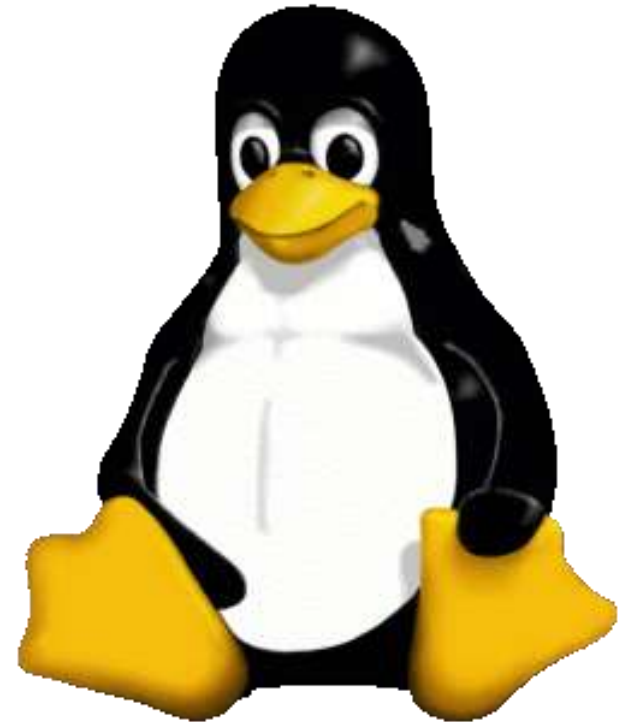# CS353 Linux Kernel

**Chentao Wu吴晨涛**
**Associate Professor**
**Dept. of CSE, SJTU**
**wuct@cs.sjtu.edu.cn**

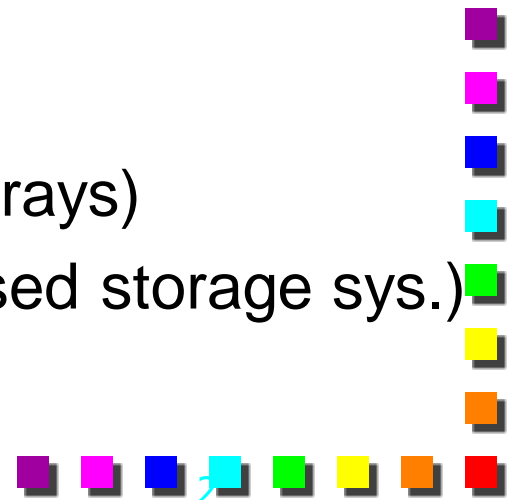# Introduction of Myself

- Dual Ph.D.
  - 2012, Electrical and Computer Engineering, *Virginia Commonwealth University (VCU)*, Richmond, VA, USA
  - 2010, Computer Architecture, *Huazhong University of Science and Technology (HUST)*, Wuhan, China

- Research Interest: Data Storage Systems
  - Storage management for Big Data
  - Cloud storage, Green storage
  - Reliable storage systems (e.g., disk arrays)
  - Semantic file systems (e.g., object-based storage sys.)
  - Cache Algorithms in storage systems

上海交通大學

# Our Lab

- Leader: Prof. Minyi Guo (Dean of CSE Dept.)
  - http://epcc.sjtu.edu.cn
  - Parallel and Distributed Computing
    - Parallel and Distributed Systems/Networks
    - High Performance Computing
    - Cloud Computing
    - Big Data
  - Welcome to participate in our lab
  - 15+ master students per year, 5+ doctoral students per year

# Download Lectures and Upload Projects

- **ftp://public.sjtu.edu.cn**
- User: wuct
- Password: wuct123456

# Teaching Assistant

- Chao Tan谭超
    - Email: 345243921@qq.com
    - Mobile Phone: 15821274485
- Help me to review Projects and the Final Exam

# **Books**

- No textbook
- References
  - Understanding the Linux Kernel
    - 3<sup>rd</sup> Edition
  - Linux Kernel Drivers
    - 3<sup>rd</sup> Edition
  - Linux Kernel Development
    - 3<sup>rd</sup> Edition

# Syllabus (1)

- Requirements:
  - Computer Organization, Operating System
  - C/C++ Programming
- Goals: Successful course participants will:
  - Understand C programming in Linux Kernel (Module programming in Linux Kernel).
  - Understand the core concepts of operating systems, including processes, threads, synchronization, virtual memory policies, and file management.

