



Linux 内核重建





本章内容

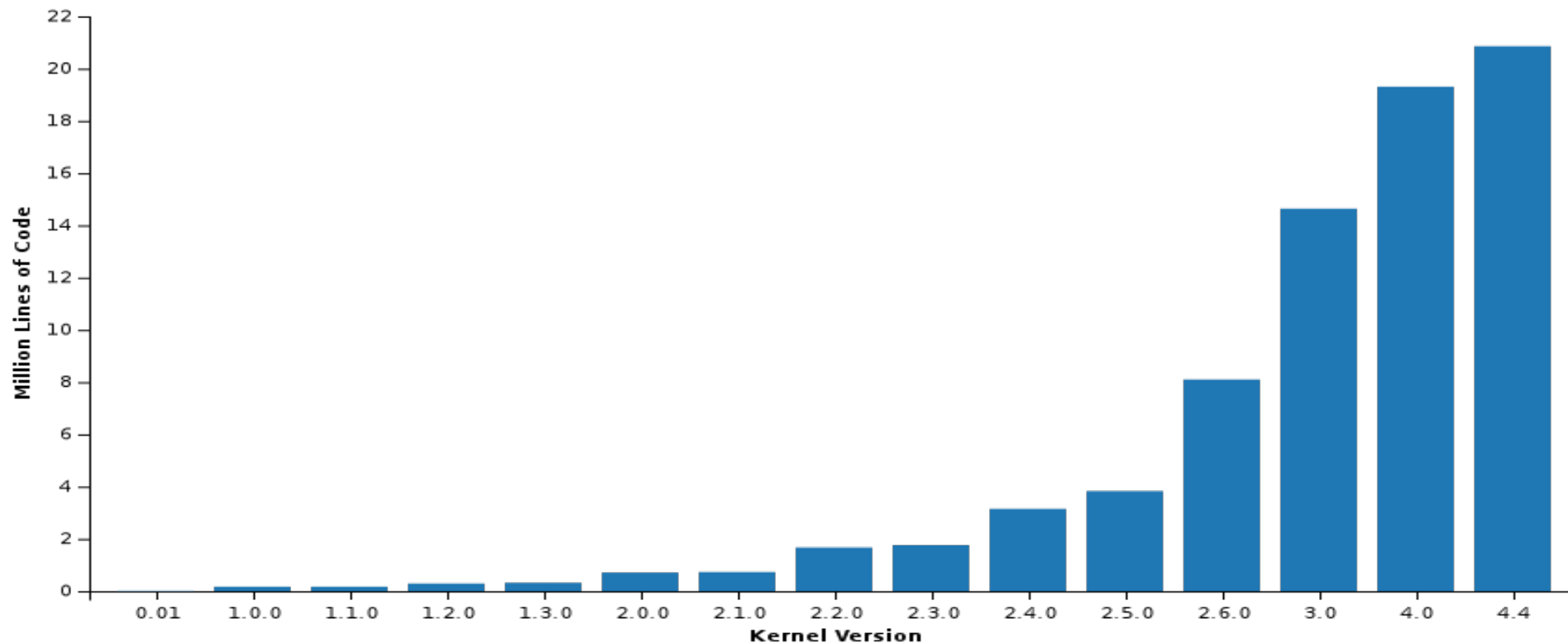
- Linux内核
- 编译Linux内核
- Linux启动*





一、Linux内核

- Linux内核近两年每**2个月左右**发布一个新版本, <http://www.kernel.org/>
 - kernel 4.19 共**29062477**行代码, 其中 “*.chS” 文件**26700224**行



在25年的时间内, Linux 内核已经从 10,000 多行代码增长到 近3千多万行代码。





Linux内核

- “Linux内核入门是不容易，它之所以难，在于庞大的规模和复杂的层面。规模一大，就不容易现出本来面目，浑然一体，自然不容易找到着手之处；层面一多，就会让人眼花缭乱，盘根错节，怎能让人提纲挈领”——《Linux内核设计与实现》译者序





Linux内核

- **Andre Morton**: 内核的学习曲线变得越来越长，也越来越陡峭。系统规模不断扩大，复杂程度不断提高。长此以往，虽然现在这一拨内核开发者对内核的掌握越发炉火纯青，但却会造成新手无法跟上内核发展步伐，出现青黄不接的断层。

Andrew Keith Paul Morton



Andrew Morton speaking at Interop, [Moscow](#), 2008

Born	1959 England
Other names	akpm
Education	Electrical engineering
Occupation	Programmer
Employer	Google
Known for	-mm tree
Spouse(s)	Kathryn Morton
Children	Victoria Morton, Micheal Morton, Matthew Morton

[http://en.wikipedia.org/wiki/Andrew_Morton_\(computer_programmer\)](http://en.wikipedia.org/wiki/Andrew_Morton_(computer_programmer))





