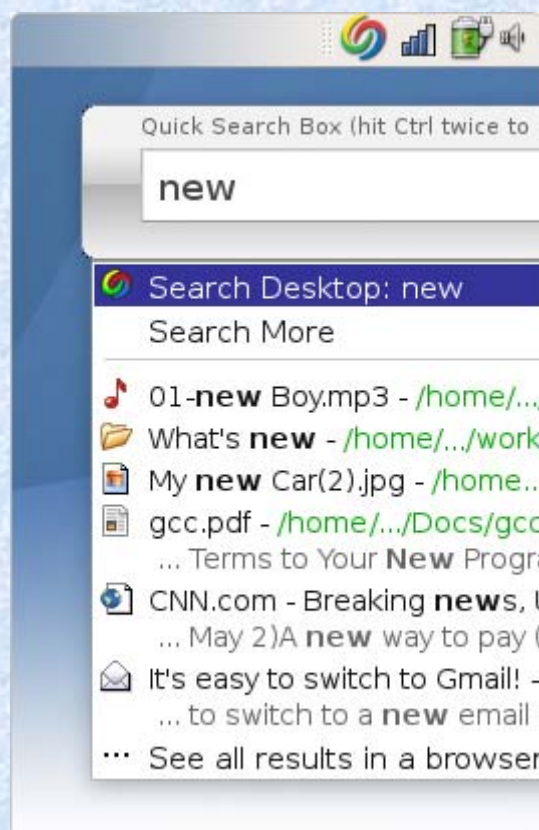


图 1.4 Google 桌面搜索 Linux 版

7月,Fcitx小企鹅输入法开源项目终止。做为Linux平台上最受欢迎的两大中文输入法之一——Fcitx小企鹅输入法,在其官方主页上宣布项目终止开发。声明中提到,有“编程高手”质疑其代码风格是项目终止的导火索。

8月,SCO在控告Linux侵犯专利权的官司中败诉,从而申请破产保护。SCO面市时以Linux销售商Caldera Systems的面目现身,然后从Santa Cruz Operation收购了Unix业务,之后重名为SCO集团。然后他们放弃了Linux业务,并开始起诉IBM、Novell及其他公司。他们认为IBM破坏了他们签署的Unix协议,将SCO特有的Unix技术在开源的

Linux 社区发布出去。Jim Zemlin 对此评论说：“如果它们把事业基础建立在协助 Linux，而不是攻击 Linux，那么它们大可享受像 RedHat 这些公司一样的成功，而不是沦落到申请破产保护的下场。”

10 月，Acacia Research 通过其子公司 IP Innovation 向 RedHat 和 Novell 提出控告，RedHat Linux 操作系统及 Novell 旗下的 SUSE Linux Enterprise Desktop 与 SUSE Linux Enterprise Server 侵犯了他所拥有的专利。随着围绕开源的纠纷不断，2007 年对于律师来说注定是“丰收”的一年。

11 月，Google 推出基于 Linux 的开源移动平台 Android。Phoronix 网站发布了 ATI 显卡在 Linux 和 Vista 下的游戏性能对比测试，结果令人鼓舞，在 Linux 下的游戏性能首次超越了 Windows！

这一年，我开始在 blog.csdn.net/fudan_abc 上写《Linux 那些事儿》。

首款Android手机

2008 年

1 月，Nokia 宣布收购了著名开源跨平台开发工具 Qt 的开发商 Trolltech。

2 月，Google 资助 Linux 版 Photoshop 的研究。

4 月，Sun 移除了 Java 最后的限制，将其彻底开源。

7 月，腾讯公司在这个月的最后一天发布了 QQ for Linux 1.0 Preview 版，这是第一次官方的版本。

9月,Google 联合 T-Mobile、HTC 正式发布了首款 Android 平台的手机 G1。Google 开源浏览器 Chrome 发布,发布仅仅几个小时,它的总体占有率就达到了 2%。

10月,OpenOffice3.0 发布,这对Linux的普及和实用化影响巨大。月底Ubuntu 8.10 发布,Fedora 10发布。

12月,各类发行版的 Linux 操作系统占据了大约三成的上网本市场份额。

对于我来说,这一年的基调是出差,大半年之后回到上海,很多地方很多事情都变得陌生起来,才发觉忘却其实也是一件很容易的事情。

Linux信用卡

2009年

1月, Linux 兼容内核正式使用项目 Unix 名称 Longene,中文别名“龙井”。兼容内核是一个自由、开源的操作系统项目,目的是要把 Linux 的内核扩充成一个既支持 Linux 应用、也支持 Windows 应用,既支持 Linux 设备驱动、也支持 Windows 设备驱动的兼容内核,使用户可以直接在 Linux 操作系统上高效运行 Windows 应用。

2月,微软起诉 GPS 设备制造商 Tomtom 侵犯其 8 项专利权, Tomtom 的 GPS 设备采用的是 Linux 系统,尽管微软声称 Linux 侵犯其专利期已有多年,但该案被视为微软状告 Linux 侵权的第一案。

3月, Adobe Reader 9.1 for Linux 发布。UltraEdit 正被移植到 Linux,名为 UEX,意即 UltraEdit for Linux。

4月, IDC 最新发表的题为《Linux 在新经济中的机会》的报告称,用户 2009 年的 Linux

开支预计将比 2008 年增长 21% ,超过整个软件市场的增长速度。整个软件市场 2009 年的增长率是 2%。

5 月 ,Nokia 宣布开放 Qt 源代码仓库 ,以便让社区的开发者能够进一步参与 Qt 的开发。

6 月 ,法国的 ENAC 开发组为 Linux 内核开发了类似 iPhone 的 Multi-touch (多点触摸) 技术支持。

7 月 ,Linux 基金会与 CardPartner 和 UMBrella 银行共同推出了 Visa 白金信用卡 ,正面印有 Linux 吉祥物 ,如图 1.5 所示。每办理一张这样的信用卡就可以为 Linux 基金会带来 50 美元的赞助 ,使用该卡每消费一次 Linux 基金会就能从中获得 1%的金额。



图 1.5 Linux 信用卡

8 月 ,微软在提交给美国证券交易委员会的年度文件中 ,将 Unbutu 列入竞争对手。

9 月 ,目前为止最新的内核版本 2.6.31 发布 ,Linux 成为首个正式支持 USB3.0 的操作

系统。

10 月，Ubuntu 9.10 发布。

11 月，Vim 的作者 Bram Moolenaar 推出了新的编程语言 Zimbu，一种不拐弯抹角直截了当的实验性编程语言。Moolenaar 表示 Zimbu 集现有语言的优点于一身，同时避开它们的不足。Zimbu 代码清晰易读，使用范围广泛——既能写 OS kernel，又能写脚本，还能写大的 GUI 程序，可以编译和运行在几乎所有系统上。

Kernel地图：Kconfig与Makefile

Makefile不是Make Love

从前在学校，混了四年，没有学到任何东西，每天就是逃课，上网，玩游戏，睡觉。毕业的时候，人家跟我说 Makefile 我完全不知，但是一说 Make Love 我就来劲了，现在想来依然觉得丢人。

毫不夸张地说，Kconfig 和 Makefile 是我们浏览内核代码时最为依仗的两个文件。基本上，Linux 内核中每一个目录下边都会有一个 Kconfig 文件和一个 Makefile 文件。对于一个希望能够在 Linux 内核的汪洋代码里看到一丝曙光的人来说，将它们放在怎么重要的地位都不过分。

我们去香港，通过海关的时候，总会有免费的地图和各种指南拿，有了它们在我们手里我们才不至于无头苍蝇般迷惘的行走在陌生的街道上。即使在内地出去旅游的时候一般来说也总是会首先找份地图，当然了，这时就是要去买，拿是拿不到的，不同的地方有不同的特色，只不过有的特色是服务，有的特色是索取。

Kconfig 和 Makefile 就是 Linux Kernel 迷宫里的地图。地图引导我们去认识一个城市，

而 Kconfig 和 Makefile 则可以让我们了解一个 Kernel 目录下面的结构。我们每次浏览 kernel 寻找属于自己的那一段代码时，都应该首先看看目录下的这两个文件。

利用Kconfig和Makefile寻找目标代码

就像利用地图寻找目的地一样，我们需要利用 Kconfig 和 Makefile 来寻找所要研究的目标代码。

比如我们打算研究 U 盘驱动的实现，因为 U 盘是一种 storage 设备，所以我们应该先进入到 drivers/usb/storage/目录。但是该目录下的文件很多，那么究竟哪些文件才是我们需要关注的？这时就有必要先去阅读 Kconfig 和 Makefile 文件。

对于 Kconfig 文件，我们可以看到下面的选项。

```
34 config USB_STORAGE_DATAFAB
35     bool "Datafab Compact Flash Reader support (EXPERIMENTAL)"
36     depends on USB_STORAGE && EXPERIMENTAL
37     help
38         Support for certain Datafab CompactFlash readers.
39         Datafab has a web page at <http://www.datafabusa.com/>.
```

显然，这个选项和我们的目的没有关系。首先它专门针对 Datafab 公司的产品，其次虽然 CompactFlash reader 是一种 flash 设备，但显然不是 U 盘。因为 drivers/usb/storage 目录下的代码是针对 usb mass storage 这一类设备，而不是针对某一种特定的设备。U 盘只是 usb mass storage 设备中的一种。再比如：

```
101 config USB_STORAGE_SDDR55
102     bool "SanDisk SDDR-55 SmartMedia support (EXPERIMENTAL)"
103     depends on USB_STORAGE && EXPERIMENTAL
104     help
105         Say Y here to include additional code to support the Sandisk SDDR-55
106         SmartMedia reader in the USB Mass Storage driver.
```

很显然这个选项是有关 SanDisk 产品的，并且针对的是 SM 卡，同样不是 U 盘，所以我们也不需要去关注。

事实上，很容易确定，只有选项 CONFIG_USB_STORAGE 才是我们真正需要关注的。

```
9 config USB_STORAGE
10 tristate "USB Mass Storage support"
11     depends on USB && SCSI
12 ---help---
13 Say Y here if you want to connect USB mass storage devices to your
14 computer's USB port. This is the driver you need for USB
15 floppy drives, USB hard disks, USB tape drives, USB CD-ROMs,
16 USB flash devices, and memory sticks, along with
17 similar devices. This driver may also be used for some cameras
18 and card readers.
19
20 This option depends on 'SCSI' support being enabled, but you
21 probably also need 'SCSI device support: SCSI disk support'
22 (BLK_DEV_SD) for most USB storage devices.
23
24 To compile this driver as a module, choose M here: the
25 module will be called usb-storage.
```

接下来阅读 Makefile 文件。

```
0 #
1 # Makefile for the USB Mass Storage device drivers.
2 #
3 # 15 Aug 2000, Christoph Hellwig <hch@infradead.org>
4 # Rewritten to use lists instead of if-statements.
5 #
6
7 EXTRA_CFLAGS := -Idrivers/scsi
8
9 obj-$(CONFIG_USB_STORAGE) += usb-storage.o
10
11 usb-storage-obj-$(CONFIG_USB_STORAGE_DEBUG) += debug.o
12 usb-storage-obj-$(CONFIG_USB_STORAGE_USBAT) += shuttle_usbat.o
13 usb-storage-obj-$(CONFIG_USB_STORAGE_SDDR09) += sddr09.o
14 usb-storage-obj-$(CONFIG_USB_STORAGE_SDDR55) += sddr55.o
15 usb-storage-obj-$(CONFIG_USB_STORAGE_FREECOM) += freecom.o
16 usb-storage-obj-$(CONFIG_USB_STORAGE_DPCM) += dpcm.o
17 usb-storage-obj-$(CONFIG_USB_STORAGE_ISD200) += isd200.o
18 usb-storage-obj-$(CONFIG_USB_STORAGE_DATAFAB) += datafab.o
19 usb-storage-obj-$(CONFIG_USB_STORAGE_JUMPSHOT) += jumpshot.o
20 usb-storage-obj-$(CONFIG_USB_STORAGE_ALAUDA) += alauda.o
21 usb-storage-obj-$(CONFIG_USB_STORAGE_ONETOUCH) += onetouch.o
22 usb-storage-obj-$(CONFIG_USB_STORAGE_KARMA) += karma.o
```

```
23
24 usb-storage-objs := scsiglue.o protocol.o transport.o usb.o \
25     initializers.o $(usb-storage-obj-y)
26
27 ifneq ($(CONFIG_USB_LIBUSUAL),)
28     obj-$(CONFIG_USB) += libusual.o
29 endif
```

前面通过 Kconfig 文件的分析,我们确定了只需要去关注 CONFIG_USB_STORAGE 选项。在 Makefile 文件里查找 CONFIG_USB_STORAGE,从第 9 行得知,该选项对应的模块为 usb-storage。

因为 Kconfig 文件里的其他选项我们都不需要关注,所以 Makefile 的 11~22 行可以忽略。第 24 行意味着我们只需要关注 scsiglue.c、protocol.c、transport.c、usb.c、initializers.c 以及它们同名的.h 头文件。

Kconfig 和 Makefile 很好的帮助我们定位到了所要关注的目标,就像我们到一个陌生的地方要随身携带地图,当我们学习 Linux 内核时,也要谨记寻求 Kconfig 和 Makefile 的帮助。

分析内核源码如何入手? (上)

透过现象看本质,兽兽门无非就是一些人体艺术展示。同样往本质里看过去,学习内核,就是学习内核的源代码,任何内核有关的书籍都是基于内核,而又不高于内核的。

既然要学习内核源码,就要经常对内核代码进行分析,而内核代码千千万,还前仆后继的不断往里加,这就让大部分人都有种雾里看花花不见的无助感。不过不要怕,孔老夫子早就留给我们了应对之策:敏于事而慎于言,就有道而正焉,可谓好学也已。这就是说,做事要踏实才是好学生好同志,要遵循严谨的态度,去理解每一段代码的实现,多问多想多记。

如果抱着走马观花，得过且过的态度，结果极有可能就是一边看一边丢，没有多大的收获。

假设全国房价上涨 1.5%，假设 80 后局长是农民子弟，……，既然我们的人生充满了假设，那么我在这里假设你现在就迫不及待的希望研究内核中 USB 子系统的实现，应该没有意见吧？那好，下面就以 USB 子系统的实现分析为标本看看分析内核源码应该如何入手。

分析README

内核中 USB 子系统的代码位于目录 `drivers/usb`，这个结论并不需要假设。于是我们进入到该目录，执行命令 `ls`，结果显示如下：

```
atm class core gadget host image misc mon serial storage Kconfig
Makefile README usb-skeleton.c
```

目录 `drivers/usb` 共包含有 10 个子目录和 4 个文件，`usb-skeleton.c` 是一个简单的 USB driver 的框架，感兴趣的可以去看看，目前来说，它还吸引不了我们的眼球。那么首先应该关注什么？如果迎面走来一个 `ppmm`，你会首先看脸、脚还是其它？当然答案依据每个人的癖好会有所不同。不过这里的问题应该只有一个答案，那就是 `Kconfig`、`Makefile`、`README`。

`README` 里有关于这个目录下内容的一般性描述，它不是关键，只是帮助你了解。再说了，面对“read 我吧 read 我吧”这么热情奔放的呼唤，善良的我们是不可能无动于衷的，所以先来看看里面都有些什么内容。

```
23 Here is a list of what each subdirectory here is, and what is contained in
24 them.
25
26 core/      - This is for the core USB host code, including the
27              usbfs files and the hub class driver ("khubd").
28
29 host/      - This is for USB host controller drivers. This
30              includes UHCI, OHCI, EHCI, and others that might
31              be used with more specialized "embedded" systems.
32
```



```

33 gadget/          - This is for USB peripheral controller drivers and
34                   the various gadget drivers which talk to them.
35
36
37 Individual USB driver directories.  A new driver should be added to the
38 first subdirectory in the list below that it fits into.
39
40 image/            - This is for still image drivers, like scanners or
41                   digital cameras.
42 input/            - This is for any driver that uses the input subsystem,
43                   like keyboard, mice, touchscreens, tablets, etc.
44 media/            - This is for multimedia drivers, like video cameras,
45                   radios, and any other drivers that talk to the v4l
46                   subsystem.
47 net/              - This is for network drivers.
48 serial/           - This is for USB to serial drivers.
49 storage/          - This is for USB mass-storage drivers.
50 class/            - This is for all USB device drivers that do not fit
51                   into any of the above categories, and work for a range
52                   of USB Class specified devices.
53 misc/             - This is for all USB device drivers that do not fit
54                   into any of the above categories.