# Syllabus (2)

- Goals (contd.)
  - The idea of the course is <span style="color:red">to learn how computers really work in Linux Kernel</span>, from the chip level up to the application level. When we finish, you will understand what is actually happening when a computer system is running a set of programs, and will be able to make informed choices as a developer, project manager, or system customer.

# Course Meeting Times

- Lectures:
  - 2 classes per week (Friday)
  - 2 classes per dual weeks (Virtual Time)
- Questions:
  - Ask me directly between/after the classes
  - Go to my office: SEIEE 3-513
  - Send me an email: wuct@cs.sjtu.edu.cn
  - Ask teaching assistant

# Final Grades

- Participation 10%
  - Randomly
- Project 20%
  - Four Projects
  - Source Code and Document
- Report (Reading Kernel Code or Update Experimental Manual) 20%
- Final Exam 50% (Close Book)

# Late Policy

- Late Policy: Deadlines will be given in each assignment. These deadlines are strict.
- Typically, project will be given on Fridays, you should submit your homework by the next three weeks.

# Cooperation

- This course is established with Intel-Shanghai Corporation.

  - Several Intel engineers will share their experiences on Linux Kernel Programming.

# 1. The Linux Kernel: Introduction

**Chentao Wu**
**Associate Professor**
**Dept. of CSE, SJTU**
**wuct@cs.sjtu.edu.cn**

# What's Linux?

- **Linux** is a *Unix-like* and *POSIX-compliant* computer operating system assembled under the model of free and open source software development and distribution. The defining component of Linux is the **Linux kernel**, an operating system kernel first released on 5 October 1991 by **Linus Torvalds**.

  (from http://www.wikipedia.org)

# History

- UNIX: 1969 Thompson & Ritchie AT&T Bell Labs.

- BSD: 1978 Berkeley Software Distribution.

- Commercial Vendors: Sun, HP, IBM, SGI, DEC.

- GNU: 1984 Richard Stallman, FSF.

- POSIX: 1986 IEEE Portable Operating System unIX.

- Minix: 1987 Andy Tannenbaum.

- SVR4:  1989 AT&T and Sun.

- Linux: 1991 Linus Torvalds Intel 386 (i386).

- Open Source: GPL.

# UNIX Family

# Linux Features

- UNIX-like operating system.
- Features:
  - Preemptive multitasking.
  - Virtual memory (protected memory, paging).
  - Shared libraries.
  - Demand loading, dynamic kernel modules.
  - Shared copy-on-write executables.
  - TCP/IP networking.
  - SMP support.
  - Open source.

# Linux Distribution

# Ubuntu

# Two Modes within Linux

| | Application software (bash, LibreOffice, Blender, 0 A.D.) | | |
|---|---|---|---|
| User mode | Complex libraries (GLib, GTK+, Qt, SDL, EFL) | | Application software |
| | Complex libraries (GLib, GTK+, Qt) | Simple libraries (*sin, opendbm*) | Application software |
| | *open, exec, sbrk, socket, fopen, calloc*<br>C standard library: *glibc* (1,187,911 lines of code) / *uClibc* (342,842 lines of code) | | |
| Kernel mode | System calls: *TRAP, CALL, BRK, INT, IOCTL* (depends on the hardware) | | |
| | **Linux kernel** (16,223,920 lines of code)<br>(device drivers, process-scheduler, networking stack, file systems)<br>ALSA, DRI, evdev, LVM, device mapper, Linux Process Scheduler, Linux Network Scheduler,Netfilter<br>Linux Security Modules: *SELinux, TOMOYO, AppArmor, Smack* | | |
| | Hardware (CPU(s), Memory, other Microprocessors, Devices etc) | | |

# What's a Kernel?

- AKA: executive, system monitor.

- Controls and mediates access to hardware.

- Implements and supports fundamental abstractions:
  - Processes, files, devices etc.

- Schedules / allocates system resources:
  - Memory, CPU, disk, descriptors, etc.

- Enforces security and protection.

- Responds to user requests for service (system calls).