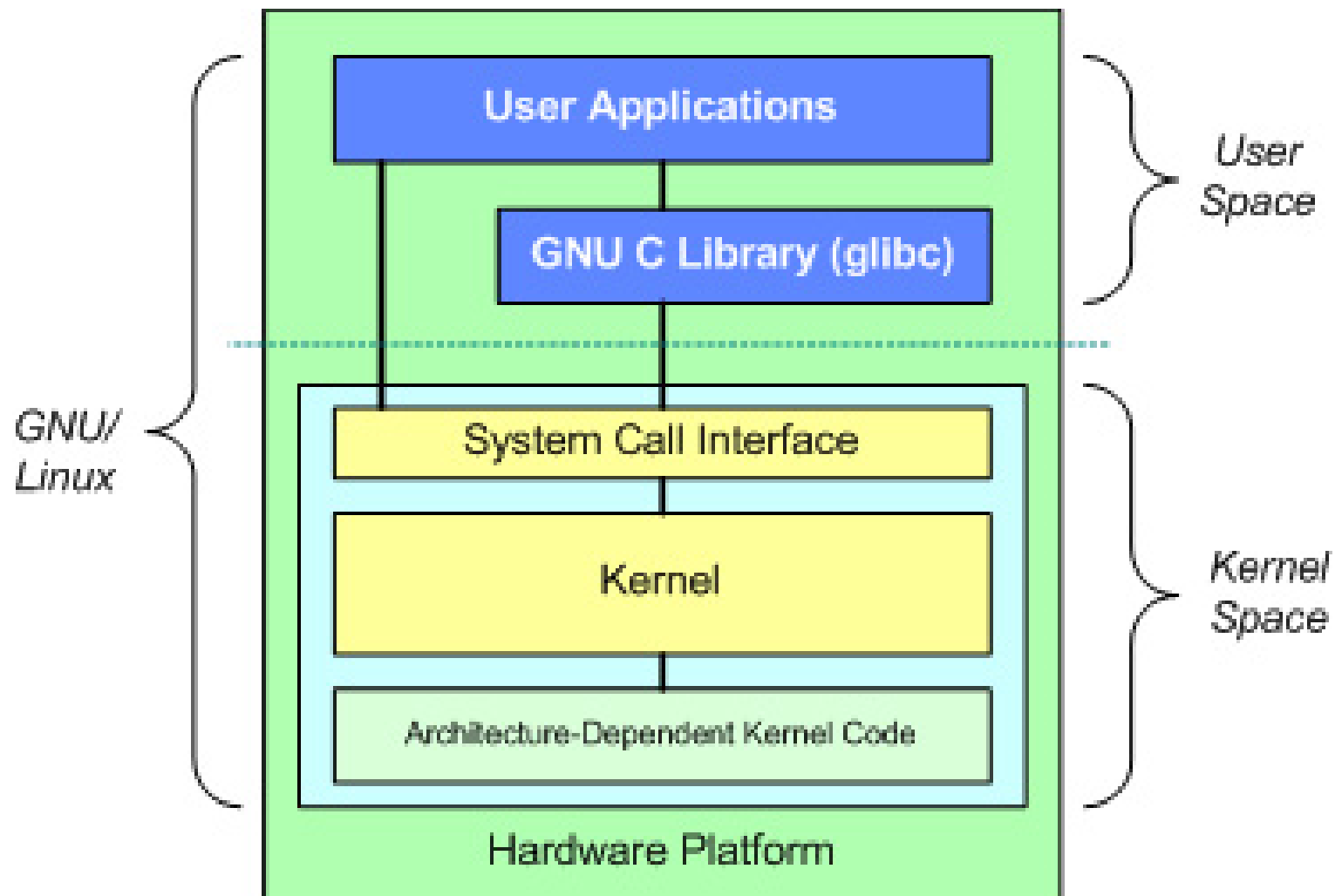
为什么要学习Linux内核

- 对学习操作系统原理有帮助外，学习Linux内核还有：
 - Linux内核是全世界最优秀的程序员写的
 - Linux内核结构非常好
 - 学习超大规模软件是如何设计的