```

这个 README 文件描述了前边使用 ls 命令列出的那 10 个文件夹的用途。那么什么是 USB Core ? Linux 内核开发者们，专门写了一些代码，负责实现一些核心的功能，为别的设备驱动程序提供服务，比如申请内存，比如实现一些所有的设备都会需要的公共的函数，并美其名曰 USB Core。

时代总在发展，当年胖杨贵妃照样迷死唐明皇，而如今人们欣赏的则是林志玲这样的魔鬼身材。同样，早期的 Linux 内核，其结构并不是如今天这般有层次感，远不像今天这般错落有致，那时候 drivers/usb/这个目录下边放了很多很多文件，USB Core 与其他各种设备的驱动程序的代码都堆砌在这里，后来，怎奈世间万千的变幻，总爱把有情的人分两端。于是在 drivers/usb/目录下面出来了一个 core 目录，就专门放一些核心的代码，比如初始化整个 USB 系统，初始化 Root Hub，初始化主机控制器的代码，再后来甚至把主机控制器相关的代码也单独建了一个目录，叫 host 目录，这是因为 USB 主机控制器随着时代的发展，

也开始有了好几种，不再像刚开始那样只有一种，所以呢，设计者们把一些主机控制器公共的代码仍然留在 core 目录下，而一些各主机控制器单独的代码则移到 host 目录下面让负责各种主机控制器的人去维护。

那么 USB gadget 那？gadget 说白了就是配件的意思，主要就是一些内部运行 Linux 的嵌入式设备，比如 PDA，设备本身有 USB 设备控制器（USB Device Controller），可以将 PC，也就是我们的主机作为 master 端，将这样的设备作为 slave 端和主机通过 USB 进行通信。从主机的观点来看，主机系统的 USB 驱动程序控制插入其中的 USB 设备，而 USB gadget 的驱动程序控制外围设备如何作为一个 USB 设备和主机通信。比如，我们的嵌入式板子上支持 SD 卡，如果我們希望在将板子通过 USB 连接到 PC 之后，这个 SD 卡被模拟成 U 盘，那么就要通过 USB gadget 架构的驱动。

剩下的几个目录分门别类的放了各种 USB 设备的驱动，比如 U 盘的驱动在 storage 目录下，触摸屏和 USB 键盘鼠标的驱动在 input 目录下，等等。

我们响应了 README 的热情呼唤，它便给予了我们想要的，通过它我们了解了 USB 目录里的那些文件夹都有着什么样的角色。到现在为止，就只剩下内核的地图——Kconfig 与 Makefile 两个文件了。有地图在手，对于在内核中游荡的我们来说，是件很愉悦的事情，不过，因为我们的目的是研究内核对 USB 子系统的实现，而不是特定设备或 host controller 的驱动，所以这里的定位很明显，USB Core 就是我们需要关注的对象，那么接下来就是要对 core 目录中的内容进行定位了。

分析Kconfig和Makefile

进入到 drivers/usb/core 目录，执行命令 ls，结果显示如下：

```
Kconfig  Makefile  buffer.c  config.c  devices.c  devio.c  driver.c
endpoint.c  file.c  generic.c  hcd-pci.c  hcd.c  hcd.h  hub.c  hub.h
```

```
inode.c message.c notify.c otg_whitelist.h quirks.c sysfs.c urb.c
usb.c usb.h
```

然后执行 wc 命令，如下所示。

```
# wc -l /*
  148 buffer.c
  607 config.c
  706 devices.c
1677 devio.c
1569 driver.c
  357 endpoint.c
  248 file.c
  238 generic.c
1759 hcd.c
  458 hcd.h
  433 hcd-pci.c
3046 hub.c
  195 hub.h
  758 inode.c
  144 Kconfig
   21 Makefile
1732 message.c
   68 notify.c
  112 otg_whitelist.h
  161 quirks.c
  710 sysfs.c
  589 urb.c
  984 usb.c
  160 usb.h
16880 total
```

drivers/usb/core 目录共包括 24 个文件，16880 行代码。core 不愧是 core，为大家默默的做这么多事。不过这么多文件里不一定都是我们所需要关注的，先拿咱们的地图来看看接下来该怎么走。先看看 Kconfig 文件，可以看到下面的选项。

```
15 config USB_DEVICEFS
16     bool "USB device filesystem"
17     depends on USB
18     ---help---
19     If you say Y here (and to "/proc file system support" in the "File
20     systems" section, above), you will get a file /proc/bus/usb/devices
21     which lists the devices currently connected to your USB bus or
```



```

22      busses, and for every connected device a file named
23      "/proc/bus/usb/xxx/yyy", where xxx is the bus number and yyy the
24      device number; the latter files can be used by user space programs
25      to talk directly to the device. These files are "virtual", meaning
26      they are generated on the fly and not stored on the hard drive.
27
28      You may need to mount the usbfs file system to see the files, use
29      mount -t usbfs none /proc/bus/usb
30
31      For the format of the various /proc/bus/usb/ files, please read
32      <file:Documentation/usb/proc_usb_info.txt>.
33
34      Usbfs files can't handle Access Control Lists (ACL), which are the
35      default way to grant access to USB devices for untrusted users of a
36      desktop system. The usbfs functionality is replaced by real
37      device-nodes managed by udev. These nodes live in /dev/bus/usb and
38      are used by libusb.

```

选项 USB_DEVICEFS 与 usbfs 文件系统有关。usbfs 文件系统挂载在 /proc/bus/usb 目录，显示了当前连接的所有 USB 设备及总线的各种信息，每个连接的 USB 设备在其中都会有一个对应的文件进行描述。比如文件 /proc/bus/usb/xxx/yyy，xxx 表示总线的序号，yyy 表示设备所在总线的地址。不过不能够依赖它们来稳定地访问设备，因为同一设备两次连接对应的描述文件可能会不同，比如，第一次连接一个设备时，它可能是 002/027，一段时间后再次连接，它可能就已经改变为 002/048。

就好比好不容易你暗恋的 mm 今天见你的时候对你抛了个媚眼，你心花怒放，赶快去买了 100 块彩票庆祝，到第二天再见到她的时候，她对你说你是谁啊，你悲痛欲绝的刮开那 100 块彩票，上面清一色的谢谢你。

因为 usbfs 文件系统并不属于 USB 子系统实现的核心部分，与之相关的代码我们可以不必关注。

```

74 config USB_SUSPEND
75     bool "USB selective suspend/resume and wakeup (EXPERIMENTAL)"
76     depends on USB && PM && EXPERIMENTAL
77     help

```

```

78     If you say Y here, you can use driver calls or the sysfs
79     "power/state" file to suspend or resume individual USB
80     peripherals.
81
82     Also, USB "remote wakeup" signaling is supported, whereby some
83     USB devices (like keyboards and network adapters) can wake up
84     their parent hub. That wakeup cascades up the USB tree, and
85     could wake the system from states like suspend-to-RAM.
86
87     If you are unsure about this, say N here.

```

这一项是有关 USB 设备的挂起和恢复。开发 USB 的人都是节电节能的好孩子，所以协议里就规定了，所有的设备都必须支持挂起状态，就是说为了达到节电的目的，当设备在指定的时间内，如果没有发生总线传输，就要进入挂起状态。当它收到一个 non-idle 的信号时，就会被唤醒。节约用电从 USB 做起。不过这个与主题也没太大关系，相关代码也可以不用关注了。

剩下的还有几项，不过似乎与咱们关系也不大，还是去看看 Makefile。

```

5 usbcore-objs := usb.o hub.o hcd.o urb.o message.o driver.o \
6               config.o file.o buffer.o sysfs.o endpoint.o \
7               devio.o notify.o generic.o quirks.o
8
9 ifeq ($(CONFIG_PCI),y)
10     usbcore-objs += hcd-pci.o
11 endif
12
13 ifeq ($(CONFIG_USB_DEVICEFS),y)
14     usbcore-objs += inode.o devices.o
15 endif
16
17 obj-$(CONFIG_USB) += usbcore.o
18
19 ifeq ($(CONFIG_USB_DEBUG),y)
20 EXTRA_CFLAGS += -DDEBUG
21 endif

```

Makefile可比Kconfig简略多了，所以看起来也更亲切点，咱们总是拿的money越多越好，看的代码越少越好。这里之所以会出现CONFIG_PCI，是因为通常USB的Root Hub包

含在一个PCI设备中。hcd-pci和hcd顾名思义就知道是说主机控制器的，它们实现了主机控制器公共部分，按协议里的说法它们就是HCDI（HCD的公共接口），host目录下则实现了各种不同的主机控制器。CONFIG_USB_DEVICEFS前面的Kconfig文件里也见到了，关于usbfs的，与咱们的主题无关，inode.c和devices.c两个文件也可以不用管了。

那么我们可以得出结论，为了理解内核对USB子系统的实现，我们需要研究buffer.c、config.c、driver.c、endpoint.c、file.c、generic.c、hcd.c、hcd.h、hub.c、message.c、notify.c、otg_whitelist.h、quirks.c、sysfs.c、urb.c 和usb.c文件。这么看来，好像大都需要关注的样子，没有减轻多少压力，不过这里本身就是USB Core部分，是要做很多的事为咱们分忧的，所以多点也是可以理解的。

分析内核源码如何入手？（下）

下面的分析，米卢教练说了，内容不重要，重要的是态度。就像韩局长对待日记的态度那样，严谨而细致。

只要你使用这样的态度开始分析内核，那么无论你选择内核的哪个部分作为切入点，比如USB，比如进程管理，在花费相对不算很多的时间之后，你就会发现你对内核的理解会上升到另外一个高度，一个抱着情景分析，抱着0.1内核完全注释，抱着各种各样的内核书籍翻来覆去的看很多遍又忘很多遍都无法达到的高度。请相信我！

让我们在Linux社区里发出号召：学习内核源码，从学习韩局长开始！

态度决定一切：从初始化函数开始

任小强们说房价高涨从现在开始，股评家们说牛市从5000点开始。他们的开始需要我

们的钱袋，我们的开始只需要一台电脑，最好再有一杯茶，伴着几支小曲儿，不盯着钱总是会比较惬意的。生容易，活容易，生活不容易，因为总要盯着钱。

有了地图 Kconfig 和 Makefile，我们可以在庞大复杂的内核代码中定位以及缩小了目标代码的范围。那么现在，为了研究内核对 USB 子系统的实现，我们还需要在目标代码中找到一个突破口，这个突破口就是 USB 子系统的初始化代码。

针对某个子系统或某个驱动，内核使用 `subsys_initcall` 或 `module_init` 宏指定初始化函数。在 `drivers/usb/core/usb.c` 文件中，我们可以发现下面的代码。

```
940 subsys_initcall(usb_init);
941 module_exit(usb_exit);
```

我们看到一个 `subsys_initcall`，它也是一个宏，我们可以把它理解为 `module_init`，只不过因为这部分代码比较核心，开发者们把它看作一个子系统，而不仅仅是一个模块。这也很好理解，`usbcore` 这个模块它代表的不是某一个设备，而是所有 USB 设备赖以生存的模块，Linux 中，像这样一个类别的设备驱动被归结为一个子系统。比如 PCI 子系统，比如 SCSI 子系统，基本上，`drivers/` 目录下面第一层的每个目录都算一个子系统，因为它们代表了一类设备。

`subsys_initcall(usb_init)` 的意思就是告诉我们 `usb_init` 是 USB 子系统真正的初始化函数，而 `usb_exit()` 将是整个 USB 子系统的结束时的清理函数。于是为了研究 USB 子系统在内核中的实现，我们需要从 `usb_init` 函数开始看起。

```
865 static int __init usb_init(void)
866 {
867     int retval;
868     if (!usb) {
869         pr_info("%s: USB support disabled\n", usbcore_name);
870         return 0;
871     }
872 }
```

```

873     retval = ksuspend_usb_init();
874     if (retval)
875         goto out;
876     retval = bus_register(&usb_bus_type);
877     if (retval)
878         goto bus_register_failed;
879     retval = usb_host_init();
880     if (retval)
881         goto host_init_failed;
882     retval = usb_major_init();
883     if (retval)
884         goto major_init_failed;
885     retval = usb_register(&usbfs_driver);
886     if (retval)
887         goto driver_register_failed;
888     retval = usb_devio_init();
889     if (retval)
890         goto usb_devio_init_failed;
891     retval = usbfs_init();
892     if (retval)
893         goto fs_init_failed;
894     retval = usb_hub_init();
895     if (retval)
896         goto hub_init_failed;
897     retval = usb_register_device_driver(&usb_generic_driver, THIS_MODULE);
898     if (!retval)
899         goto out;
900
901     usb_hub_cleanup();
902 hub_init_failed:
903     usbfs_cleanup();
904 fs_init_failed:
905     usb_devio_cleanup();
906 usb_devio_init_failed:
907     usb_deregister(&usbfs_driver);
908 driver_register_failed:
909     usb_major_cleanup();
910 major_init_failed:
911     usb_host_cleanup();
912 host_init_failed:
913     bus_unregister(&usb_bus_type);
914 bus_register_failed:
915     ksuspend_usb_cleanup();
916 out:

```

```
917     return retval;
918 }
```

(1) __init 标记。

关于 usb_init，第一个问题是，第 865 行的 __init 标记具有什么意义？

写过驱动的应该不会陌生，它对内核来说就是一种暗示，表明这个函数仅在初始化期间使用，在模块被装载之后，它占用的资源就会释放掉用作它处。它的暗示你懂，可你的暗示，她却不懂或者懂装不懂，多么让人感伤。它在自己短暂的一生中一直从事繁重的工作，吃的是草吐出的是牛奶，留下的是整个 USB 子系统的繁荣。

受这种精神所感染，我觉得有必要为它说的更多些。__init 的定义在 include/linux/init.h 文件里

```
43 #define __init      __attribute__((__section__(".init.text")))
```

好像这里引出了更多的疑问，__attribute__ 是什么？Linux 内核代码使用了大量的 GNU C 扩展，以至于 GNU C 成为能够编译内核的唯一编译器，GNU C 的这些扩展对代码优化、目标代码布局、安全检查等方面也提供了很强的支持。而 __attribute__ 就是这些扩展中的一个，它主要被用来声明一些特殊的属性，这些属性主要被用来指示编译器进行特定方面的优化和更仔细的代码检查。GNU C 支持十几个属性，section 是其中的一个，我们查看 GCC 的手册可以看到下面的描述

```
'section ("section-name")'
Normally, the compiler places the code it generates in the `text'
section. Sometimes, however, you need additional sections, or you
need certain particular functions to appear in special sections.
The `section' attribute specifies that a function lives in a
particular section. For example, the declaration:

extern void foobar (void) __attribute__((section ("bar")));
```



```
puts the function 'foobar' in the 'bar' section.
```

Some file formats do not support arbitrary sections so the 'section' attribute is not available on all platforms. If you need to map the entire contents of a module to a particular section, consider using the facilities of the linker instead.