- Etc…etc…

# Kernel Design Goals

- Performance: efficiency, speed.
  - Utilize resources to capacity with low overhead.
- Stability: robustness, resilience.
  - Uptime, graceful degradation.
- Capability: features, flexibility, compatibility.
- Security, protection.
  - Protect users from each other & system from bad users.
- Portability.
- Extensibility.

# Example "Core" Kernel

| Applications |
|---|

| System Libraries (libc) |
|---|

**Modules**

| System Call Interface |
|---|

| I/O Related | Process Related |
|---|---|
| File Systems | Scheduler |
| Networking | Memory Management |
| Device Drivers | IPC |

| Architecture-Dependent Code |
|---|

| Hardware |
|---|

# Architectural Approaches

- Monolithic.

- Layered.

- Modularized.

- Micro-kernel.

- Virtual machine.

# Board View of Linux Kernel

# Various Types of Hardware Are Running on Linux

# Software Components of Linux Desktop

# History of Linux Kernel Versions

- 0.01 -- 17 Sep. 1991
- 1.0 -- 14 Mar. 1994
- 2.0 -- 9 June 1996
- 2.4 -- 4 Jan. 2001
- 2.6 -- 18 Dec. 2003
- 3.0 -- 21 July 2011
- 4.0 -- 12 Apr. 2015
- 4.5 – 21 Feb. 2016

# Lines of Source Code in Linux

# kernel.org

## The Linux Kernel Archives

About   Contact us   FAQ   Releases   Signatures   Site news

| Protocol | Location |
|---|---|
| HTTP | https://www.kernel.org/pub/ |
| GIT | https://git.kernel.org/ |
| RSYNC | rsync://rsync.kernel.org/pub/ |

**Latest Stable Kernel:**
⬇ **4.4.2**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| mainline: | 4.5-rc5 | 2016-02-20 | [tar.xz] | [pgp] | [patch] | | [view diff] | [browse] | |
| stable: | 4.4.2 | 2016-02-17 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| stable: | 4.3.6 [EOL] | 2016-02-19 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 4.1.18 | 2016-02-15 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.18.27 | 2016-02-15 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.14.61 | 2016-02-17 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.12.54 | 2016-02-15 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.10.97 | 2016-02-19 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.4.110 | 2015-10-22 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 3.2.77 | 2016-02-13 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| longterm: | 2.6.32.70 | 2016-01-29 | [tar.xz] | [pgp] | [patch] | [inc. patch] | [view diff] | [browse] | [changelog] |
| linux-next: | next-20160225 | 2016-02-25 | | | | | | [browse] | |

# Linux Source Tree Layout (2.6 or earlier)

```
                    ┌──────────────┐              ┌─────────┐
 ┌──────────────┐   │ /usr/src/linux│ ──────────► │ scripts │
 │ Documentation│ ◄─┤              │              └─────────┘
 └──────────────┘   └──────────────┘
```