GNU/Linux 操作系统基本层次结构



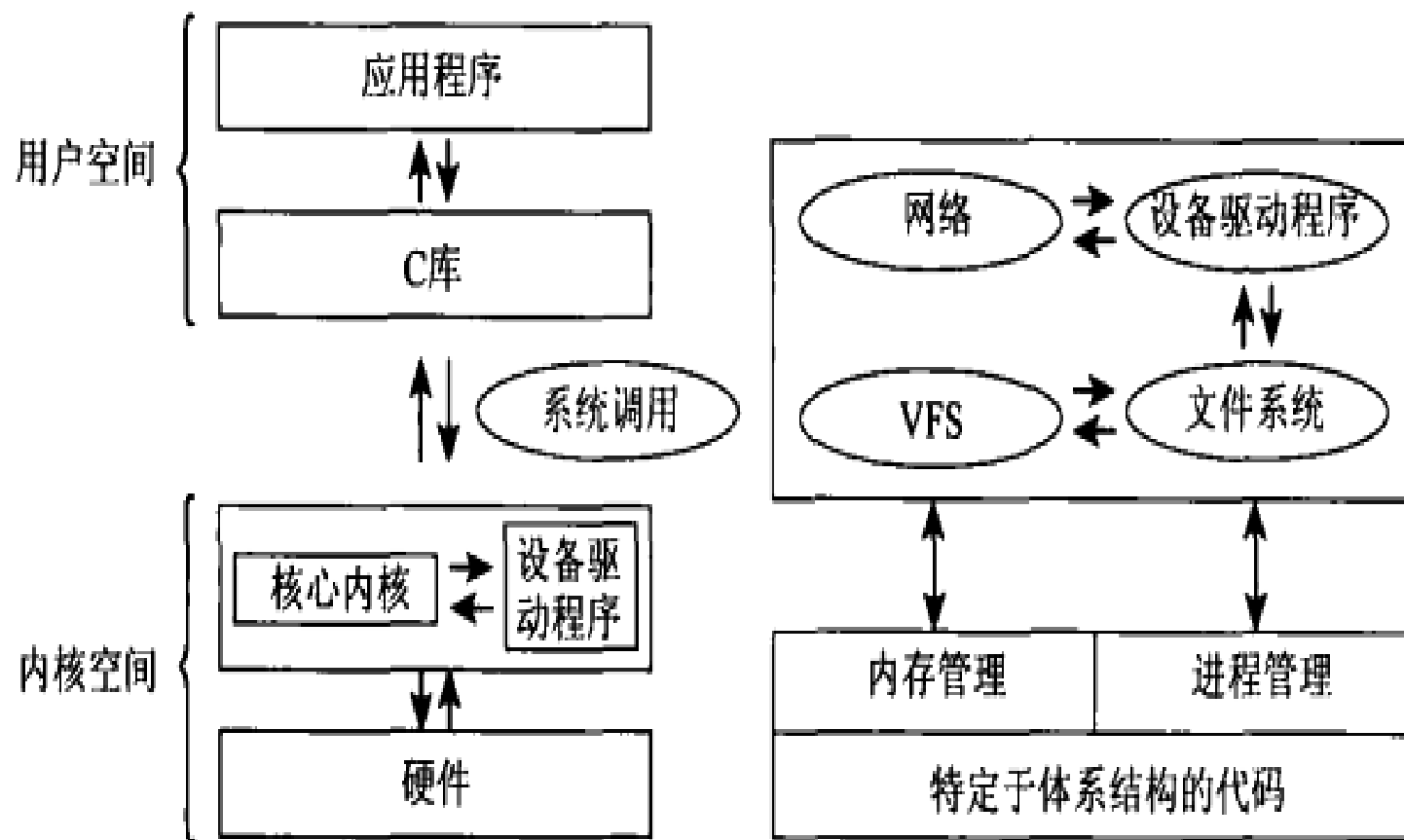


Linux系统层次结构

- **GNU C Library**（glibc）：它提供了连接内核的系统调用接口，还提供了在用户空间应用程序和内核之间进行转换的机制。
- **Linux内核**——系统调用接口、进程管理、内存管理、虚拟文件系统、网络堆栈、设备驱动程序、硬件架构的相关代码。



Linux 内核(续)



Linux 内核的体系结构图





Linux内核子系统(续)

- **系统调用接口**：SCI 层提供了某些机制执行从用户空间到内核的函数调用。它实现了一些基本的功能，例如 read 和 write。
- **进程调度**：控制着进程对CPU的访问。
- **内存管理**：允许多个进程安全地共享内存区域
- **虚拟文件系统（VFS）**：隐藏各种不同硬件的具体细节，为所有设备提供统一的接口。
- **网络**：提供了对各种网络标准协议的存取和各种网络硬件的支持。
- **进程间通信(IPC)**：支持进程间各种通信机制，包括共享内存、消息队列及管道等。





Linux内核(续)