通常编译器将函数放在.text 节，变量放在.data 或.bss 节，使用 section 属性，可以让编译器将函数或变量放在指定的节中。那么前面对__init 的定义便表示将它修饰的代码放在.init.text 节。连接器可以把相同节的代码或数据安排在一起，比如__init 修饰的所有代码都会被放在.init.text 节里，初始化结束后就可以释放这部分内存。

问题可以到此为止，也可以更深入，即内核又是如何调用到这些__init修饰的初始化函数？要回答这个问题，还需要回顾一下subsys_initcall宏，它也在include/linux/init.h里定义

```
125 #define subsys_initcall(fn)          __define_initcall("4",fn,4)
```

这里又出现了一个宏__define_initcall，它用于将指定的函数指针fn放到initcall.init节里 而对于具体的subsys_initcall宏，则是把fn放到.initcall.init的子节.initcall4.init里。要弄清楚.initcall.init、.init.text和.initcall4.init这样的东东，我们还需要了解一点内核可执行文件相关的概念。

内核可执行文件由许多链接在一起的对象文件组成。对象文件有许多节，如文本、数据、init 数据、bass 等等。这些对象文件都是由一个称为链接器脚本的文件链接并装入的。这个链接器脚本的功能是将输入对象文件的各节映射到输出文件中；换句话说，它将所有输入对象文件都链接到单一的可执行文件中，将该可执行文件的各节装入到指定地址处。vmlinux.lds 是存在于 arch/<target>/ 目录中的内核链接器脚本，它负责链接内核的各个节并将它们装入内存中特定偏移量处。

我可以负责的告诉你，要看懂 vmlinux.lds 这个文件是需要一番功夫的，不过大家都是聪明人，聪明人做聪明事，所以你需要做的只是搜索 `initcall.init`，然后便会看到似曾相识的内容

```
__initcall_start = .;
.initcall.init : AT(ADDR(.initcall.init) - 0xC0000000) {
    *(.initcall1.init)
    *(.initcall2.init)
    *(.initcall3.init)
    *(.initcall4.init)
    *(.initcall5.init)
    *(.initcall6.init)
    *(.initcall7.init)
}
__initcall_end = .;
```

这里的 `__initcall_start` 指向 `.initcall.init` 节的开始，`__initcall_end` 指向它的结尾。

而 `.initcall.init` 节又被分为了 7 个子节，分别是

```
.initcall1.init
.initcall2.init
.initcall3.init
.initcall4.init
.initcall5.init
.initcall6.init
.initcall7.init
```

我们的 `subsys_initcall` 宏便是将指定的函数指针放在了 `.initcall4.init` 子节。其它的比如 `core_initcall` 将函数指针放在 `.initcall1.init` 子节，`device_initcall` 将函数指针放在了 `.initcall6.init` 子节等等，都可以从 `include/linux/init.h` 文件找到它们的定义。各个字节的顺序是确定的，即先调用 `.initcall1.init` 中的函数指针再调用 `.initcall2.init` 中的函数指针，等等。`__init` 修饰的初始化函数在内核初始化过程中调用的顺序和 `.initcall.init` 节里函数指针的顺序有关，不同的初始化函数被放在不同的子节中，因此也就决定了它们的调用顺序。

至于实际执行函数调用的地方，就在 `/init/main.c` 文件里，内核的初始化么，不在那里

还能在哪里，里面的 `do_initcalls` 函数会直接用到这里的 `__initcall_start`、`__initcall_end` 来进行判断。

(2) 模块参数。

关于 `usb_init` 函数，第二个问题是，第 868 行的 `nousb` 表示什么？

知道 C 语言的人都会知道 `nousb` 是一个标志，只是不同的标志有不一样的精彩，这里的 `nousb` 是用来让我们在启动内核的时候通过内核参数去掉 USB 子系统的，Linux 社会是一个很人性化的世界，它不会去逼迫我们接受 USB，一切都只关乎我们自己的需要。不过我想我们一般来说是不会去指定 `nousb` 的吧。如果你真的指定了 `nousb`，那它就只会幽怨的说一句 “USB support disabled”，然后退出 `usb_init`。

`nousb` 在 `drivers/usb/core/usb.c` 文件中定义为：

```
static int nousb; /* Disable USB when built into kernel image */
module_param_named(autosuspend, usb_autosuspend_delay, int, 0644);
MODULE_PARM_DESC(autosuspend, "default autosuspend delay");
```

从中可知 `nousb` 是个模块参数。关于模块参数，我们都知道可以在加载模块的时候可以指定，但是如何在内核启动的时候指定？

打开系统的 `grub` 文件，然后找到 `kernel` 行，比如：

```
kernel /boot/vmlinuz-2.6.18-kdb root=/dev/sda1 ro splash=silent vga=0x314
```

其中的 `root`，`splash`，`vga` 等都表示内核参数。当某一模块被编译进内核的时候，它的模块参数便需要在 `kernel` 行来指定，格式为 “模块名.参数=值”，比如：

```
modprobe usbcore autosuspend=2
```

对应到 `kernel` 行，即为：

```
usbcore.autosuspend=2
```


通过命令 “modinfo -p \${modulename}” 可以得知一个模块有哪些参数可以使用。

同时，对于已经加载到内核里的模块，它们的模块参数会列举在 /sys/module/\${modulename}/parameters/ 目录下面，可以使用 “echo -n \${value} > /sys/module/\${modulename}/parameters/\${parm}” 这样的命令去修改。

(3) 可变参数宏。

关于 usb_init 函数，第三个问题是，pr_info 如何实现与使用？

pr_info 只是一个打印信息的可辨参数宏，printf 的变体，在 include/linux/kernel.h 里定义：

```
242 #define pr_info(fmt,arg...) \
243     printk(KERN_INFO fmt,##arg)
```

99 年的 ISO C 标准里规定了可变参数宏，和函数语法类似，比如

```
#define debug(format, ...) fprintf (stderr, format, __VA_ARGS__)
```

里面的 “...” 就表示可变参数，调用时，它们就会替代宏体里的 __VA_ARGS__。GCC 总是会显得特立独行一些，它支持更复杂的形式，可以给可变参数取个名字，比如

```
#define debug(format, args...) fprintf (stderr, format, args)
```

有了名字总是会容易交流一些。是不是与 pr_info 比较接近了？除了 ‘##’，它主要是针对空参数的情况。既然说是可变参数，那传递空参数也总是可以的，空即是多，多即是空，股市里的哲理这里同样也是适合的。如果没有 ‘##’，传递空参数的时候，比如

```
debug ("A message");
```

展开后，里面的字符串后面会多个多余的逗号。这个逗号你应该不会喜欢，而 ‘##’ 则会使预处理器去掉这个多余的逗号。

关于 usb_init 函数，上面的三个问题之外，余下的代码分别完成 usb 各部分的初始化，

接下来就需要围绕它们分别进行深入分析。因为这里只是演示如何入手分析，展示的只是一种态度，所以具体的深入分析就免了吧。

内核学习的心理问题

对于学习来说，无论是在学校的课堂学习，还是这里说的内核学习，效果好或者坏，最主要取决于两个方面——方法论和心理。注意，我无视了智商的差异，这玩意儿玄之又玄，岔开了说，属于迷信的范畴。

前面又是 Kernel 地图，又是如何入手，说的都是方法论的问题，那么这里要面对的就主要是心理上的问题。

而心理上的问题主要有两个，一个是盲目，就是在能够熟练适用 Linux 之前，对 Linux 为何物还说不出来个道道来，就迫不及待的盲目的去研究内核的源代码。这一部分人会觉得既然是学习内核，那么耗费时间在熟悉 Linux 的基本操作上纯粹是浪费宝贵的时间和感情。不过这样虽然很有韩峰同志的热情和干劲儿，但明显走入了一种心理误区。重述 Linus 的那句话：要先会使用它。

第二个就是恐惧。人类进化这么多年，面对复杂的物体和事情还是总会有天生的惧怕感，体现在内核学习上面就是：那么庞大复杂的内核代码，让人面对起来该情何以堪啊！

有了这种恐惧无力感存在，心理上就会去排斥面对接触内核源码，宁愿去抱着情景分析，搜集各种各样五花八门的内核书籍放在那里屯着，看了又忘，忘了又看，也不大情愿去认真细致得浏览源码。

这个时候，我们在心理上是脆弱得，我们忘记了芙蓉姐姐，工行女之所以红起来，不是

她们有多好，而是因为她们得心理足够坚强。是的，除了向韩局长学习态度，我们还要向涌现出来的无数个芙蓉姐姐和工行女学习坚强的心理。

有必要再强调一次，学习内核，就是学习内核的源代码，任何内核有关的书籍都是基于内核，而又不高于内核的。内核源码本身就是最好的参考资料，其他任何经典或非经典的书最多只是起到个辅助作用，不能也不应该取代内核代码在我们学习过程中的主导地位。

内核学习的相关资源

“世界上最缺的不是金钱，而是资源。”当我在一份报纸上看到这句大大标题时，我的第一反应是——作者一定是个自然环保主义者，然后我在羞愧得反省自身的同时油然而生出一股对这样的无产主义理想者无比崇敬的情绪来。

于是，我继续往下看，“因此在 XXX 还未正式面市之时，前来咨询的客户已经不少，这些有眼光的购房者明白，谁能在目前最好的购房机会下最大化地占有绝版资源，谁就掌控了未来财富流向。”（为了避免做广告的嫌疑，请允许我使用 XXX 代替该楼盘的名字。）顿时，我悟道了！

其实，韩峰同志已经在日记里告诉了我们资源的重要性，因此我们在学习韩峰同志严谨细致的态度同时，还要领悟他对资源的灵活运用。只有在以内核源码为中心，坚持各种学习资源的长期建设不动摇，才能达到韩局长那样的高度，俯视 Linux 内核世界里的人生百态。

注意，这个观点与前面所说的学习效果主要取决于方法论和心理两个方面并不矛盾，它们属于不同层次上的问题。

内核文档

内核代码中包含有大量的文档,这些文档对于学习理解内核有着不可估量的价值,记住,在任何时候,它们在我们心目中的地位都应该高于那些各式的内核参考书。下面是一些内核新人所应该阅读的文档。

README

这个文件首先简单介绍了 Linux 内核的背景,然后描述了如何配置和编译内核,最后还告诉我们出现问题时应该怎么办。

Documentation/Changes

这个文件给出了用来编译和使用内核所需要的最小软件包列表。

Documentation/CodingStyle

这个文件描述了内核首选的编码风格,所有代码都应该遵守里面定义的规范。

Documentation/SubmittingPatches

Documentation/SubmittingDrivers

Documentation/SubmitChecklist

这三个文件都是描述如何提交代码的,其中 SubmittingPatches 给出创建和提交补丁的过程,SubmittingDrivers 描述了如何将设备驱动提交给 2.4、2.6 等不同版本的内核树,SubmitChecklist 则描述了提交代码之前需要 check 自己的代码应该遵守的某些事项。

Documentation/stable_api_nonsense.txt

这个文件解释了为什么内核没有一个稳定的内部 API 到用户空间的接口——系统调用

——是稳定的), 它对于理解 Linux 的开发哲学至关重要, 对于将开发平台从其他操作系统转移到 Linux 的开发者来说也很重要。

Documentation/stable_kernel_rules.txt

解释了稳定版内核 (stable releases) 发布的规则, 以及如何将补丁提交给这些版本。

Documentation/SecurityBugs

内核开发者对安全性问题非常关注, 如果你认为自己发现了这样的问题, 可以根据这个文件中给出的联系方式提交 bug, 以便能够尽可能快的解决这个问题。

Documentation/kernel-docs.txt

这个文件列举了很多内核相关的文档和书籍, 里面不乏经典之作。

Documentation/applying-patches.txt

这个文件回答了如何为内核打补丁。

Documentation/bug-hunting

这个文件是有关寻找、提交、修正 bug 的。

Documentation/HOWTO

这个文件将指导你如何成为一名内核开发者, 并且学会如何同内核开发社区合作。它尽可能不包括任何关于内核编程的技术细节, 但会给你指引一条获得这些知识的正确途径。

经典书籍

待到山花烂漫时, 还是那些经典在微笑。

有关内核的书籍可以用汗牛充栋来形容，不过只有一些经典的神作经住了考验。首先是 5 本久经考验的神作（个人概括为“2+1+2”，第一个 2 是指 2 本全面讲解内核的书，中间的 1 指 1 本讲解驱动开发的书，后面的 2 则指 2 本有关内核具体子系统的书，你是否想到了某某广告里三个人突然站起单臂齐举高呼“1 比 1 比 1”的场景？）。

《Linux 内核设计与实现》

简称 LKD，从入门开始，介绍了诸如进程管理、系统调用、中断和中断处理程序、内核同步、时间管理、内存管理、地址空间、调试技术等方面，内容比较浅显易懂，个人认为是内核新人首先必读的书籍。新人得有此书，足矣！

《深入理解 Linux 内核》

简称 ULK，相比于 LKD 的内容不够深入、覆盖面不广，ULK 要深入全面得多。

前面这两本，一本提纲挈领，一本全面深入。

《Linux 设备驱动程序》

简称 LDD，驱动开发者都要人手一本了。

《深入理解 Linux 虚拟内存管理》

简称 LVMM，是一本介绍 Linux 虚拟内存管理机制的书。如果你希望深入的研究 Linux 的内存管理子系统，仔细的研读这本书无疑是最好的选择。

《深入理解 LINUX 网络内幕》

一本讲解网络子系统实现的书，通过这本书，我们可以了解到 Linux 内核是如何实现复杂的网络功能的。

这 5 本书各有侧重，正如下面的图所展示的那样，恰好代表了个人一直主张的内核学习方法：首先通过 LKD 或 ULK 了解内核的设计实现特点，对内核有个整体全局的认识和理解，然后可分为两个岔路，如果从事驱动开发，则钻研 LDD，如果希望对内核不是泛泛而谈而是有更深入的理解，则可以选择一个自己感兴趣的子系统，仔细分析它的代码，不懂的地方就通过社区、邮件列表或者直接发 Email 给 maintainer 请教等途径弄懂，切勿得过且过，这样分析下来，对同步、中断等等内核的很多机制也同样会非常了解，俗话说的一通则百通就是这个道理。当然，如果你选择研究的是内存管理或者网络，则可以有上面的两本书可以学习，如果是其他子系统，可能就没有这么好的运气了。



内核社区

最近几年，社区网站非常的热火，不过此社区非彼社区。

Linux最大的一个优势就是它有一个紧密团结了众多使用者和开发者的社区，它的目标就是提供尽善尽美的内核。内核社区的中心是内核邮件列表（Linux Kernel Mailing List，LKML），我们可以在<http://vger.kernel.org/vger-lists.html#linux-kernel>上面看到订阅这个邮件列表的细节。

内核邮件列表的流量很大，每天都有几百条消息，这里是大牛们的战场，小牛们的天堂，任何一个内核开发者都可以从中受益非浅。

除了LKML，大多数子系统也有自己独立的邮件列表来协调各自的开发工作，比如USB子系统的邮件列表可以在<http://www.linux-usb.org/mailling.html>上面订阅。

其他网络资源

除了内核邮件列表，还有很多其他的论坛或网站值得我们经常关注。我们要知道，网络上不仅有兽兽和凤姐，也不仅有犀利哥和韩局长。

<http://www.kernel.org/>

可以通过这个网站上下载内核的源代码和补丁、跟踪内核bug等。

<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，创建于1997年底的一个Linux新闻站点。

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

这是内核开发的中文邮件列表，里面活跃着很多内核开发领域的华人，比如Herbert Xu、Mingming Cao、Bryan Wu等。

<http://linux.chinaunix.net/>

全球最大的Linux/Unix中文技术社区。

.....

模块机制与 “Hello World!”