Documentation    /usr/src/linux    scripts

init    ipc    kernel    net

arch    drivers    fs    include    lib    mm

| arch | drivers | fs | include | lib | net |
|------|---------|-----|---------|-----|-----|
| alpha | acorn | adfs | asm-alpha | adfs | 802 |
| arm | atm | affs | asm-arm | affs | appletalk |
| i386 | block | autofs | asm-generic | autofs | atm |
| ia64 | cdrom | autofs4 | asm-i386 | autofs4 | ax25 |
| m68k | char | bfs | asm-ia64 | bfs | bridge |
| mips | dio | code | asm-m68k | code | core |
| mips64 | fc4 | cramfs | asm-mips | cramfs | decnet |
| ppc | i2c | devfs | asm-mips64 | devfs | econet |
| s390 | i2o | devpts | linux | devpts | ethernet |
| sh | ide | efs | math-emu | efs | ipv4 |
| sparc | ieee1394 | ext2 | net | ext2 | ipv6 |
| sparc64 | isdn | fat | pcmcia | fat | ipx |
|  | macintosh | hfs | scsi | hfs | irda |
|  | misc | hpfs | video … | hpfs … | khttpd |
|  | net | … |  |  | lapb |
|  | … |  |  |  | … |
```

# Linux kernel map



| functions / layers | system | processing | memory | storage | networking | human interface |
|---|---|---|---|---|---|---|

**system** kernel/
**processing** kernel/
**memory** mm/
**storage** fs/
**networking** net/
**human interface**

## user space interfaces

**system interfaces**
linux/syscalls.h · system files
asm-i386/uaccess.h · /proc /dev /sysfs
copy_from_user
cdev — cdev_add
register_chrdev
cdev_map · sysfs_ops
sys_reboot · sys_init_module

**processes**
sys_execve · sys_fork · sys_vfork · sys_clone
fs/exec.c
linux_binfmt
sys_nanosleep

**memory access**
sys_brk
sys_mmap2
/proc/self/maps

**files & directories access**
sys_open
do_path_lookup · sys_read · sys_write
sys_mount
sys_sync

**sockets access**
sys_socketcall
sys_socket
socket_file_ops

**HI char devices**
cdev · kmsg
sys_syslog
input_fops · snd_fops
video_fops
console_fops
fb_fops
input

## virtual

drivers/base/ · kobject
**Device Model**
kset
subsystem_register
bus_type
class
device
device_create
class_device
class_device_create
device_driver
probe · driver_register

work_struct · **threads** · workqueue_struct
create_workqueue
kthread_create
kernel_thread
current
thread_info · do_fork

**Virtual memory**
virtually continues memory
find_vma_prepare · vmalloc
vmlist
vm_struct
virt_to_page

**Virtual File System**
file · vfs_read
vfs_write
inode · vfs_create
inode_operations · file_operations
file_system_type
get_sb
super_block
ramfs_fs_type

**protocol families**
sock_create — socket
inet_family_ops
inet_create · unix_family_ops
proto_ops
inet_dgram_ops · inet_stream_ops

security/ · linux/security.h
**security**
may_open
security_socket_create
security_inode_create
printk

## bridges

**Synchronization**
wake_up · mutex_lock · completion
add_timer · mutex · down
msleep · spin_lock

**Memory Mapping**
mm/memar.c
do_mmap_pgoff
vma_link
mm_struct
vm_area_struct

**Page Cache**
do_sync
address_space
pdflush

**Swap**
kswapd
do_swap_page
wakeup_kswapd

**networking storage**
nfs_file_operations
smb_fs_type
cifs_file_ops
iscsi_tcp_transport

security_ops
selinux_ops

© 2007 Constantine Shulyupin
www.LinuxDriver.co.il/kernel_map

## functional

init/ · kernel/ **system run**
kernel_restart · load_module
kernel_power_off
init/main.c · module
start_kernel
do_initcalls · run_init_process

kernel/sched.c
**Scheduler**
schedule_timeout
setup_timer
process_thread
activate_task — rq
task_struct
schedule
context_switch

mm/slab.c · /proc/slabinfo
**logical memory**
physically mapped memory
kfree · kmalloc
kmem_cache_alloc
kmem_cache
slab

**Logical File Systems**
ext3_file_operations
ext3_get_sb

**protocols**
proto · /proc/net/protocols
udp_prot · tcp_prot
ip_rcv
alloc_skb · ip_queue_xmit
ip_forward
sk_buff

**HI subsystems**
tty · log_buf
mousedev_handler
oss
alsa

drivers/ 

## devices control

**generic HW access**
pci_driver
usb_driver · pci_register_driver
usb_hcd_giveback_urb · pci_request_regions
request_region · usb_submit_urb
request_mem_region
ioremap · usb_hcd

**interrupt context**
timer_list
jiffies ++ · tasklet_struct
setup_irq
timer_interrupt · do_softirq
do_IRQ - irq_desc

mm/page_alloc.c · /proc/buddyinfo
**Page Allocator**
get_free_pages
zone · __alloc_pages
page
get_page_from_freelist
try_to_free_pages

**Block devices**
block/ · gendisk
block_device_operations
init_scsi · request_queue
scsi_device
scsi_driver