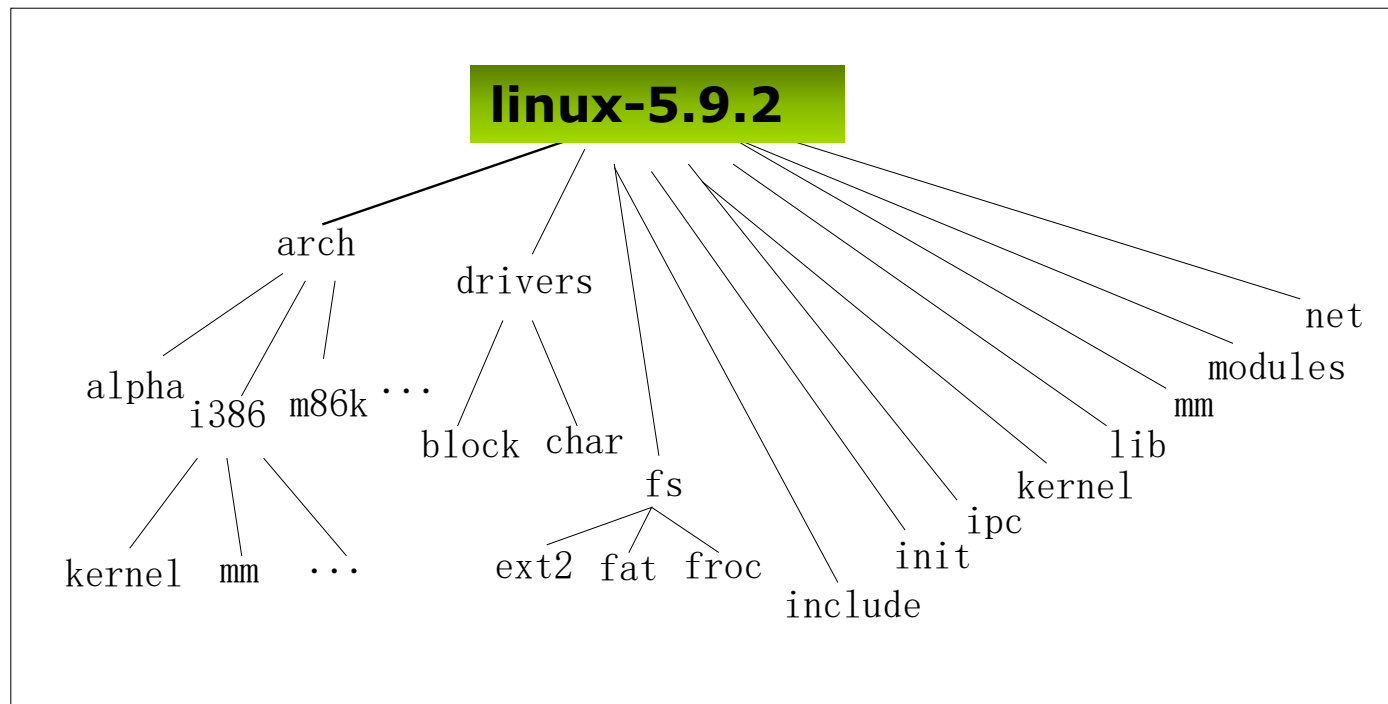
- 内核结构：
 - 单内核/宏内核(monolithic kernel): Unix、Linux
 - 微内核(Microkernel) : windows NT、Mac OS
- 从操作系统内核结构上看，Linux是一个单内核，Linux内核运行在单独的内核地址空间。
- Linux吸取了微内核的精华：其引以为豪的是模块化设计、抢占式内核、支持内核线程以及动态装载和卸载内核模块。





Linux内核源码目录结构

■ Linux内核代码可以安装在 `~/linux-x.x.x` 目录下



➤ Linux内核源代码包含在64312个C语言和汇编等文件中，5.9.2代码大约有2670万行，占用1.1G多空间。





Linux内核源码目录结构(续)

- **arch**: 该子目录包括了所有和体系结构相关的内核代码。它的每一个子目录都代表一种支持的体系结构, 例如i386就是关于intel cpu及与之相兼容体系结构的子目录。PC机一般都基于此目录。
- **Include**: 该子目录包括编译内核所需要的大部分头文件。与平台无关的头文件在include/linux子目录下, 与intel cpu相关的头文件在include/asm-i386子目录下, 而include/scsi目录则是有关scsi设备的头文件目录。
- **init**: 该子目录包含内核的初始化代码, 包含两个文件main.c和version.c。
- **mm**: 该子目录包括所有独立于cpu体系结构的内存管理代码, 如页式存储管理内存的分配和释放等; 而和体系结构相关的内存管理代码则位于arch/*/mm/, 例如arch/i386/mm/fault.c





Linux内核源码目录结构(续)

- **kernel**: 主要的内核代码, 此目录下的文件实现了大多数linux系统的内核函数, 其中最重要的文件当属sched.c; 同样, 和体系结构相关的代码在arch/*/kernel中。
- **drivers**: 放置系统所有的设备驱动程序; 每类驱动程序又各占用一个子目录: 如, /block下为块设备驱动程序, 比如ide (ide.c)。设备初始化程序在drivers/block/genhd.c中的device_setup()。
- **lib**: 放置内核的库代码。
- **net**: 内核与网络相关的代码。
- **ipc**: 这个目录包含内核的进程间通讯的代码。
- **fs**: 所有的文件系统代码和各种类型的文件操作代码, 它的每一个子目录支持一个文件系统, 例如fat和ext2;
- **scripts**: 此目录包含用于编译内核的脚本文件等。

在大多数目录下, 都有一个**Kconfig**文件和一个**Makefile**文件, 这两个文件都是编译时使用的辅助文件; 而且, 在有的目录下还有readme文件, 它是对该目录下的文件的一些说明, 同样有利于对内核源码的理解。





内核开发的特点

- 内核编程时不能访问C库；
- 内核编程时必须使用GNU C；
- 内核编程时缺乏像用户空间那样的内存保护机制；
- 内核编程时浮点数很难使用；
- 内核只有很小的定长堆栈；
- 由于内核支持异步中断、抢占式和SMP，因此必须时刻注意同步和并发；
- 要考虑可移植性的重要性