有一种感动,叫泪流满面,有一种机制,叫模块机制。显然,这种模块机制给那些 Linux 的发烧友们带来了方便,因为模块机制意味着人们可以把庞大的 Linux 内核划分为许许多多小的模块。对于编写设备驱动程序的开发者来说,从此以后他们可以编写设备驱动程序而不需要把她编译进内核,不用 reboot 机器,她只是一个模块,当你需要她的时候,你可以把她抱入怀中 (insmod),当你不再需要她的时候,你可以把她一脚踢开 (rmmod)。

于是,忽如一夜春风来,内核处处是模块。让我们从一个伟大的例子去认识模块。这就是传说中的“Hello World!”,这个梦幻般的名字我们看过无数次了,每一次她出现在眼前,就意味着我们开始接触一种新的计算机语言了。(某程序员对书法十分感兴趣,退休后决定在这方面有所建树。于是花重金购买了上等的文房四宝。一日,饭后突生雅兴,一番磨墨拟纸,并点上了上好的檀香,颇有王羲之风范,又具颜真卿气势,定神片刻,泼墨挥毫,郑重地写下一行字: hello world)

请看下面这段代码,她就是 Linux 下的一个最简单的模块。当你安装这个模块的时候,她会用她特有的语言向你表白:“Hello, world !”,而后来你卸载了这个模块,你无情抛弃了她,她很伤心,她很绝望,但她没有抱怨,她只是淡淡地说,“Goodbye, cruel world !” (再见,残酷的世界!)

```
/****** hello.c *****/

1 #include <linux/init.h> /* Needed for the macros */
2 #include <linux/module.h> /* Needed for all modules */
3 MODULE_LICENSE("Dual BSD/GPL");
4 MODULE_AUTHOR("fudan_abc");
5
6 static int __init hello_init(void)
7 {
8     printk(KERN_ALERT "Hello, world!\n");
9     return 0;
```



```

10 }
11
12 static void __exit hello_exit(void)
13 {
14     printk(KERN_ALERT "Goodbye, cruel world\n");
15 }
16
17 module_init(hello_init);
18 module_exit(hello_exit);

```

你需要使用 `module_init()` 和 `module_exit()`，你可以称它们为函数，不过实际上它们是一些宏，你可以不用去知道她们背后的故事，只需要知道，在 Linux Kernel 2.6 的世界里，你写的任何一个模块都需要使用它们来初始化或退出，或者说注册以及后来的注销。

当你用 `module_init()` 为一个模块注册了之后，在你使用 `insmod` 这个命令去安装的时候，`module_init()` 注册的函数将会被执行。而当你用 `rmmod` 这个命令去卸载一个模块的时候，`module_exit()` 注册的函数将会被执行。`module_init()` 被称为驱动程序的初始化入口 (driver initialization entry point)。

怎么样演示以上代码的运行呢？没错，你需要一个 Makefile。

```

1 # To build modules outside of the kernel tree, we run "make"
2 # in the kernel source tree; the Makefile these then includes this
3 # Makefile once again.
4 # This conditional selects whether we are being included from the
5 # kernel Makefile or not.
6 ifeq ($(KERNELRELEASE),)
7
8     # Assume the source tree is where the running kernel was built
9     # You should set KERNELDIR in the environment if it's elsewhere
10    KERNELDIR ?= /lib/modules/$(shell uname -r)/build
11    # The current directory is passed to sub-makes as argument
12    PWD := $(shell pwd)
13
14 modules:
15     $(MAKE) -C $(KERNELDIR) M=$(PWD) modules
16
17 modules_install:

```

```
18 $(MAKE) -C $(KERNELDIR) M=$(PWD) modules_install
19
20 clean:
21 rm -rf *.o *~ core .depend *.cmd *.ko *.mod.c .tmp_versions
22
23 .PHONY: modules modules_install clean
24
25 else
26 # called from kernel build system: just declare what our modules are
27 obj-m := hello.o
28 endif
```

在 lwn.net 上可以找到这个例子，你可以把以上两个文件放在你的某个目录下，然后执行 make，也许你不一定能成功，因为 Linux Kernel 2.6 要求你编译模块之前，必须先在内核源代码目录下执行 make，换言之，你必须先配置过内核，执行过 make，然后才能 make 你自己的模块。原因就不细说了，你按着要求的这么去做就行了。在内核顶层目录 make 过之后，你就可以在你当前放置 Makefile 的目录下执行 make 了。make 之后你就应该看到一个叫做 hello.ko 的文件生成了，恭喜你，这就是你将要测试的模块。

执行命令，

```
#insmod hello.ko
```

同时在另一个窗口，用命令 tail -f /var/log/messages 察看日志文件，你会看到 Hello world 被打印了出来。再执行命令，

```
#rmmod hello.ko
```

此时，在另一窗口你会看到 Goodbye，cruel world！被打印了出来。

到这里，我该恭喜你，因为你已经能够编写 Linux 内核模块了。这种感觉很美妙，不是吗？你可以嘲笑秦皇汉武略输文采唐宗宋祖稍逊风骚，还可以嘲笑一代天骄成吉思汗只识弯弓射大雕了。是的，阿娇姐姐告诉我们，只要我喜欢，还有什么不可以。

日后我们会看到，2.6 内核中，每个模块都是以 module_init 开始，以 module_exit

结束。对大多数人来说没有必要知道这是为什么，记住就可以了，对大多数人来说，这就像是 1+1 为什么等于 2 一样，就像是两点之间最短的是直线，不需要证明，如果一定要证明两点之间直线最短，可以扔一块骨头在 B 点，让一条狗从 A 点出发，你会发现狗走的是直线，是的，狗都知道，咱还能不知道吗？

设备模型（上）

对于驱动开发来说，设备模型的理解是根本，毫不夸张得说，理解了设备模型，再去看那些五花八门的驱动程序，你会发现自己站在了另一个高度，从而有了一种俯视的感觉，就像凤姐俯视知音和故事会，韩峰同志俯视女下属。

顾名思义就知道设备模型是关于设备的模型，既不是任小强们的房模，也不是张导的炮模。对咱们写驱动的和不用写驱动的人来说，设备的概念就是总线和与其相连的各种设备了。电脑城的 IT 工作者都会知道设备是通过总线连到计算机上的，而且还需要对应的驱动才能用，可是总线是如何发现设备的，设备又是如何和驱动对应起来的，它们经过怎样的艰辛才找到命里注定的那个他，它们的关系如何，白头偕老型的还是朝三暮四型的，这些问题就不是他们关心的了，而是咱们需要关心的。在房市股市千锤百炼的咱们还能够惊喜的发现，这些疑问的中心思想中心词汇就是总线、设备和驱动，没错，它们就是咱们这里要聊的 Linux 设备模型的名角。

总线、设备、驱动，也就是 bus、device、driver，既然是名角，在内核里都会有它们自己专属的结构，在 include/linux/device.h 里定义。

```
52 struct bus_type {
53     const char          * name;
54     struct module        * owner;
55 }
```



```

56     struct kset          subsys;
57     struct kset          drivers;
58     struct kset          devices;
59     struct klist         klist_devices;
60     struct klist         klist_drivers;
61
62     struct blocking_notifier_head bus_notifier;
63
64     struct bus_attribute  * bus_attrs;
65     struct device_attribute * dev_attrs;
66     struct driver_attribute * drv_attrs;
67     struct bus_attribute drivers_autoprobe_attr;
68     struct bus_attribute drivers_probe_attr;
69
70     int (*match)(struct device * dev, struct device_driver * drv);
71     int (*uevent)(struct device *dev, char **envp,
72     int num_envp, char *buffer, int buffer_size);
73     int      (*probe)(struct device * dev);
74     int      (*remove)(struct device * dev);
75     void      (*shutdown)(struct device * dev);
76
77     int (*suspend)(struct device * dev, pm_message_t state);
78     int (*suspend_late)(struct device * dev, pm_message_t state);
79     int (*resume_early)(struct device * dev);
80     int (*resume)(struct device * dev);
81
82     unsigned int drivers_autoprobe:1;
83 };

124 struct device_driver {
125     const char          * name;
126     struct bus_type      * bus;
127
128     struct kobject       kobj;
129     struct klist         klist_devices;
130     struct klist_node    knode_bus;
131
132     struct module        * owner;
133     const char          * mod_name; /* used for built-in modules */
134     struct module_kobject * mkobj;
135
136     int      (*probe)      (struct device * dev);
137     int      (*remove)     (struct device * dev);
138     void      (*shutdown)  (struct device * dev);

```

```

139     int      (*suspend)      (struct device * dev, pm_message_t state);
140     int      (*resume)       (struct device * dev);
141 };

407 struct device {
408     struct klist      klist_children;
409     struct klist_node  knode_parent;    /* node in sibling list */
410     struct klist_node  knode_driver;
411     struct klist_node  knode_bus;
412     struct device      *parent;
413
414     struct kobject kobj;
415     char bus_id[BUS_ID_SIZE];    /* position on parent bus */
416     struct device_type *type;
417     unsigned    is_registered:1;
418     unsigned    uevent_suppress:1;
419
420     struct semaphore sem; /* semaphore to synchronize calls to
421                          * its driver.
422                          */
423
424     struct bus_type    * bus;        /* type of bus device is on */
425     struct device_driver *driver;    /* which driver has allocated this
426                          device */
427     void      *driver_data; /* data private to the driver */
428     void      *platform_data; /* Platform specific data, device
429                          core doesn't touch it */
430     struct dev_pm_info power;
431
432 #ifdef CONFIG_NUMA
433     int      numa_node; /* NUMA node this device is close to */
434 #endif
435     u64      *dma_mask; /* dma mask (if dma'able device) */
436     u64      coherent_dma_mask; /* Like dma_mask, but for
437                          alloc_coherent mappings as
438                          not all hardware supports
439                          64 bit addresses for consistent
440                          allocations such descriptors. */
441
442     struct list_head dma_pools; /* dma pools (if dma'ble) */
443
444     struct dma_coherent_mem *dma_mem; /* internal for coherent mem
445                          override */
446     /* arch specific additions */

```

```

447     struct dev_archdata    archdata;
448
449     spinlock_t             devres_lock;
450     struct list_head       devres_head;
451
452     /* class_device migration path */
453     struct list_head       node;
454     struct class            *class;
455     dev_t                   devt;           /* dev_t, creates the sysfs "dev" */
456     struct attribute_group  **groups; /* optional groups */
457
458     void (*release)(struct device * dev);
459 };

```

有没有发现它们的共性是什么？对，不是很傻很天真，而是很长很复杂。不过不妨把它们看成艺术品，既然是艺术，当然不会让你那么容易的就看懂了，不然怎么称大师称名家。这么想想咱们就会比较的宽慰了，阿 Q 是鲁迅对咱们 80 后最大的贡献。

我知道进入了 21 世纪，最缺的就是耐性，房价股价都让咱们没有耐性，内核的代码也让人没有耐性。不过做为最没有耐性的一代人，还是要平心静气的扫一下上面的结构，我们会发现，struct bus_type中有成员struct kset drivers 和struct kset devices，同时struct device中有两个成员struct bus_type * bus和struct device_driver *driver，struct device_driver中有两个成员struct bus_type * bus和struct klist klist_devices。先不说什么是klist、kset，光从成员的名字看，它们就是一个完美的三角关系。我们每个人心中是不是都有两个她？一个梦中的她，一个现实中的她。

凭一个男人的直觉，我们可以知道，struct device中的bus表示这个设备连到哪个总线上，driver表示这个设备的驱动是什么，struct device_driver中的bus表示这个驱动属于哪个总线，klist_devices表示这个驱动都支持哪些设备，因为这里device是复数，又是list，更因为一个驱动可以支持多个设备，而一个设备只能绑定一个驱动。当然，struct bus_type中的drivers和devices分别表示了这个总线拥有哪些设备和哪些驱动。

单凭直觉,张钰红不了。我们还需要看看什么是 klist、kset。还有上面 device 和 driver 结构里出现的 kobject 结构是什么?作为一个五星红旗下长大的孩子,我可以肯定的告诉你,kobject 和 kset 都是 Linux 设备模型中最基本的元素,总线、设备、驱动是西瓜,kobject、klist 是种瓜的人,没有幕后种瓜人的汗水不会有清爽解渴的西瓜,我们不能光知道西瓜的甜,还要知道种瓜人的辛苦。kobject 和 kset 不会在意自己的得失,它们存在的意义在于把总线、设备和驱动这样的对象连接到设备模型上。种瓜的人也不会在意自己的汗水,在意的只是能不能送出甜蜜的西瓜。

一般来说应该这么理解,整个 Linux 的设备模型是一个 OO 的体系结构,总线、设备和驱动都是其中鲜活存在的对象,kobject 是它们的基类,所实现的只是一些公共的接口,kset 是同种类型 kobject 对象的集合,也可以说是对象的容器。只是因为 C 里不可能会有 C++ 里类的 class 继承、组合等的概念,只有通过 kobject 嵌入到对象结构里来实现。这样,内核使用 kobject 将各个对象连接起来组成了一个分层的结构体系,就好像马列主义将我们 13 亿人也连接成了一个分层的社会体系一样。kobject 结构里包含了 parent 成员,指向了另一个 kobject 结构,也就是这个分层结构的上一层结点。而 kset 是通过链表来实现的,这样就可以明白,struct bus_type 结构中的成员 drivers 和 devices 表示了一条总线拥有两条链表,一条是设备链表,一条是驱动链表。我们知道了总线对应的数据结构,就可以找到这条总线关联了多少设备,又有哪些驱动来支持这类设备。

那么klist呢?其实它就包含了一个链表和一个自旋锁,我们暂且把它看成链表也无妨,本来在 2.6.11 内核里,struct device_driver结构的devices成员就是一个链表类型。这么一说,咱们上面的直觉都是正确的,如果买股票,摸彩票时直觉都这么管用,就不会有咱们这被压扁的一代了。

现在的人都知道，三角关系很难处。那么总线、设备和驱动之间是如何和谐共处那？先说说总线中的那两条链表是怎么形成的。内核要求每次出现一个设备就要向总线汇报，或者说注册，每次出现一个驱动，也要向总线汇报，或者说注册。比如系统初始化的时候，会扫描连接了哪些设备，并为每一个设备建立起一个 struct device 的变量，每一次有一个驱动程序，就要准备一个 struct device_driver 结构的变量。把这些变量统统加入相应的链表，device 插入 devices 链表，driver 插入 drivers 链表。这样通过总线就能找到每一个设备，每一个驱动。然而，假如计算机里只有设备却没有对应的驱动，那么设备无法工作。反过来，倘若只有驱动却没有设备，驱动也起不了任何作用。在他们遇见彼此之前，双方都如同路埂的野草，一个飘啊飘，一个摇啊摇，谁也不知道未来在哪里，只能在生命的风里飘摇。于是总线上的两张表里就慢慢的就挂上了那许多孤单的灵魂。devices 开始多了，drivers 开始多了，他们像是两个来自世界，devices 们彼此取暖，drivers 们一起狂欢，但他们有一点是相同的，都只是在等待属于自己的那个另一半。

现在，总线上的两条链表已经有了，这个三角关系三个边已经有了两个，剩下的那个那？链表里的设备和驱动又是如何联系那？先有设备还是先有驱动？很久很久以前，在那激情燃烧的岁月里，先有的是设备，每一个要用的设备在计算机启动之前就已经插好了，插放在它应该在的位置上，然后计算机启动，然后操作系统开始初始化，总线开始扫描设备，每找到一个设备，就为其申请一个 struct device 结构，并且挂入总线中的 devices 链表中来，然后每一个驱动程序开始初始化，开始注册其 struct device_driver 结构，然后它去总线的 devices 链表中去寻找(遍历)，去寻找每一个还没有绑定驱动的设备，即 struct device 中的 struct device_driver 指针仍为空的设备，然后它会去观察这种设备的特征，看是否是他所支持的设备，如果是，那么调用一个叫做 device_bind_driver 的函数，然后他们就结为了秦晋之好。换句话说，把 struct device 中的 struct device_driver driver 指向这个驱动，

而 struct device_driver driver 把 struct device 加入他的那张 struct klist klist_devices 链表中来。就这样，bus、device 和 driver，这三者之间或者说他们中的两两之间，就给联系上了。知道其中之一，就能找到另外两个。一荣俱荣，一损俱损。

但现在情况变了，在这红莲绽放的日子里，在这樱花伤逝的日子里，出现了一种新的名词，叫热插拔。设备可以在计算机启动以后在插入或者拔出计算机了。因此，很难再说是先有设备还是先有驱动了。因为都有可能。设备可以在任何时刻出现，而驱动也可以在任何时刻被加载，所以，出现的情况就是，每当一个 struct device 诞生，它就会去 bus 的 drivers 链表中寻找自己的另一半，反之，每当一个 struct device_driver 诞生，它就去 bus 的 devices 链表中寻找它的那些设备。如果找到了合适的，那么 OK，和之前那种情况一下，调用 device_bind_driver 绑定好。如果找不到，没有关系，等待吧，等到昙花再开，等到风景看透，心中相信，这世界上总有一个人是你所等的，只是还没有遇到而已。