sd_fops
usb_storage_driver

**virtual network device**
net_device
netif_rx
dev_queue_xmit
alloc_etherdev · alloc_ieee80211
ether_setup · ieee80211_rx
ieee80211_xmit

drivers/input/ · drivers/media/ · sound/
**abstract devices and HID class drivers**
console
kbd
mousedev · fb_ops
video_device

## hardware interfaces

include/asm/ · arch/i386/
**devices access and bus drivers**
ehci_irq · pci_read
writew · pci_write
readw
outw · ehci_urb_enqueue
inw · usb_hcd_irq

arch/i386/kernel/
**CPU specific**
start_thread
/proc/interrupts · interrupt
atomic_t
system_call
show_regs · trap_init · switch_to
pt_regs
cli · sti

arch/i386/mm/
**physical memory operations**
die
do_page_fault

**disk controllers drivers**
Scsi_Host · libata
ahci_pci_driver
aic94xx_init · ide_intr
idedisk_ops
ide_do_request

drivers/net/
**network devices drivers**
e100_open · ipw2100_open
rtl8139_open · zd1201_net_open

**HI peripherals device drivers**
vga_con
atkbd_drv
psmouse
i8042_driver · ac97_driver
drivers/video/

## electronics

**I/O**
I/O mem · PCI controller
I/O ports · USB controller · DMA

**CPU**
registers · APIC · interrupt controller

**memory**
RAM · MMU

**disk controllers**
SCSI · IDE · SATA

**network controllers**
Ethernet · WiFi

**user peripherals**
keyboard · cam
mouse · audio
graphics card

# Linux Source Tree (1)

- All directories are grouped under the root entry "/"
- **root** - The home directory for the root user
- **home** - Contains the user's home directories along with directories for services
  - ftp
  - HTTP
  - samba

# Linux Source Tree (2)

- **bin** - Commands needed during booting up that might be needed by normal users

- **sbin** - Like bin but commands are not intended for normal users. Commands run by LINUX.

- **proc** - This filesystem is not on a disk. It is a virtual filesystem that exists in the kernels imagination which is memory

  - 1 - A directory with info about process number 1. Each process has a directory below proc.

# Linux Source Tree (3)

- **usr** - Contains all commands, libraries, man pages, games and static files for normal operation.
  - **bin** - Almost all user commands. some commands are in /bin or /usr/local/bin.
  - **sbin** - System admin commands not needed on the root filesystem. e.g., most server programs.
  - **include** - Header files for the C programming language. Should be below /user/lib for consistency.
  - **lib** - Unchanging data files for programs and subsystems
  - **local** - The place for locally installed software and other files.
  - **man** - Manual pages
  - **info** - Info documents
  - **doc** - Documentation
  - **tmp**
  - **X11R6** - The X windows system files. There is a directory similar to usr below this directory.
  - **X386** - Like X11R6 but for X11 release 5

# Linux Source Tree (4)

- **boot** - Files used by the bootstrap loader. Kernel images are often kept here.

- **lib** - Shared libraries needed by the programs on the root filesystem

- **modules** - Loadable kernel modules, especially those needed to boot the system after disasters.

- **dev** - Device files

- **etc** - Configuration files specific to the machine.

- **skel** - When a home directory is created it is initialized with files from this directory

- **sysconfig** - Files that configure the linux system for devices.

# Linux Source Tree (5)

- **var** - Contains files that change for mail, news, printers log files, man pages, temp files
    - **file**
    - **lib** - Files that change while the system is running normally
    - **local** - Variable data for programs installed in /usr/local.
    - **lock** - Lock files. Used by a program to indicate it is using a particular device or file
    - **log** - Log files from programs such as login and syslog which logs all logins and logouts.
    - **run** - Files that contain information about the system that is valid until the system is next booted
    - **spool** - Directories for mail, printer spools, news and other spooled work.
    - **tmp** - Temporary files that are large or need to exist for longer than they should in /tmp.
    - **catman** - A cache for man pages that are formatted on demand

# Linux Source Tree (6)

- **mnt** - Mount points for temporary mounts by the system administrator.

- **tmp** - Temporary files. Programs running after bootup should use /var/tmp