Linux内核源代码分析工具

- LXR (Linux Cross Reference) , 代码交叉检索工具。
 - <https://git.kernel.org/cgit/linux/kernel/git/stable/linux-stable.git/>
- Windows平台下的源代码阅读工具
 - 如: Source Insight
- Linux平台下的源代码阅读工具
 - 如: SourceNavigator





内核映像

- Linux系统引导过程使用内核映像，/boot目录下文件名称如：**vmlinuz-2.6.15.5**：
 - 普通内核映像：**zImage**（Image compressed with gzip），大小不能超过512k
 - 大内核映像：**bzImage**（big Image compressed with gzip），包含了大部分系统核心组件：系统初始化、进程调度、内核管理模块





内核构建

- 从Linux内核2.6开始，Linux内核的编译采用Kbuild系统，kbuild基于**GNU Make**，是一套非常复杂的系统。
- Linux内核的每个目录一个的Makefile，同它对应的有一个Kconfig，其内容是一些默认的编译选项。每一个子目录都有一个Kbuild的Makefile文件，用来执行从其上层目录传递下来的命令。
- Kbuild的Makefile并不直接被当做Makefile执行，而是从.config文件中提取信息，生成Kbuild完成内核编译所需的文件列表。





内核构建步骤

■ 第一次编译链接

- kbuild分别编译各个子目录下程序为目标文件，如built-in.o、lib.a等，然后将他们链接为**ELF格式的vmlinux**，并存放在顶层目录中。

■ 第二次编译链接

- kbuild使用工具objcopy，对ELF格式的vmlinux、piggy.S、piggy.o、head_32.c、misc.c进行处理、编译、封装、压缩等，在arch/x86/boot目录下，生成裸二进制格式的**vmlinux.bin**文件。





内核构建步骤

■ 第三次编译链接

- kbuild将arch/x86/boot下的a20.o、bioscall.o等目标文件链接为setup.elf，使用objcopy将其转换为裸二进制格式，并命名为setup.bin。