设备模型（下）

设备模型拍得再玄幻，它也只是个模型，必须得落实在具体的子系统，否则就只能抱着个最佳技术奖空遗憾。既然前面已经以 USB 子系统的实现分析示例了分析内核源码应该如何入手，那么这里就仍然以 USB 子系统为例，看看设备模型是如何软着陆的。

内核中USB子系统的结构

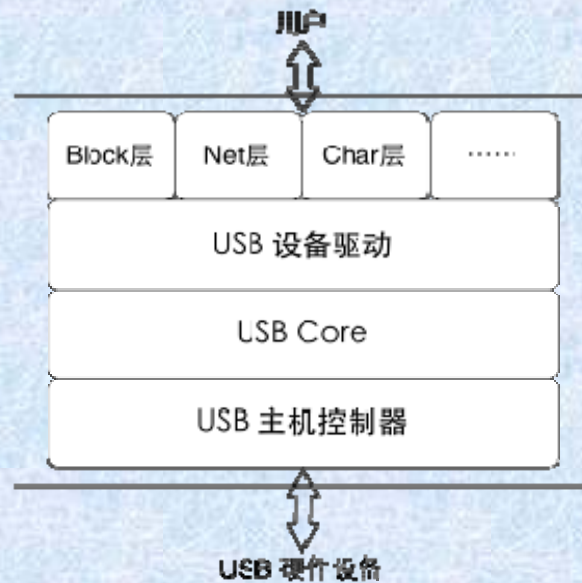
我们已经知道了 USB 子系统的代码都位于 drivers/usb 目录下面，也认识了一个很重要的目录——core 子目录。现在，我们再来看一个很重要的模块——usbcore。你可以使用“lsmod”命令看一下，在显示的结果里能够找到有一个模块叫做 usbcore。

```
localhost:/usr/src/linux-2.6.23/drivers/usb/core # lsmod
```


Module	Size	Used by
af_packet	55820	2
raw	89504	0
nfs	230840	2
lockd	87536	2 nfs
nfs_acl	20352	1 nfs
sunrpc	172360	4 nfs,lockd,nfs_acl
ipv6	329728	36
button	24224	0
battery	27272	0
ac	22152	0
apparmor	73760	0
aamatch_pcre	30720	1 apparmor
loop	32784	0
usbhid	60832	0
dm_mod	77232	0
ide_cd	57120	0
hw_random	22440	0
ehci_hcd	47624	0
cdrom	52392	1 ide_cd
uhci_hcd	48544	0
shpchp	61984	0
bnx2	157296	0
usbcore	149288	4 usbhid,ehci_hcd,uhci_hcd
e1000	130872	0
pci_hotplug	44800	1 shpchp
reiserfs	239616	2
edd	26760	0
fan	21896	0
.....		

找到了 usbcore 那一行吗？core 就是核心，基本上你要在你的电脑里用 USB 设备，那么两个模块是必须的：一个是 usbcore，这就是核心模块；另一个是主机控制器的驱动程序，比如这里 usbcore 那一行我们看到的 ehci_hcd 和 uhci_hcd，你的 USB 设备要工作，合适的 USB 主机控制器模块也是必不可少的。

usbcore 负责实现一些核心的功能，为别的设备驱动程序提供服务，提供一个用于访问和控制 USB 硬件的接口，而不用去考虑系统当前存在哪种主机控制器。至于 core、主机控制器和 USB 驱动三者之间的关系，如下图所示。



USB 驱动和主机控制器就像 core 的两个保镖，协议里也说了，主机控制器的驱动（HCD）必须位于 USB 软件的最下一层。HCD 提供主机控制器硬件的抽象，隐藏硬件的细节，在主机控制器之下是物理的 USB 及所有与之连接的 USB 设备。而 HCD 只有一个客户，对一个人负责，就是 usbcore。usbcore 将用户的请求映射到相关的 HCD，用户不能直接访问 HCD。

core 为咱们完成了大部分的工作，因此咱们写 USB 驱动的时候，只能调用 core 的接口，core 会将咱们的请求发送给相应的 HCD。

USB子系统与设备模型

关于设备模型，最主要的问题就是，bus、device、driver 是如何建立联系的？换言之，这三个数据结构中的指针是如何被赋值的？绝对不可能发生的事情是，一旦为一条总线申请了一个 struct bus_type 的数据结构之后，它就知道它的 devices 链表和 drivers 链表会包含哪些东西，这些东西一定不会是先天就有的，只能是后天填进来的。

具体到 USB 子系统，完成这个工作的就是 USB core。USB core 的代码会进行整个 USB 系统的初始化，比如申请 struct bus_type usb_bus_type，然后会扫描 USB 总线，看

线上连接了哪些 USB 设备 ,或者说 Root Hub 上连了哪些 USB 设备 ,比如说连了一个 USB 键盘 ,那么就为它准备一个 struct device ,根据它的实际情况 ,为这个 struct device 赋值 ,并插入 devices 链表中来。

又比如 Root Hub 上连了一个普通的 Hub ,那么除了要为这个 Hub 本身准备一个 struct device 以外 ,还得继续扫描看这个 Hub 上是否又连了别的设备 ,有的话继续重复之前的事情 ,这样一直进行下去 ,直到完成整个扫描 ,最终就把 usb_bus_type 中的 devices 链表给建立了起来。

那么 drivers 链表呢 ? 这个就不用 bus 方面主动了 ,而该由每一个 driver 本身去 bus 上面登记 ,或者说挂牌。具体到 USB 子系统 ,每一个 USB 设备的驱动程序都会对应一个 struct usb_driver 结构 ,其中有一个 struct device_driver driver 成员 ,USB core 为每一个设备驱动准备了一个函数 ,让它把自己的这个 struct device_driver driver 插入到 usb_bus_type 中的 drivers 链表中去。而这个函数正是我们此前看到的 usb_register。而与之对应的 usb_deregister 所从事的正是与之相反的工作 ,把这个结构体从 drivers 链表中删除。

而 struct bus_type 结构的 match 函数在 USB 子系统里就是 usb_device_match 函数 ,它充当了一个红娘的角色 ,在 USB 总线的 USB 设备和 USB 驱动之间牵线搭桥 ,类似于交大 BBS 上的鹊桥版 ,虽然它们上面的条件都琳琅满目的 ,但明显这里 match 的条件不是那么的苛刻 ,要更为实际些。

可以说 ,USB core 的确是用心良苦 ,为每一个 USB 设备驱动做足了功课 ,正因为如此 ,作为一个实际的 USB 设备驱动 ,它在初始化阶段所要做的事情就很少 ,很简单了 ,直接调用 usb_register 即可。事实上 ,没有人是理所当然应该为你做什么的 ,但 USB core 这么

做了。所以每一个写 USB 设备驱动的人应该铭记，USB 设备驱动绝不是一个人在工作，在他身后，是 USB core 所提供的默默无闻又不可或缺的支持。

驱动开发三件宝：spec、datasheet与内核源码

设备模型之外，对于驱动程序的开发来说，有三样东西是不可缺少的：第一是协议或标准的 spec，也就是规范，比如 usb 协议规范；第二是硬件的 datasheet，即你的驱动要支持的硬件的手册；第三就是内核里类似驱动的源代码，比如你要写触摸屏驱动的话，就可以参考内核里已经有的一些触摸屏驱动。

spec、datasheet、内核源代码这三样东西对于每个开发设备驱动的人来说都是再寻常不过了，但正是因为它们的普通，所以在很多人眼里都被归为被忽视的群体。于是大家开发驱动的过程中，遇到问题的时候首先想到的可能还是“问问牛人怎么解决吧”、“旁边要是有个牛人该多好”，因为牛人的稀有，所以知道牛人的价值，而又因为 spec、datasheet 和内核源代码的唾手可得，所以常常体会不到它们在解决问题时的重要性。

当然我并不是贬低牛人的价值，宣扬依赖牛人不好，如果你很幸运身边真就有牛人这种稀缺资源的话，自然是要好好利用，也可以少走很多弯路，节省很多摸索的时间。只是人生不如意十之八九，多数人还是没有这份幸运的，所以与其遍寻牛人讨教，不如多依赖自己，多利用自己身边有的资源去寻找解决问题的途径。

对这三样看似普通的东西，关键在于很好的去利用，而不是拥有。就说 USB 吧，USB 驱动和 USB 设备如何进行交流，交流的方式，交流过程中出现了什么问题是什么引起的等等都在 USB spec 里有描述，而你的 USB 设备支持多少种配置包含多少端点只有设备的 datasheet 才知道。协议的 spec 和设备的 datasheet 是最好的参考资料，驱动开发调试中

出现的问题绝大部分都能在它们的某个角落里找到答案。而内核中类似设备的驱动源代码是最好的模版，对很多硬件设备，你都可以在内核找到同种设备的驱动代码进行参考实现，甚至于可以拷贝或共享大部分的代码，只进行局部的修改，比如说位于 `drivers/input/touchscreen` 目录下的各个触摸屏驱动，它们之间的代码很多都是类似的甚至是相同的。

如果你不仅仅只是打算写驱动，而是还想阅读内核中实现某种总线、设备的源代码，钻研它们的实现机制，那协议的 spec 就尤为重要，它们在代码里的体现无处不在，你需要在阅读代码前就对协议规范有个整体的理解。形象点说，spec 是理论基础，内核代码是具体实现，理论懂了，看代码就和看故事会差不多了。

Linux内核问题门——学习问题、经验集锦

陈宪章说：“学贵有疑，小疑则小进，大疑则大进。疑者，觉悟之机也，一番觉悟一番长进。” 培根说：“多问的人将多得。” 还在学校的时候导师在激情讲演之后对着会议室里形态各异但均静默不语的我们痛心疾首的说：“会提问题很重要啊，同志们！不会提问题怎么有资格做研究！”

这样铿锵有力的训诫今日想起仍觉深受刺激，于是就要不可避免得要做出一些反应来。不过一是因为咱这年代还没有非主流的说法，二是因为也没有冯仰妍同学的性别优势，不可能受到刺激就整出个门来。咱能够做到的最大反应也就是在这里开贴专门探讨探讨内核学习的相关问题，为了稍微增加那么一些广告效应，就称为“问题门”吧。

使用“问题门”的称呼，一是内心里潜藏的那点低级趣味想去沾点近些年层出不穷各种各样的“门”的仙气，二是在内核的学习过程中的确实实实在在的存在着这样的一个“门”，

横亘在我们的面前，跨过去便海阔天空是另一番世界，但却是让无数人竞折腰，百思不得其钥匙。

另外，这个“问题门”也算是为拙作《Linux 内核修炼之道》制作的一个小插曲，以感谢精华篇发布过程中很多朋友的关心与支持，希望通过对大家内核学习过程中遇到的问题与经验心得做一番展示，来帮助还在门外的朋友寻找到这扇门的钥匙。

我先整理一些在与很多网友交流过程中遇到的部分内核学习问题和自己的一些经验，来作为这个“问题门”的雏形，大家也可以在评论里提出自己的问题和分享自己的学习心得，我会及时地对其进行整理汇总，大家一起将“问题门”逐步完善，帮助后来者和有需要地人不再为“入门”而苦恼。希望这是一篇能够成长得文章！

Linux内核学习常见问题

2010 年 3 月 24 日更新

“问题门”第 5 回：学习 Linux 内核，应该从 Linux 哪个版本代码开始阅读更好呢？

fudan_abc 的回答：

个人建议从新的内核开始，固然新内核的代码非常庞大，但并没有说非要求大求全，追求每个部分都要理解。

学内核忌讳求大而全，如果对哪部分比较感兴趣，研究相关的源码和 change 就行了，当然仁者见仁智者见智，自己如果觉得从低版本开始更好更适合，那采用这种方式也未尝不可，毕竟各人的路还是各自走的。

2010 年 3 月 23 日更新

“问题门”第 3 回：通常，语言及其库的学习分为几个层次，1.熟练使用，2.阅读源码，了

解实现原理，3.对源码进行扩展。那么 linux kernel 怎么划分层次，每个层次如何达到？

(hust_tulip 提出)

fudan_abc 的回答：

问题中的三个层次对应到 linux 内核的学习上：“熟练使用”就是要能够熟练的使用 linux 系统；“阅读源码”就是指“学习内核就是学习内核源代码”，必须勇敢的去学习内核源码；“对源码进行扩展”可以对应于融入内核社区，参与内核的开发。

这也正好在一定程度上契合了本书前言里对内核学习划分的几个层次：全面了解抓基本，兴趣导向深钻研；融入社区做贡献，坚持坚持再坚持。

——详见[修炼之道 之 前言](#)

“问题门”第 4 回：每个层次的学习都有什么对应的参考资料以及网络资源？(hust_tulip 提出)

fudan_abc 的回答：

首先是“全面了解抓基本”，这个层次，最好的书自然就是 lkd 和 ulk 了，这两本书，一本提纲挈领，一本全面深入，都能很好的帮助全面的理解内核的整体机制。新人的话，一本 lkd 就足亦了。

第二个层次“兴趣导向深钻研”，这个层次就是要以内核源码为中心，选择内核中一个自己感兴趣的部分，以韩峰同志对待日记的态度，严谨而细致的仔细分析它的代码，不懂的地方就通过社区、邮件列表或者直接发 Email 给 maintainer 请教等途径弄懂，切勿得过且过。至于这个层次的参考书么，网络子系统的有《深入理解 LINUX 网络内幕》，内存管理的有《深入理解 Linux 虚拟内存管理》，推荐看英文版，呵呵，usb 的可以看我们的《Linux

那些事儿》，其它子系统的还没注意到有什么专门讲解的，不过内核源码本身就是最好的参考资料了。

至于第三个层次“融入社区做贡献”，就是要努力融入到内核的开发社区，经过前两个层次的修炼，此时你已经不会再是社区中潜水小白的角色，而是会针对某个问题发表自己的见解。可以尝试参与到内核的开发中去。相关资源有很多，详细可以参考[修炼之道精华篇的\(9\)](#)

[内核学习资源](#)

最后一层就是坚持了，不管遇到什么挫折都不放弃，就像咱们的袁教授不管遭受到什么样的辱骂都要坚持不断的发疯一样，有这样的精神，何愁在 linux 内核的学习道路上修不成大道那？

2010 年 3 月 22 日更新

“问题门”第 1 回：我是一个初学者，两眼一抹黑，我该如何学习内核？

fudan_abc 的回答：

这个问题每个初学者都无法回避，它非常之大，完全可以做为整个“问题门”的框架而存在，其他的各种问题都不过是在这个框架上装饰和完善。

同时这个问题并没有一个标准的答案，只有一些学习的脉络可以遵循，祝早日“入门”。

第一步：先会使用它。连 Linux 是什么、基本操作都不会就去研究内核 纯属扯淡，“门”都没有。

第二步：看懂内核源码需要一些操作系统、C 语言等的基础。

第三步：找本合适的内核参考书，让它帮助你对内核有个整体的理解和认识，

第四步：要能够动手配置编译内核，还要基本看得懂内核中的 Kconfig 和 Makefile 文件。

最后，记住：“学习内核，就是学习内核的源代码，任何内核有关的书籍都是基于内核，而又不高于内核的。内核源码本身就是最好的参考资料，其他任何经典或非经典的书最多只是起到个辅助作用，不能也不应该取代内核代码在我们学习过程中的主导地位。”

因此你要做得是选择内核的一个部分或子系统，以韩峰同志对待日记的态度，严谨而细致得理解每一段代码的实现，多问多想多记。切勿抱着走马观花，得过且过的态度。

“问题门”第 2 回：学习内核需要什么基础知识？

albcamus 的回答：

(1) 需要掌握操作系统理论的最初级的知识。

不需要通读并理解《操作系统概念》《现代操作系统》等巨著，但总要知道分时 (time-shared) 和实时 (real-time) 的区别是什么，进程是个什么东西，CPU 和系统总线、内存的关系 (很粗略即可)，等等。