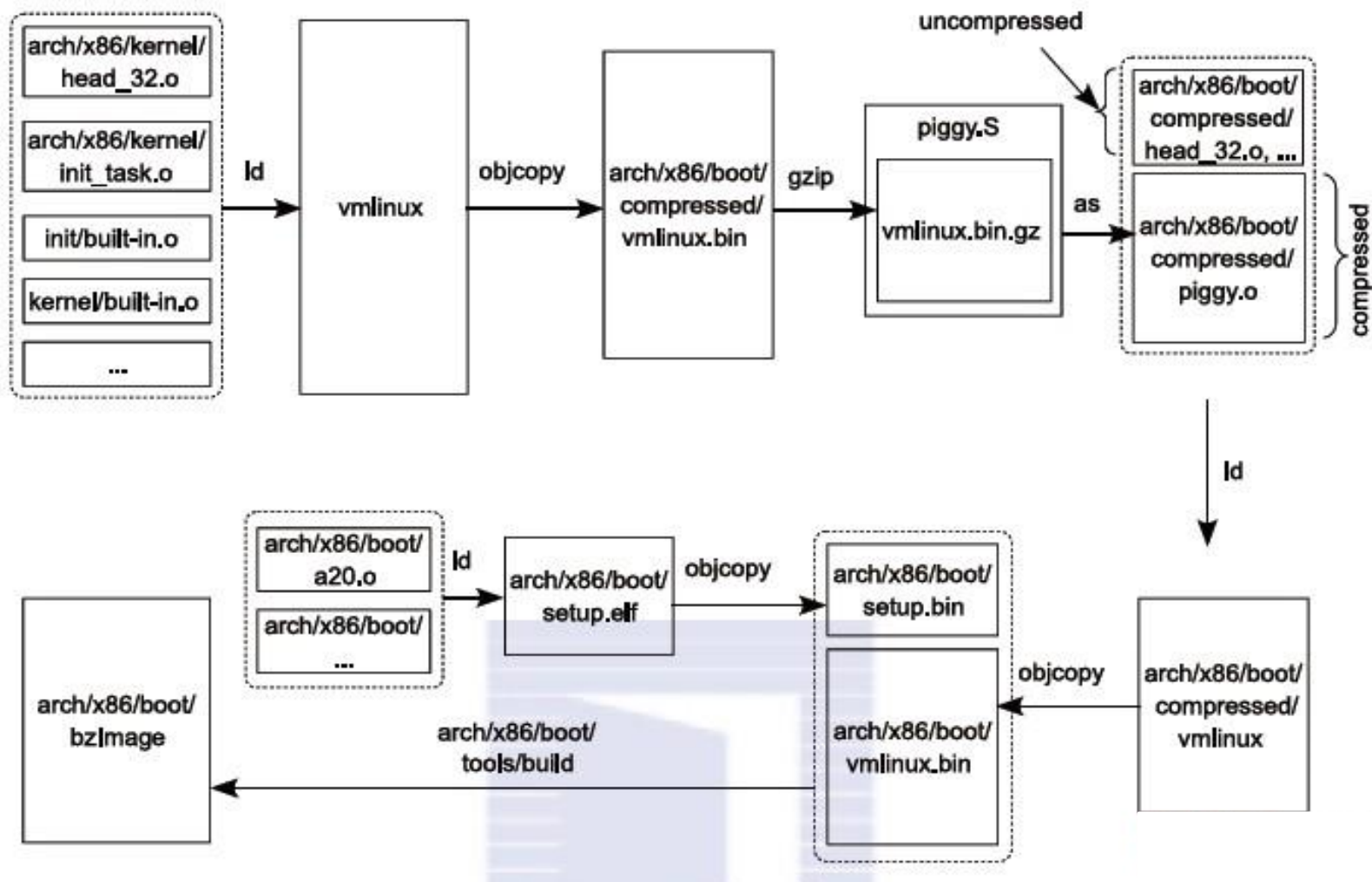
■ 组合

- kbuild调用内核自带的程序build，将vmlinux.bin和setup.bin合并为bzImage。





内核构建过程





Linux内核重建过程





(VirtualBox, ubuntu 19.10 , kernel 5.9.2)

重建内核步骤如下：

- 第一次要安装一些软件包：`sudo apt install libssl-dev libncurses5-dev`

/*编译内核和make menuconfig需要使用的软件包*/；也可能需要安装**make、gcc、bison、flex**等包

(1) 部署内核：

A. 官网内核

- 下载内核源代码：`linux-5.9.2.tar.xz` 文件。下载地址：
`http://www.kernel.org/` 或 `http://mirrors.aliyun.com/linux-kernel/`
- 部署内核代码
 - ▶ 把内核代码文件`linux-5.9.2.tar.xz` 存放在主目录（~）中。
 - ▶ `tar -xvf linux-5.9.2.tar.xz` /*解压内核包，生成的内核源代码放在linux.5.9.2 目录中*/





配置内核

B. 发行版本内核 (Ubuntu)

- `apt-cache search linux-source` // 查看内核版本
- `apt-get install linux-source-5.9.2` // 下载对应版本的内核，默认安装在 `/usr/src` 目录下





配置内核

(2) 配置内核:

- `cd linux-5.9.2`
- 如果需要将内核源代码树置于一种完整和一致的状态。命令`make mrproper`, 将清除目录下所有配置文件和先前生成核心时产生的.o文件:
 - `make mrproper`
- 用`uname -r`命令查看当前的内核版本号, 得知为5.3.0-59。把/usr/src/linux-headers-5.3.0-59-generic里面的.config文件复制到linux-5.9.2文件夹中, 使用系统的原配置文件:
 - `cp /usr/src/linux-headers-5.3.0-59-generic/.config .`
 - 或: `cp /usr/src/linux-headers-`uname -r`/.config .`





配置内核

■ `make menuconfig` // 同时生成.config文件

- 进行配置时，大部分选项可以使用其缺省值，只有小部分需要根据用户不同的需要选择。例如，如果硬盘分区采用ext2文件系统（或ext3文件系统），则配置项应支持ext2文件系统（ext3文件系统）。又例如，系统如果配有SCSI总线及设备，需要在配置中选择SCSI卡的支持。
- 对每一个配置选项，用户有三种选择，它们分别代表的含义如下：
 - ▶ “<*>”或 “[*]” — 将该功能编译进内核
 - ▶ “[]” — 不将该功能编译进内核
 - ▶ “[M]” — 将该功能编译成可以在需要时动态插入到内核中的模块
 - 与核心其它部分关系较远且不经常使用的部分功能代码编译成为可加载模块，有利于减小内核的长度，减小内核消耗的内存，简化该功能相应的环境改变时对内核的影响。许多功能都可以这样处理，例如像上面提到的对SCSI卡的支持，等等。





重建内核

(3) 编译内核

- `make` 或 `make -j2`

(4) 安装模块

- `sudo make modules_install`

(5) 安装内核

- `sudo make install`

在/boot目录下生成initrd.img-5.9.2、vmlinuz-5.9.2等启动文件

耐心等待

吃饭去咯!





启动

■ 查看启动选项

- `sudo gedit /boot/grub/grub.cfg`

■ 重新启动

- `sudo reboot` //启动时忽略错误信息提示

■ 启动后查看内核版本号

- `uname -r`

5.9.2

进一步学习内核构件过程：深度探索Linux操作系统系统构建和原理解析，第3章





四、Linux系统启动*

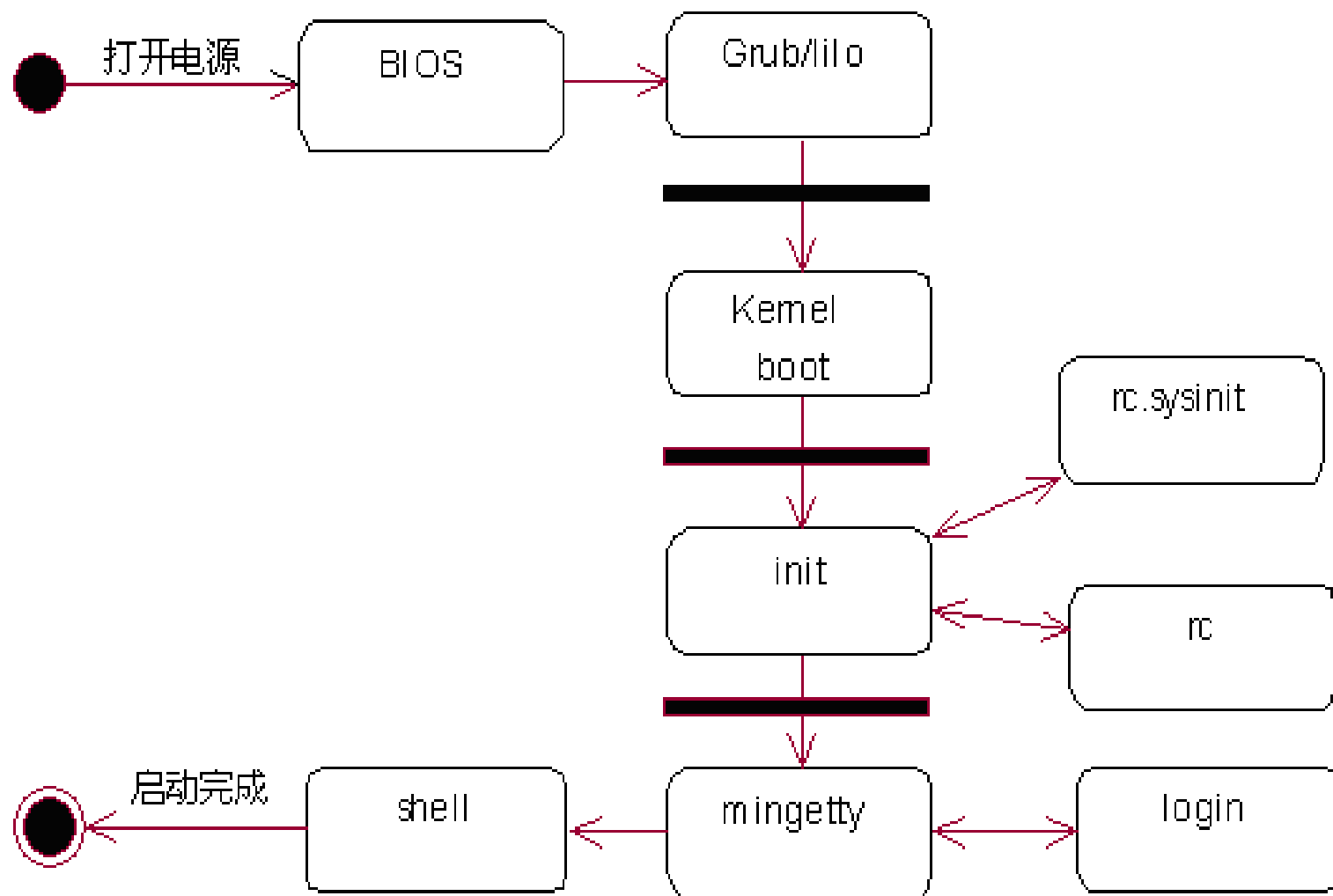
■ Linux系统的启动过程大致可分成以下几个阶段：

0. 硬件检测（自检）；
1. 由BIOS加载操作系统引导装入程序；
2. 由操作系统引导装入程序加载操作系统内核；
3. 内核代码解压缩；
4. 内核初始化；
5. 生成init进程；
6. 系统初始化，shell命令文本的执行；
 - ▶ 通过/etc/inittab文件进行初始化
 - ▶ 如设置键盘、字体、装载模块，设置网络等。
7. 生成各终端进程；
8. 用户登录。
 - ▶ 执行/bin/login程序





Linux系统启动



Linux系统启动过程





内核引导阶段

- bootsect
- setup
- head.S
- main.c





Init阶段

■ init有两项重要的工作：

- 解释/etc/inittab文件，定义runlevels
- 运行系统初始化脚本，/etc/rc.d目录定义各个runlevels上运行的服务





inittab文件的格式

■ inittab文件中每一行有下列格式

id:runlevel:action:process

- id: 配置记录标识符, tty为序号
- runlevel: 运行级
- action: process的执行方式
- process: 要运行的程序的路径和命令选项





8种执行方式

- sysinit: 提供init的路径
- respawn: 每当一个命令结束后, 就重启该命令
- askfirst: 类似respawn, 但是要先问一下用户
- wait: 阻塞式命令, init要等待其运行完毕
- once: 只运行一次, 不必等待
- ctrlaltdel: 三键齐按时, 要执行的命令
- shutdown: 系统关闭时执行
- restart: 系统重启时执行, 通常就是init