(2) 熟练使用 C 语言。

不需要已经很精通 C 语言，只要能熟练编写 C 程序，能看懂链表、散列表等数据结构的 C 实现，用过 gcc 编译器，就可以了。当然，如果已经精通 C 语言显然是大占便宜的。

(3) 了解 CPU 的相关知识。

Linux 内核学习经验

1. 内核学习的心理误区

心理上的问题主要有两个，一个是盲目，就是在能够熟练使用 Linux 之前，对 Linux 为何物还说不出的道道来，就迫不及待的盲目的去研究内核的源代码。重述 Linus 的那句话：要先会使用它。

第二个就是恐惧。人类进化这么多年，面对复杂的物体和事情还是总会有天生的惧怕感，体现在内核学习上面就是：那么庞大复杂的内核代码，让人面对起来该情何以堪啊！

有了这种恐惧无力感存在，心理上就会去排斥面对接触内核源码，宁愿去抱着情景分析，搜集各种各样五花八门的内核书籍放在那里屯着，看了又忘，忘了又看，也不大情愿去认真细致得浏览源码。

——详见[修炼之道精华篇（9）内核学习的心理问题](#)

2. 学习内核就是学习内核的源代码

学习内核，就是学习内核的源代码，任何内核有关的书籍都是基于内核，而又不高于内核的。内核源码本身就是最好的参考资料，其他任何经典或非经典的书籍最多只是起到个辅助作用，不能也不应该取代内核代码在我们学习过程中的主导地位。

3. 要抱着严谨细致的态度分析内核源码

既然要学习内核源码，就要经常对内核代码进行分析，而内核代码千千万，还前仆后继的不断往里加，这就让大部分人都有种雾里看花花不见的无助感。不过不要怕，孔老夫子早就留给我们了应对之策：敏于事而慎于言，就有道而正焉，可谓好学也已。这就是说，做事要踏实才是好学生好同志，要遵循严谨的态度，去理解每一段代码的实现，多问多想多记。如果抱着走马观花，得过且过的态度，结果极有可能就是一边看一边丢，没有多大的收获。

只要你使用这样的态度开始分析内核，那么无论你选择内核的哪个部分作为切入点，比

如 USB，比如进程管理，在花费相对不算很多的时间之后，你就会发现你对内核的理解会上升到另外一个高度，一个抱着情景分析，抱着 0.1 内核完全注释，抱着各种各样的内核书籍翻来覆去的看很多遍又忘很多遍都无法达到的高度。

——详见修炼之道[精华篇（6）](#)与[精华篇（7）](#)分析内核源码如何入手？

4. 通过 Kconfig 与 Makefile 定位目标代码

毫不夸张地说，Kconfig 和 Makefile 是我们浏览内核代码时最为依仗的两个文件。基本上，Linux 内核中每一个目录下边都会有一个 Kconfig 文件和一个 Makefile 文件。对于一个希望能够在 Linux 内核的汪洋代码里看到一丝曙光的人来说，将它们放在怎么重要的地位都不过分。

Kconfig 和 Makefile 就是 Linux Kernel 迷宫里的地图。地图引导我们去认识一个城市，而 Kconfig 和 Makefile 则可以让我们了解一个 Kernel 目录下面的结构。我们每次浏览 kernel 寻找属于自己的那一段代码时，都应该首先看看目录下的这两个文件。就像利用地图寻找目的地一样，我们需要利用 Kconfig 和 Makefile 来寻找所要研究的目标代码。

——详见[修炼之道精华篇（5）Kernel地图：Kconfig与Makefile](#)

《Linux内核修炼之道》精华分享与讨论

（14）——内核中的链表

早上上班坐地铁要排队，到了公司楼下等电梯要排队，中午吃饭要排队，下班了追求一个女孩子也要排队，甚至在网上传下载个什么门的短片也要排队，每次看见

人群排成一条长龙时，才真正意识到自己是龙的传人。那么下面咱们就说说队列（链表）。

使用链表的目的很明确，因为有很多事情要做，于是就把它放进链表里，一件事一件事的处理。比如在 USB 子系统里，U 盘不停的提交 urb 请求，USB 键盘也提交，USB 鼠标也提交，那 USB 主机控制器咋应付得过来呢？很简单，建一个链表，然后你每次提交就是往里边插入，然后 USB 主机控制器再统一去调度，一个一个来执行。这里有力得证明了，谭浩强大哥的 C 程序设计是我们学习 Linux 的有力武器，书中对链表的介绍无疑是英明的，谭大哥，您不是一个人在战斗！

内核中链表的实现位于 include/linux/list.h 文件，链表数据结构的定义也很简单。

```
21 struct list_head {  
  
22     struct list_head *next, *prev;  
  
23 };
```

list_head 结构包含两个指向 list_head 结构的指针 prev 和 next，由此可见，内核中的链表实际上都是双链表（通常都是双循环链表）。

通常，我们在数据结构课堂上所了解的链表定义方式是这样的（以单链表为例）：

```
struct list_node {  
  
    struct list_node *next;
```



```
ElemType data;
```

```
};
```

通过这种方式使用链表，对每一种数据类型，都要定义它们各自的链表结构。而内核中的链表却与此不同，它并没有数据域，不是在链表结构中包含数据，而是在描述数据类型的结构中包含链表。

比如在 hub 驱动中使用 struct usb_hub 来描述 hub 设备，hub 需要处理一系列的事件，比如当探测到一个设备连进来时，就会执行一些代码去初始化该设备，所以 hub 就创建了一个链表来处理各种事件，这个链表的结构如下图。



(1) 声明与初始化。

链表的声明可以使用两种方式，一种为使用 LIST_HEAD 宏在编译时静态初始化，一种为使用 INIT_LIST_HEAD() 在运行时进行初始化。

```
25 #define LIST_HEAD_INIT(name) { &(amp;name), &(name) }
```

```
26
```

```
27 #define LIST_HEAD(name) \
```

```
28 struct list_head name = LIST_HEAD_INIT(name)
```

29

```
30 static inline void INIT_LIST_HEAD(struct list_head *list)
```

```
31 {
```

```
32     list->next = list;
```

```
33     list->prev = list;
```

```
34 }
```

无论采用哪种方式，新生成的链表头的两个指针 next、prev 都初始化为指向自己。

(2) 判断链表是否为空。

```
298 static inline int list_empty(const struct list_head *head)
```

```
299 {
```

```
300     return head->next == head;
```

```
301 }
```

(3) 插入。

有了链表，自然就要往里面加东西、减东西。就像我们每个人每天都在不停的走进去，又走出来，似是梦境又不是梦境。一切都是不经意的。走进去是一年四季，

走出来是春夏秋冬。list_add()和 list_add_tail()这两个函数就是往队列里加东西。

```
67 static inline void list_add(struct list_head *new, struct list_head *head)
```

```
68 {
```

```
69     __list_add(new, head, head->next);
```

```
70 }
```

```
84 static inline void list_add_tail(struct list_head *new, struct list_head  
*head)
```

```
85 {
```

```
86     __list_add(new, head ->prev, head);
```

```
87 }
```

其中，list_add()将数据插入在 head 之后，list_add_tail()将数据插入在 head->prev 之后。其实对于循环 链表来说，表头的 next、prev 分别指向链表中的第一个和最后一个节点，所以，list_add()和 list_add_tail()的区别并不大。

(4) 删除。

搞懂谭浩强那本书之后看这些链表的代码那就是小菜一碟。再来看下一个 list_del_init()，里的元素不能只加不减，没用了的元素就该删除掉，把空间腾出

来给别人。郭敬明说过，我生命里的温暖就那么多，我全部给了 你，但是你离开了我，你叫我以后怎么再对别人笑.....

链表里的元素不能只加不减，没用了的元素就应该删除掉。

```
254 static inline void list_del_init(struct list_head *entry)

255 {

256     __list_del(entry->prev, entry->next);

257     INIT_LIST_HEAD(entry);

258 }
```

list_del_init()从链表里删除一个元素，并且将其初始化。

(5) 遍历。

内核中的链表仅仅保存了 list_head 结构的地址，我们如何通过它或取一个链表节点真正的数据项？这就要提到有关链表的所有操作里面，最为重要 超级经典的 list_entry 宏了，我们可以通过它很容易地获得一个链表节点的数据。

```
425 #define list_entry(ptr, type, member) \

426     container_of(ptr, type, member)
```

我相信，list_entry()这个宏在 Linux 内核代码中的地位，就相当于广告词中的任静付笛生的洗洗更健康，相当于大美女关之琳的一分钟轻松做女人，这都是耳

熟能详妇孺皆知的，是经典中的经典。如果你说你不知道 `list_entry()`，那你千万别跟人说你懂 Linux 内核，就好比你不知道 陈文登不知道任汝芬你就根本不好意思跟人说你考过研，要知道每个考研人都是左手一本陈文登右手一本任汝芬。

可惜，关于 `list_entry`，这个谭浩强老师的书里就没有了，当然你不能指责谭浩强的书不行，再好的书也不可能包罗万象。

关于 `list_entry()`，让我们结合实例来看，还是 hub 驱动的那个例子，当我们真的要处理 hub 的事件的时候，我们当然需要知道具体是哪个 hub 触发了这起事件。而 `list_entry` 的作用就是，从 `struct list_head event_list` 得到它所对应的 `struct usb_hub` 结构体变量。比如以下四行代码：

```
struct list_head *tmp;

struct usb_hub *hub;

tmp = hub_event_list.next;

hub = list_entry(tmp, struct usb_hub, event_list);
```

从全局链表 `hub_event_list` 中取出一个来，叫做 `tmp`，然后通过 `tmp`，获得它所对应的 `struct usb_hub`。

《Linux内核修炼之道》精华分享与讨论

(15) ——子系统的初始化：内核选项解析

首先感谢国家。其次感谢上大的钟莉颖，让我知道了大学不仅有校花，还有校鸡，而且很多时候这两者其实没什么差别。最后感谢清华女刘静，让我深刻体会到了素质教育的重要性，让我感到有责任写写子系统的初始化。

各个子系统的初始化是内核整个初始化过程必然要完成的基本任务，这些任务按照固定的模式来处理，可以归纳为两个部分：内核选项的解析以及那些子系统入口（初始化）函数的调用。

内核选项

Linux 允许用户传递内核配置选项给内核，内核在初始化过程中调用 `parse_args` 函数对这些选项进行解析，并调用相应的处理函数。

`parse_args` 函数能够解析形如“变量名=值”的字符串，在模块加载时，它也会被调用来解析模块参数。

内核选项的使用格式同样为“变量名=值”，打开系统的 `grub` 文件，然后找到 `kernel` 行，比如：


```
kernel /boot/vmlinuz-2.6.18 root=/dev/sda1 ro splash=silent  
vga=0x314 pci=noacpi
```

其中的 “pci=noacpi” 等都表示内核选项。

内核选项不同于模块参数，模块参数通常在模块加载时通过 “变量名=值” 的形式指定，而不是内核启动时。如果希望在内核启动时使用模块参数，则必须添加模块名做为前缀，使用 “模块名.参数=值” 的形式，比如，使用下面的命令在加载 usbcore 时指定模块参数 autosuspend 的值为 2。

```
$ modprobe usbcore autosuspend=2
```

若是在内核启动时指定，则必须使用下面的形式：

```
usbcore.autosuspend=2
```

从 Documentation/kernel-parameters.txt 文件里可以查询到某个子系统已经注册的内核选项，比如 PCI 子系统注册 的内核选项为：

```
pci=option[,option...] [PCI] various PCI subsystem options:
```

```
off [X86-32] don't probe for the PCI bus
```

```
bios [X86-32] force use of PCI BIOS, don't access  
the hardware directly. Use this if your machine  
has a non-standard PCI host bridge.
```

```
nobios [X86-32] disallow use of PCI BIOS, only direct  
hardware access methods are allowed. Use this
```

if you experience crashes upon bootup and you suspect they are caused by the BIOS.

conf1 [X86-32] Force use of PCI Configuration

Mechanism 1.

conf2 [X86-32] Force use of PCI Configuration

Mechanism 2.

nommconf [X86-32,X86_64] Disable use of MMCONFIG for PCI Configuration

nomsi [MSI] If the PCI_MSI kernel config parameter is enabled, this kernel boot option can be used to disable the use of MSI interrupts system-wide.

nosort [X86-32] Don't sort PCI devices according to order given by the PCI BIOS. This sorting is done to get a device order compatible with older kernels.

biosirq [X86-32] Use PCI BIOS calls to get the interrupt routing table. These calls are known to be buggy on several machines and they hang the machine when used, but on other computers it's the only way to get the interrupt routing table. Try this option if the kernel is unable to allocate IRQs or discover secondary PCI buses on your

motherboard.

rom [X86-32] Assign address space to expansion ROMs.

Use with caution as certain devices share address decoders between ROMs and other resources.

irqmask=0xMMMM [X86-32] Set a bit mask of IRQs allowed to be assigned automatically to PCI devices. You can make the kernel exclude IRQs of your ISA cards this way.

pirqaddr=0xAAAAA [X86-32] Specify the physical address of the PIRQ table (normally generated by the BIOS) if it is outside the F0000h-100000h range.

lastbus=N [X86-32] Scan all buses thru bus #N. Can be useful if the kernel is unable to find your secondary buses and you want to tell it explicitly which ones they are.

assign-busses [X86-32] Always assign all PCI bus numbers ourselves, overriding whatever the firmware may have done.

useirqmask [X86-32] Honor the possible IRQ mask stored in the BIOS \$PIR table. This is needed on

some systems with broken BIOSes, notably
some HP Pavilion N5400 and Omnibook XE3
notebooks. This will have no effect if ACPI
IRQ routing is enabled.

`noacpi` [X86-32] Do not use ACPI for IRQ routing
or for PCI scanning.

`routeirq` Do IRQ routing for all PCI devices.

This is normally done in `pci_enable_device()`,
so this option is a temporary workaround
for broken drivers that don't call it.

`firmware` [ARM] Do not re-enumerate the bus but instead
just use the configuration from the
bootloader. This is currently used on
IXP2000 systems where the bus has to be
configured a certain way for adjunct CPUs.

`noearly` [X86] Don't do any early type 1 scanning.

This might help on some broken boards which
machine check when some devices' config space
is read. But various workarounds are disabled
and some IOMMU drivers will not work.

`bfsort` Sort PCI devices into breadth-first order.

This sorting is done to get a device

order compatible with older (≤ 2.4) kernels.

nobfsort Don't sort PCI devices into breadth-first order.

cbiosize=nn[KMG] The fixed amount of bus space which is reserved for the CardBus bridge's IO window.

The default value is 256 bytes.

cbmemsize=nn[KMG] The fixed amount of bus space which is reserved for the CardBus bridge's memory window. The default value is 64 megabytes.

注册内核选项

就像我们不需要明白钟莉颖是如何走上校鸡的修炼之道，我们也不必理解 `parse_args` 函数的实现细节。但我们必须知道如何注册内核选项：模块参数使用 `module_param` 系列的宏注册，内核选项则使用 `__setup` 宏来注册。

`__setup` 宏在 `include/linux/init.h` 文件中定义。

```
171 #define __setup(str, fn)  \
172  __setup_param(str, fn, 0)
```

`__setup` 需要两个参数，其中 `str` 是内核选项的名字，`fn` 是该内核选项关联的处理函数。`__setup` 宏告诉内核，在启动时如果检测到内核选项 `str`，则执行函数 `fn`。`str` 除了包括内核选项名字之外，必须以 `"="` 字符结束。

不同的内核选项可以关联相同的处理函数，比如内核选项 netdev 和 ether 都关联了 netdev_boot_setup 函数。

除了__setup 宏之外，还可以使用 early_param 宏注册内核选项。它们的使用方式相同，不同的是，early_param 宏注册的内核选项必须要在其他内核选项之前被处理。

两次解析

相应于__setup 宏和 early_param 宏两种注册形式，内核在初始化时，调用了两次 parse_args 函数进行解析。

```
parse_early_param();  
  
parse_args("Booting kernel", static_command_line, __start__param,  
          __stop__param - __start__param,  
          &unknown_bootoption);
```

parse_args 的第一次调用就在 parse_early_param 函数里面，为什么会出现两次调用 parse_args 的情况？这是因为内核选项又分成了两种，就像现实世界中的我们，一种是普普通通的，一种是有特权的，有特权的需要在普通选项之前进行处理。

现实生活中特权的定义好像很模糊，不同的人有不同的诠释，比如哈医大二院的纪委书记在接受央视的采访“老人住院费 550 万元”时如是说：“我们就是 一所人民医院.....就是一所贫下中农的医院，从来不用特权去索取自己身外的任何

利益.....我们不但没有多收钱还少收了。”

人生就是如此的复杂和奇怪。 内核选项相对来说就要单纯得多，特权都是阳光下的，不会藏着掖着，直接使用 `early_param` 宏去声明，让你一眼就看出它是有特权的。使用 `early_param` 声明的那些选项就会首先由 `parse_early_param` 去解析。

《Linux内核修炼之道》精华分享与讨论

(16) ——子系统的初始化：那些入口函数

内核选项的解析完成之后，各个子系统的初始化即进入第二部分——入口函数的调用。通常 USB、PCI 这样的子系统都会有一个名为 `subsys_initcall` 的入口，如果你选择它们作为研究内核的切入点，那么就请首先找到它。

朱德庸在《关于上班这件事》里说，要花前半生找入口，花后半生找出口。可见寻找入口对于咱们这一生，对于看内核代码这件事儿都是无比重要的。

但是很多时候，入口并不仅仅只有 `subsys_initcall` 一个，比如 PCI。

```
117 #define pure_initcall(fn)          __define_initcall("0",fn,1)
118
119 #define core_initcall(fn)           __define_initcall("1",fn,1)
120 #define core_initcall_sync(fn)      __define_initcall("1s",fn,1s)
```

```

121 #define postcore_initcall(fn)      __define_initcall("2",fn,2)
122 #define postcore_initcall_sync(fn)  __define_initcall("2s",fn,2s)
123 #define arch_initcall(fn)          __define_initcall("3",fn,3)
124 #define arch_initcall_sync(fn)      __define_initcall("3s",fn,3s)
125 #define subsys_initcall(fn)         __define_initcall("4",fn,4)
126 #define subsys_initcall_sync(fn)    __define_initcall("4s",fn,4s)
127 #define fs_initcall(fn)             __define_initcall("5",fn,5)
128 #define fs_initcall_sync(fn)        __define_initcall("5s",fn,5s)
129 #define rootfs_initcall(fn)         __define_initcall("rootfs",fn,rootfs)
130 #define device_initcall(fn)         __define_initcall("6",fn,6)
131 #define device_initcall_sync(fn)     __define_initcall("6s",fn,6s)
132 #define late_initcall(fn)           __define_initcall("7",fn,7)
133 #define late_initcall_sync(fn)       __define_initcall("7s",fn,7s)
134
135 #define __initcall(fn) device_initcall(fn)

```

这些入口有个共同的特征，它们都是使用__define_initcall 宏定义的。它们的调用也不是随便的，而是按照一定顺序的，这个顺序就取决于__define_initcall 宏。__define_initcall 宏用来将指定的函数指针放到.initcall.init 节里。

.initcall.init 节

内核可执行文件由许多链接在一起的对象文件组成。对象文件有许多节，如文本、数据、init 数据、bss 等等。这些对象文件都是由一个称为链接器 脚本的文件

链接并装入的。这个链接器脚本的功能是将输入对象文件的各节映射到输出文件中；换句话说，它将所有输入对象文件都链接到单一的可执行文件中，将该可执行文件的各节装入到指定地址处。vmlinux.lds 是存在于 arch/<target>/目录中的内核链接器脚本，它负责链接内核的各个节并将它们装入内存中特定偏移量处。在 vmlinux.lds 文件里查找 initcall.init 就可以看到下面的内容

```
__initcall_start = .;

.initcall.init : AT(ADDR(.initcall.init) - 0xC0000000) {

*(.initcall1.init)

*(.initcall2.init)

*(.initcall3.init)

*(.initcall4.init)

*(.initcall5.init)

*(.initcall6.init)

*(.initcall7.init)

}

__initcall_end = .;
```

这就告诉我们.initcall.init 节又分成了 7 个子节，而 xxx_initcall 入口函数指针具体放在哪一个子节里边儿是由 xxx_initcall 的定义中，__define_initcall 宏的参数决定的，比如 core_initcall 将函数指针放在.initcall1.init 子节，device_initcall 将函数指针放在了.initcall6.init 子节等等。各个子节的顺序是确

定的,即先调用.initcall1.init 中的函数指针再调用.initcall2.init 中的函数指针,等等。不同的入口函数被放在不同的子节中,因此也就决定了它们的调用顺序。

do_initcalls()函数

那些入口函数的调用由 do_initcalls 函数来完成。

do_initcall 函数通过 for 循环,由__initcall_start 开始,直到__initcall_end 结束,依次调用识别到的初始化函数。而位于__initcall_start 和__initcall_end 之间的区域组成了.initcall.init 节,其中保存了由 xxx_initcall 形式的宏标记的函数地址,do_initcall 函数可以很轻松的取得函数地址并执行其指向的函数。

.initcall.init 节所保存的函数地址有一定的优先级,越前面的函数优先级越高,也会比位于后面的函数先被调用。

由 do_initcalls 函数调用的函数不应该改变其优先级状态和禁止中断。因此,每个函数执行后,do_initcalls 会检查该函数是否做了任何变化,如果有必要,它会校正优先级和中断状态。

另外,这些被执行的函数有可以完成一些需要异步执行的任务,flush_scheduled_work 函数则用于确保 do_initcalls 函数在返回前等待这些异步任务结束。

```
666 static void __init do_initcalls(void)
667 {
668     initcall_t *call;
```

```
669 int count = preempt_count();
670
671 for (call = __initcall_start; call < __initcall_end; call++) {
672     ktime_t t0, t1, delta;
673     char *msg = NULL;
674     char msgbuf[40];
675     int result;
676
677     if (initcall_debug) {
678         printk("Calling initcall 0x%p", *call);
679         print_fn_descriptor_symbol(": %s()",
680             (unsigned long) *call);
681         printk("\n");
682         t0 = ktime_get();
683     }
684
685     result = (*call)();
686
687     if (initcall_debug) {
688         t1 = ktime_get();
689         delta = ktime_sub(t1, t0);
690
```

```

691     printk("initcall 0x%p", *call);
692     print_fn_descriptor_symbol(":%s()",
693         (unsigned long) *call);
694     printk(" returned %d.\n", result);
695
696     printk("initcall 0x%p ran for %Ld msecs: ",
697         *call, (unsigned long long)delta.tv64 >> 20);
698     print_fn_descriptor_symbol("%s()\n",
699         (unsigned long) *call);
700 }
701
702 if (result && result != -ENODEV && initcall_debug) {
703     sprintf(msgbuf, "error code %d", result);
704     msg = msgbuf;
705 }
706 if (preempt_count() != count) {
707     msg = "preemption imbalance";
708     preempt_count() = count;
709 }
710 if (irqs_disabled()) {
711     msg = "disabled interrupts";
712     local_irq_enable();

```



```
713 }  
714 if (msg) {  
715     printk(KERN_WARNING "initcall at 0x%p", *call);  
716     print_fn_descriptor_symbol(": %s()",  
717         (unsigned long) *call);  
718     printk(": returned with %s\n", msg);  
719 }  
720 }  
721  
722 /* Make sure there is no pending stuff from the initcall sequence */  
723 flush_scheduled_work();  
724 }
```

《Linux内核修炼之道》精华分享与讨论

(17)——子系统的初始化：以PCI子系统为例

由 Kconfig 这张地图的分布来看，PCI 这块儿的代码应该分布在两个地方，drivers/pci 和 arch/i386/pci，两岸三地都属于一个中国，不管是 drivers/pci 那儿的，还是 arch/i386/pci 那儿的，也都只属于一个 PCI 子系统，本着一个中国的原则，咱们要统筹 的全面的考察分析位于两个地方的代码，于是，这些远

远突破了五位数的代码左看右看横看竖看都显得那么的阴森恐怖,不过人家咋说也是整个一 PCI 子系统,就 像走在 T 台上的芙蓉姐姐和杨二车那姆一样,看起来恐怖但也是很有内涵的,岂能够让人三眼两眼三言两语就给看透了说透了?

那现在咱们就高瞻远瞩统筹全面的扫视一下这两个地方的代码,根据前面的内容可以推测对于 USB、PCI 这样的子系统都应该有一个 subsys_initcall 这样的入口,咱们得先找到它。朱德庸在《关于上班这件事》里说了,要花前半生找入口,花后半生找出口。可见寻找入口对于咱们这一生,对于看内核代码这件事儿都是无比重要的,当然寻找 subsys_initcall 这个入口是不用花前半生那么久的。下边儿俺就把找到的给列出来,为什么说“列”出来?难道还会有很多么?你猜对了,PCI 这边儿入口格外多,而且是有预谋有组织成系列的,不单单有 subsys_initcall, 还有 arch_initcall、postcore_initcall 等等等等。

文 件	函 数	入
口	内存位置	
arch/i386/pci/acpi.c	pci_acpi_init	subsys_initcall
.initcall4.init		
arch/i386/pci/common.c	pcibios_init	subsys_initcall
ll .initcall4.init		
arch/i386/pci/i386.c	pcibios_assign_resources	fs_initcall
.initcall5.init		
arch/i386/pci/legacy.c	pci_legacy_init	subsys_initcall
.initcall4.init		

drivers/pci/pci-acpi.c	acpi_pci_init	arch_initcall
.initcall3.init		
drivers/pci/pci-driver.c	pci_driver_init	postcore_initc
all .initcall2.init		
drivers/pci/pci-sysfs.c	pci_sysfs_init	late_initcall
.initcall7.init		
drivers/pci/pci.c	pci_init	device_initcall
.initcall6.init		
drivers/pci/probe.c	pcibus_class_init	postcore_initca
ll .initcall2.init		
drivers/pci/proc.c	pci_proc_init	__initcall
.initcall6.init		
arch/i386/pci/init.c	pci_access_init	arch_initcall
.initcall3.init		

看看那一列入口，形尽而意不同的种种 xxx_initcall 让人眼花缭乱的，真不知道该从哪儿下手，应了 keso 那句话：所有的痛苦都来自选择，所谓幸福，就是没有选择。像 USB 子系统那样子简简单单一个 subsys_initcall，没得选择，傻强都知道怎么走。不过你迷惘一阵儿就可以了，可别真的被绕进去了。要知道“多少事，从来急；天地转，光阴迫。一万年太久，只争朝夕。四海翻腾云水怒，五洲震荡风雷激。要看清一切入口，全无敌。”咱们要只争朝夕看清一切入口的。

咱们已经知道对这些 xxx_initcall 函数的调用是必须按照一定顺序的，先调用.initcall1.init 中的再调用.initcall2.init 中的，很明显，表里列出来的应该最先被调用的是.initcall2.init 子节中的两个函数 pcibus_class_init 和 pci_driver_init。现在问题出现了，对于处于同一子节中的那些函数，比如 pcibus_class_init 和 pci_driver_init 这两个函数来说又是哪个会最先被调用？当然，你可以说处在前面地址的会最先被调用，这是大实话，因为 do_initcalls 函数的实现就是在.initcall.init 所处的地址上来回的 for 循环。可你怎么知道同一子节的函数 哪个在前面哪个在后面？

别的不多说，先看看 gcc 的 Using the GNU Compiler Collection 中的一段话：

the linker searches and processes libraries and object files in the order they are specified. Thus, 'foo.o -lz bar.o' searches library 'z' after file 'foo.o' but before 'bar.o' .

看完这段话，希望会听到你说：我悟道了！更希望会看到你翻出来 drivers/pci/Makefile 文件，瞅到下边儿这两行

```
5 obj-y      += access.o bus.o probe.o remove.o pci.o quirks.o \  
6            pci-driver.o search.o pci-sysfs.o rom.o setup-res.o
```

probe.o 在 pci- driver.o 的前面，那么 probe.c 里的 pcibus_class_init 函数也会在 pci- driver.c 里的 pci_driver_init 函数之前被调用。再

给你看一句话，Documents/kbuild/makefile.txt 的 3.2 中的：

The order of files in \$(obj-y) is significant.

对于 pcibus_class_init 函数和 pci_driver_init 函数这样位于同一目录位置的可以通过该目录 Makefile 文件指定的链接顺序来判断,而对于 initcall3.init 子节中的 acpi_pci_init 函数和 pci_access_init 函数则不能使用这个方法。

acpi_pci_init 在 drivers/pci/pci-acpi.c 文件里,而 pci_access_init 在 arch/i386/pci/init.c 文件里,它俩根本就不在同一个目录下面,所以前边儿判断

pcibus_class_init 和 pci_driver_init 的顺序的技巧并不适用,那有什么方法可以让咱们找出它们的顺序?看看王冉怎么说:“昨天是五一劳动节,可是全国都在放大假绝大多数人不劳动。可见,庆祝一件事的最好的方法就是不去做这件事。譬如,庆祝世界杯的最好的方式就是不去参加世界杯——中国队几乎一直都是这么做的。再譬如,庆祝情人节的最好的方式就是不去找情人——于是,很多中国的男人把情人节的前一天(2月13日)过成了情人节。”按他这说法,认清这两函数之间顺序的最好方法就是不去管它们的顺序,俺可以点兵点将的随便点一个出来先说,不过作为一个很清楚自己责任和使命的80后,俺还是决定去发掘一下它们的顺序。

其实这个问题可以转化为 arch/i386/pci 下面的 Makefile 和 drivers/pci 下面的 Makefile 谁先谁后的问题,往大的方面说,就是内核是怎么构建的,也就是 kbuild 的问题。

内核里的 Makefile 主要有三种:第一种是根目录里的 Makefile,它虽然只有一个,但地位远远凌驾于其它 Makefile 之上,里面定义了所有与体系结构无关的

变量和目标；第二种是 arch/*/Makefile，看到 arch 就知道它是与特定体系结构相关的，它包含在根目录下的 Makefile 中，为 kbuild 提供体系结构的特定信息，而它里面又包含了 arch/*/下面各级子目录的那些 Makefile；第三种就是密密麻麻 躲在 drivers/等各个子目录下边儿的那些 Makefile 了。

而 kbuild 构建内核的过程中，是首先从根目录 Makefile 开始执行，从中获得与体系结构无关的变量和依赖关系，并同时从 arch/* /Makefile 中获得体系结构特定的变量等信息，用来扩展根目录 Makefile 所提供的变量。此时 kbuild 已经拥有了构建内核需要的所有变量和 目标，然后，Make 进入各个子目录，把部分变量传递给子目录里的 Makefile，子目录 Makefile 根据配置信息决定编译哪些源文件，从而构建出 一个需要编译的文件列表。

然后，然后还有很漫长的路，你编译内核要耗多久，它就有多漫长，不过说到这儿前面问题的答案就已经浮出水面了，很明显，arch/i386/pci 下面的 Makefile 是处在 drivers/pci 下面的 Makefile 前面的，也就是说，pci_access_init 处在 acpi_pci_init 前面。

掌握了这些潜规则，我们在研究某个子系统时，就可以获得初始化函数的执行顺序，并按照该顺序使用韩峰同志对待日记的态度进行深入的分析。

《Linux内核修炼之道》精华分享与讨论

(18) ——选择发行版

学习内核首先要会使用它，依照一个由上至下循序渐进的过程，在能够熟练的使用Linux操

作系统之后再去研究内核中的实现。因此，了解并选择一个发行版进行安装使用便是一个不能回避的过程。

目前已经有超过 600 个 Linux 发行版，可以在 http://en.wikipedia.org/wiki/List_of_Linux_distributions 上看到它们的列表，其中，有多于 300 个正处于活跃的开发中，不断的改进。