- 一个可能的inittab如下（id和runlevel都为空）：

```
::sysinit:/etc/init.d/rcS
```

设置/etc/init.d/rcS作为系统初始化文件

```
::respawn:/sbin/getty 115200 ttyS0
```

在串口（115200波特率）启动一个登录会话

```
::respawn:/control-module/bin/init
```

启动控制模块定制的系统初始化脚本

```
::restart:/sbin/initt
```

设置/sbin/init为重启时运行的命令

```
::shutdown:/bin/umount -a -r
```

系统关闭时，运行umount





init阶段工作

1、确定用户登录模式

- 在“`/etc/inittab`”中列出了如下所示的登录模式，主要有单人维护模式、多用户无网络模式、文字界面多用户模式、X-Windows多用户模式等。其中的单人维护模式（run level为1）是类似于Windows中的“安全模式”，在这种情况下，系统不加载复杂的模式从而使系统能够正常启动。

`cat /etc/inittab`

```
# Default runlevel. The runlevels used by RHS are:
# 0 - halt (Do NOT set initdefault to this)
# 1 - Single user mode
# 2 - Multiuser, without NFS (The same as 3, if you do not have networking)
# 3 - Full multiuser mode （文本界面启动模式）
# 4 - unused
# 5 - X11 （图形界面启动模式）
# 6 - reboot (Do NOT set initdefault to this)
#
id:5:initdefault:
.....
```





init阶段工作

2、执行脚本/etc/rc.d/rc.sysinit

- 将Linux的主机信息读入Linux系统，其内容就是文件“**/etc/rc.d/rc.sysinit**”中的。查看此文件可以看出，在这里确定了默认路径、主机名称、“/etc/sysconfig/network”中所记录的网络信息等。

```
# System initialization.
```

```
si::sysinit:/etc/rc.d/rc.sysinit
```





init阶段工作

3、启动内核的外挂模块及各运行级的脚本

- 读取模块加载配置文件 ([/etc/modules.conf](#))，以确认需要加载哪些模块。接下来会根据不同的运行级 (run level)，通过带参数 (运行级) 运行 “[/etc/rc.d/rc](#)”脚本，加载不同的模块，启动系统服务。init进程会等待 (wait) “[/etc/rc.d/rc](#)”脚本的返回。

10:0:wait:/etc/rc.d/rc 0

11:1:wait:/etc/rc.d/rc 1

12:2:wait:/etc/rc.d/rc 2

13:3:wait:/etc/rc.d/rc 3

14:4:wait:/etc/rc.d/rc 4

15:5:wait:/etc/rc.d/rc 5

16:6:wait:/etc/rc.d/rc 6





init阶段工作

4. 进入用户登录界面

- 系统还需要配置一些异常关机的处理部分。
- 最后通过“`/sbin/mingetty`”打开几个虚拟终端（tty1~tty6），用于用户登录。
- 如果运行级为5（图形界面启动），则运行`xdm`程序，给用户提供xdm图形界面的登录方式。





例：IA32 Linux* Boot

IA32 as an example

■ Power on CPU

- Registers will be set to default values
- cs:eip is initialized to 0xffffffff0

■ BIOS

- BIOS has code at 0xffffffff0(EEPROM or similar)
- Initialize devices
- Prepare tables/info for OS (E820, MPS ...)
- Copy boot device's first sector(MBR) to 0x7c00
- Jump to 0x7c00

■ Boot loader

- LILO/GRUB
- Boot loader loads kernel image and initrd
- Put kernel image into specific physical address
- Jump to kernel real mode initialization code





Linux Boot

■ Real mode initialization ([setup.S](#))

- Entry point is in 0x200 offset of kernel image
- Initialize memory (E820 ...)
- ...
- **Switch to protected mode**
- Copy kernel to 0x100000
- Jump to startup_32

■ Protected mode initialization([head.S::startup_32](#))

- Entry point is at physical address 0x100000
- Clear BSS
- Copy boot parameters to kernel data
- Initialize registers/GDT/IDT ...
- Jump to C code 'start_kernel'





Linux Boot

■ Start_kernel ([main.c](#))

- Setup_arch (architecture dependent setup,)
- Setup direct memory mapping/initialize buddy ...
- Kernel subsystem initialize/...
- Boot other CPUs
- Populate_rootfs (from initrd)
- Initcalls initialize
- Prepare files for init (0, 1, 2)
- Mount root fs
- Fork a new task and execute init process (/sbin/init) in the task





Linux Boot

- User space init (**init** process)
 - Read /etc/inittab （系统配置文件）
 - Execute init scripts (/etc/rc.d/rc*.d/, ...) per current run level
 - Init script will mount fs/start daemons/start x ...
 - Start getty process, open console
 - Type user name/password, login check them, and start shell
- Run level
 - 1 – single user mode
 - **3 – multiple user mode**
 - **5 – X mode**
 - 6 – reboot





作业

■ Linux 内核重建