对新人来说，在众多的发行版中选择一个是一件很让人头痛的事情，还是那句话：所有的痛苦都来自选择，所谓幸福，就是没有选择。这里介绍几种挑选的方法以供参考。

1. 参考发行版排行榜

在网站 <http://distrowatch.com/> 主页的右侧，依据页面点击次数给出了排名前 100 的 Linux 发行版，默认是最近 6 个月的点击次数，你可以选择数据统计的时间范围。这虽然不够科学，但也足够说明问题。



页面点击次数排名		
资料范围:		
Last 6 months		
刷新		
名次	发行	H.P.D*
1	Ubuntu	2180>
2	Fedora	1662>
3	Mint	1334>
4	openSUSE	1301>
5	Mandriva	989>
6	Debian	906>
7	Puppy	828>
8	Sabayon	751<
9	PCLinuxOS	726>
10	Arch	699>

2. Linux发行版比较服务

<http://polishlinux.org/choose/comparison/> 上提供的这个服务可以让你指定两个发行版，然后进行直观的比较，比较的项目非常细致，下面为选择Ubuntu和Mandriva时的部分显示结果。



3. Linux Distribution Chooser

<http://www.zegeniastudios.net/ldc/index.php?firsttime=true> 上提供了“Linux Distribution Chooser”服务通过问卷的形式帮助你挑选合适的Linux发行版。你只需要回答一些简单的问题，然后系统会给出几个符合你要求的发行版，最后还会给出一些不完全符合但也值得考虑的发行版。

[!\[\]\(99f58673407353e96a019fbca558fd72_img.jpg\) .gif" border="0" alt="发行版 03" width="244" height="195" />](#)

4. 听听别人的意见

通过上面的挑选，如果你已经决定了最终使用的那个发行版，那么恭喜你，你可以去下载或者购买该发行版的安装盘安装体验了。如果仍然在几个发行版之间徘徊，那么就要问问别人的意见，或者到网上查找一下相关的资料，到社区看看别人的评价。如果还是决定不了，那还是试用一下吧，相信并不会有什么损失。

《Linux内核修炼之道》精华分享与讨论

(19) ——不稳定的内核API

刚才欣闻在 SB 会 试运行期间，参观的上海市民情绪非常稳定的，很好很舒服的展示了自己的风采。于是我们要在这里要反思一下，为什么内核的 API 就不能同样的稳定？

开源社区正以极快的速度向内核中添加新功能，同时又在努力让修补 bug 的步伐跟上去，放慢开发速度看上去是不太可能的：首先 Linux 不能在技术上落后，否则就会失去要求越来越苛刻的商业用户；其次是因为 Linux 需要推动开发者社区的发展，不断增加新功能可以使开发者不感到厌倦，否则他们就可能转移到其它项目，另外也能在现有开发者年老或退出的时候吸引新人才。

在这样的快节奏下，内核开发人员一旦在当前的接口中找到 bug，或者更好的实现方式，他们就会很快的去修改当前的接口，这就意味着，函数名可能会改变，结构体可能被扩充或者删减，函数的参数也可能发生改变。一旦接口被修改，内核中使用这些接口的地方需要同时得到修正，这样才能保证所有的部分继续正常工作。

比如，内核 中的 USB 接口到目前为止至少经历了三次重写，解决了下面的问题：

把数据流从同步模式改成异步模式，这就减少了许多驱动程序的复杂度，提高了所有 USB 驱动程序的吞吐量(throughput)，结果就是几乎所有的 USB 设备都能以最大速率工作了。

修改了从 USB Core 中分配数据包内存的方式，以至于为了修正许多死锁问题，所有驱动都必须提供更多的 参数给 USB Core 代码

这和一些封闭源代码的操作系统形成鲜明的对比，在那些操作系统上，不得不额外的维护旧的 USB 接口。这就导致了一个可能性，新的开发者依然会不小心使用旧的接口，以不恰当的方式编写代码，进而影响到操作系统的稳定性。

在上面的例子中，所有的开发者都同意这些改动是重要的不得不进行的，在这样的情况下修改代价很低。如果 Linux 保持一个稳定的内核接口，那么就不得不创建一个新的接口，同时旧的有问题的接口也必须一直维护，这就会给 USB 开发者带来额外的工作。既然所有的 USB 开发者都是利用自己的时间工作，那么要求他们去做这些毫无意义的免费的额外工作，是不可能的。

安全问题对 Linux 来说是十分重要的，一个安全问题被发现，就会在非常短的短时间内得到修复。在很多情况下，这将导致内核中的一些接口被重写，以从根本上避免安全问题的发生。一旦内核接口被重写，所有使用这些接口的驱动程序，必须同时得到修正，以确定安全问题已经得到修复并且不可能在未来还有同样的安全问题。如果内核内部接口不允许改变，那么就不可能修复这样的安全问题，也不可能确认这样的安全问题以后不会发生。

开发者一直在清理内核接口。如果一个接口没有人在使用了，它就会被删除。这样可以确保内核尽可能的小，而且所有潜在的接口都会得到尽可能完整的测试（没有人使用的接口是不可能得到良好的测试的）。

《Linux内核修炼之道》精华分享与讨论

(20) ——学会使用Git

作为一名人民的好干部，如果希望被惦记，可以学我们的郑书记，将自己和蔼可亲的光辉形象搬上台 历；作为一名有梦想有追求而又不知道如何出名的人，你可以参考对岸的“超想被包养”社团。而作为一个内核爱好者，要想成为一名 内核开发者，为内核贡献自己的代码，我们必须能够与其他众多的内核开发者协同工作，这就意味着应该能够使用内核的版本控制工具 Git 管理内核代码。

1. 什么 是 Git

Git是Linux专门为内核而开发的一个开放源码的版本控制软件，如下图所示，

Git的主页<http://git-scm.com/>很好的回答了Git是什么的问题。



Git 是一款免费的、开源的、分布式的版本控制系统，旨在快速高效地处理无论规模大小的任何软件工程。每一个 Git clone（克隆）都是一个含有全部项目历史记录的文件仓库（repository），具有完整的版本修订追踪能力，它不依赖于网络连接或中心服务器，创建分支（branching）和合并分支（merging）非常的快速和简单。

目前，已经有越来越多的著名项目采用 Git 来管理项目的开发，比如 Perl、Gnome、Wine 等等。（注意，在 Ubuntu/debian 上，安装的是 git-core，而不是 git）

2. Git 的由来

Linus 于 2002 年 2 月 开始使用 BitKeeper 作为内核的版本控制工具。但是 BitMover 公司在商业版的 BitKeeper 之外，提供的 BitKeeper 只是仅可免费使用但不允许加以修改开放的 精简版，因此，包括 GNU 之父 Richard Stallman 在内的很多人，对 Linus 使用 BitKeeper 感到不满。

然而，当时 市场上并没有其他具备 BitKeeper 类似功能的自由软件可用，于是有些人就尝试对其进行逆向工程，这惹恼了 BitMover，该公司于是 决定停止提供 BitKeeper 的免费版本。为解决无工具可用的窘境，Linus 便自行开发 Git，希望在适当的工具出现前，暂时得充当解决方案。当时 Linus 曾称 Git 为愚蠢的内容管理器（the stupid content tracker）。当 Git 有了迅速成长之后，Linus 就建议能够以其作为长期的解决方案，并于之后的 2.6.12-rc3 内核第一次采用 Git 进行发布。

3. 一段 录像

在 Git 历史上有段很著名的录像，是 Linus 在 Google 的一个演讲，我们可以在 youtube 上 看到它。在这段录像中，Linus 说明了设计 Git 的原因，基本的设计哲学，以及与其他版本控制工具的比较。

从技术的观点上，Linus 非常尖锐的批判了 CVS 与 SVN。虽然 Linus 从来没有使用过 CVS 去管理内核代码，但是他在商业公司曾有过一段不短时间的使用经历，而且对其强烈的厌恶。同时 他批判 SVN 是毫无意义的，因为 SVN 尝试从各方面去改善 CVS 的一些缺点，却无法根本的解决一些基本的使用限制。具体来说就是，SVN 改善了创建分支的所耗费的成本，相对 CVS 利用了比较少的系统资源，但是 却无法解决合并分支的需求。但是许多项目的开发过程中，都时常需要为不同的新功能创建分支、合并分支，如此依赖，SVN 就成为一个没有未来的项目。

Git 作为一个分布式的版本控制工具，你可以随意的创建新分支，进行修改、测试、提交，这些在本地的提交完全不会影响到其他人，可以等到工作完成后再提交给公共的仓库。这样就可以支持离线工作，本地提交可以稍后提交到服务器上。

Linus 提到，在内核开发社区中有一种信任关系（web of trust），像 内核这样庞大的项目，每个版本参与的开发者都非常多，但是 Linus 不可能认识这么多的人，自然地，他只能信任最为熟悉的极少数人，并相信那些人的智商与能力是足以信赖的，于是他只需要信赖这些人的成果，而同时这些人又在自己的信

任圈力找到他可以信赖的人,于是利用这样的信赖机制扩展成了网状的内核开发社区。

实际上,社区中也会演化出几个角色,比如司令(dictator)、副官(lieutenants)、开发者,少数的副官只需要专注在他们熟悉的领域,整合开发者的成果,并提交给司令做最后的整合决策,这样一来,各种不同的领域都可以交给最为熟悉的开发者去管理,而项目开发本身不会被限制阻塞在某个角色身上,相对而言是一种比较高效的开发社区结构。

4. 一些 资源

本书并不会 也不需要 对 Git 的具体使用过程进行介绍,下面仅推荐几个好的网站或资料以供学习参考。

<http://www.ibm.com/developerworks/cn/linux/l-git/>

这篇文章详细描述了如何使用 Git 来管理内核代码。

<http://www.youtube.com/watch?v=4XpnKHJAok8>

这就是上面提到的那段著名的录像,相信认真的消化之后会有不小的收获。

<http://book.opensourceproject.org.cn/versioncontrol/git/gittutorcn.htm>

Git 的中文教程。

<http://zh-cn.whygitisbetterthanx.com/#github>

这是 Kanru 翻译的《为什么 Git 比 X 更好》，简要的说明了 Git 与 其他版本控制工具的比较，可以让你了解各种工具之间的差异细节。

<http://github.com/>

很多人说正是 GitHub 让他们选择了 Git，相比其他的项目托管网站，它更象一个社交网络，可以追踪别人的状态，不过追踪的不是朋友 发出的信息，而是朋友写出的代码。人们可以在 GitHub 上找到与他们 在做的事相关的其他开发人员或项目，然后轻松地 fork 和贡献，这样形 成了一个以 Git 和各种项目为中心的活跃社区。

《Linux内核修炼之道》精华分享与讨论

(21) ——二分法与printf()

人生就是一个茶几，上面摆满了杯具。内核也是一个大茶几，不过它上面的杯具是一个个的 bug。确定 bug 什么时候被引入是一个很关键的步骤， 在这个定位 bug 的过程中，不论有意或无意，都会很自然地用到二分查找的方法。

二分查找法的基本原理

对于二分查找法，我们不会也不应该会感到陌生。作为 一种高效的查找算法，它曾出现在我们的数据结构课堂里，出现在一次又一次的面试里，更是会频繁地应用在我们的代码里。在我们所接触到的各种算法里，它可以 说是最为大众化、最充满生活智慧的一个，很多人并不知道二分查找法的概念，却能够在生活中熟练的去应用。

比如，一个工人要维修一条 10km 长的电话线，首先他 需要定位出故障所在，如果沿着线路一小段一小段地查找，显然非常得困难，每查一个点都要爬一次电线杆，10km 长 的距离会有大约 200 多根电线杆。假设电线两端分别为 A、B，这时他会很自然地首先从中间的 C 开始查起，用话机 向两端测试时，发现 AC 段正常，故而断定故障在 BC 段， 再到 BC 段的中点 D，如果发现 BD 段正常，则故障在 CD 段，然后再到 CD 的中间点 E 查找，这样 每查一次，就可以把待查线路的 长度缩减一半，因而经过 7 次查找，就可以将故障发生的 范围缩小到 50 ~ 100m 左右，即在一两根电线杆附近。如此一来要 节省很多的精力与时间。

这是二分查找法在生活中的一个典型应用，实际上，查 找内核的 bug 与查找电话线的故障 相比，本质上都是相同的，并没有高深到哪里去，都是首先要定位出 故障的位置，然后去 解决它。

比如你在使用某个版本的内核时，发现了一个内核 bug， 这时你需要知道它究竟是在应用 哪个补丁时被引入的，如果一个一个的去还原那些补丁，每还原一个补丁就要测试一次内核， 那么必然会浪费过多的时间，而应用二 分查找法，首先确定一个肯定没有出现该 bug 的内 核版本，然后去测试位于这两个版本中间的那个版 本，这样重复筛选，就能够很容易的定 位出是从哪个版本开始出现了这个 bug。

printk()

printk()应该是每一个驱动开发者最为亲密的 伙伴了，我们常常将它与二分查找法结合在一 起寻找代码中发生问题的位置。

通常情况下，对于代码中的两个 printk()语句， 如果一个正常执行，而另一个没有被执行， 就说明问题发生在这两个 printk()之间，接下来就可 以在这个范围内应用二分查找法定位

有问题的代码。

1. printk()与 printf()

用户空间有 `printf()`，内核空间有 `printk()`，它们就如代表善与恶的命运双生子，即使长相功能如何的接近，都不能在代码中共存。

对于我们来说，最容易犯的错误是，在需要 `printk()` 的地方误用了 `printf()`，而在需要 `printf()` 的地方却又误用了 `printk()`，通常这都不会是因为不知道它们的区别，而只是习惯使然。

民间流传有这样的说法：当你在编写用户空间应用程序的时候，下意识写出的都是 `printk()`，那么就说明你是个标准的内核开发者了。

2. printk()的消息级别

`printk()`与 `printf()`的一个重要区别就是前者可以指定消息的打印级别，内核根据这个指定的级别来决定是否将消息打印到终端上。如下表所示，`printk()`共有 8 个级别。

级别	描述
KERN_EMERG	紧急情况，系统可能会崩溃
KERN_ALERT	必须立即响应
KERN_CRIT	临界情况
KERN_ERR	错误信息
KERN_WARNING	警告信息
KERN_NOTICE	普通的但可能需要注意的信息
KERN_INFO	提示性信息

KERN_DEBUG	调试信息
------------	------

如果没有指定消息的级别，`printk()`会使用默认的 `DEFAULT_MESSAGE_LOGLEVEL` (通常是 `KERN_WARNING`)。

3. 控制台的日志级别 (`console_loglevel`)

当 `printk` 指定的消息级别小于指定的控制台日志级别时，消息的内容就会显示在该控制台上。控制台的日志级别定义在 `include/linux/kernel.h` 文件中，默认为 `DEFAULT_CONSOLE_LOGLEVEL` (值等于 7)，也就是说默认情况下，比 `KERN_DEBUG` 级别高的 `printk()` 消息内容都可以在控制台上显示。

我们可以执行下面的命令使任何级别的 `printk()` 消息都被打印在终端上

```
$ echo 8 > /proc/sys/kernel/printk
```

4. `printk()` 的变体

内核在 `include/linux/kernel.h` 文件中提供了两个 `printk()` 的变体 `pr_debug` 和 `pr_info`，它们的定义为：

```
235 #define pr_debug(fmt,arg...) \
236     printk(KERN_DEBUG fmt,##arg)
244 #define pr_info(fmt,arg...) \
245     printk(KERN_INFO fmt,##arg)
```


5. printk()不是万能的

printk()虽然很好用，但它并不是万能的，在 系统启动时，终端还没有初始化之前，它并不能被使用，不过如果不是在调试系统的启动过程的话，这并不能算是个问题。

其实内核提供了一个 printk()的变 体 early_printk()，专门用于在系统启动的初期在终端上打印消息，它与 printk()的区别仅仅在于名字的不同以及它能够更早地工作。