# linux/arch

- Subdirectories for each current port.
- Each contains **kernel**, **lib**, **mm**, **boot** and other directories whose contents override code stubs in architecture independent code.
- **lib** contains highly-optimized common utility routines such as memcpy, checksums, etc.
- **arch** as of 2.4:
    - alpha, arm, i386, ia64, m68k, mips, mips64.
    - ppc, s390, sh, sparc, sparc64.

# linux/drivers

- Largest amount of code in the kernel tree (~1.5M).
- device, bus, platform and general directories.
- drivers/char – n_tty.c is the default line discipline.
- drivers/block – elevator.c, genhd.c, linear.c, ll_rw_blk.c, raidN.c.
- drivers/net –specific drivers and general routines Space.c and net_init.c.
- drivers/scsi – scsi_*.c files are generic; sd.c (disk), sr.c (CD-ROM), st.c (tape), sg.c (generic).
- General:
  - cdrom, ide, isdn, parport, pcmcia, pnp, sound, telephony, video.
- Buses – fc4, i2c, nubus, pci, sbus, tc, usb.
- Platforms – acorn, macintosh, s390, sgi.

# linux/fs

- Contains:
  - virtual filesystem (VFS) framework.
  - subdirectories for actual filesystems.
- vfs-related files:
  - exec.c, binfmt_*.c - files for mapping new process images.
  - devices.c, blk_dev.c – device registration, block device support.
  - super.c, filesystems.c.
  - inode.c, dcache.c, namei.c, buffer.c, file_table.c.
  - open.c, read_write.c, select.c, pipe.c, fifo.c.
  - fcntl.c, ioctl.c, locks.c, dquot.c, stat.c.

# linux/include

- include/asm-*:
    - Architecture-dependent include subdirectories.
- include/linux:
    - Header info needed both by the kernel and user apps.
    - Usually linked to /usr/include/linux.
    - Kernel-only portions guarded by #ifdefs
        - #ifdef __KERNEL__
        - /* kernel stuff */
        - #endif
- Other directories:
    - math-emu, net, pcmcia, scsi, video.

# linux/init

- Just two files: version.c, main.c.

- version.c – contains the version banner that prints at boot.

- main.c – architecture-independent boot code.

- start_kernel is the primary entry point.

# linux/ipc

- System V IPC facilities.
- If disabled at compile-time, util.c exports stubs that simply return –ENOSYS.
- One file for each facility:
  - sem.c – semaphores.
  - shm.c – shared memory.
  - msg.c – message queues.

# linux/kernel

- The core kernel code.
- sched.c – "the main kernel file":
    - scheduler, wait queues, timers, alarms, task queues.
- Process control:
    - fork.c, exec.c, signal.c, exit.c etc…
- Kernel module support:
    - kmod.c, ksyms.c, module.c.
- Other operations:
    - time.c, resource.c, dma.c, softirq.c, itimer.c.
    - printk.c, info.c, panic.c, sysctl.c, sys.c.

# linux/lib

- kernel code cannot call standard C library routines.
- Files:
    - brlock.c – "Big Reader" spinlocks.
    - cmdline.c – kernel command line parsing routines.
    - errno.c – global definition of errno.
    - inflate.c – "gunzip" part of gzip.c used during boot.
    - string.c – portable string code.
        - Usually replaced by optimized, architecture-dependent routines.
    - vsprintf.c – libc replacement.

# linux/mm

- Paging and swapping:
    - swap.c, swapfile.c (paging devices), swap_state.c (cache).
    - vmscan.c – paging policies, kswapd.
    - page_io.c – low-level page transfer.
- Allocation and deallocation:
    - slab.c – slab allocator.
    - page_alloc.c – page-based allocator.
    - vmalloc.c – kernel virtual-memory allocator.
- Memory mapping:
    - memory.c – paging, fault-handling, page table code.
    - filemap.c – file mapping.
    - mmap.c, mremap.c, mlock.c, mprotect.c.

# linux/scripts

- Scripts for:
    - Menu-based kernel configuration.
    - Kernel patching.
    - Generating kernel documentation.

# Summary

- Linux is a modular, UNIX-like monolithic kernel.

- Kernel is the heart of the OS that executes with special hardware permission (kernel mode).

- "Core kernel" provides framework, data structures, support for drivers, modules, subsystems.

- Architecture dependent source sub-trees live in /arch.

# Project 1:
# Compile the Linux Kernel

# Preparation

- **For each student**
  - **Prepare one laptop for each group**
  - **Install Virtual Machine Software: VMware Workstation or VirtualBox**
  - **Install the Linux Operating System in a virtual machine: RedHat Enterprise Linux (RHEL), Fedora, Gentoo, Ubuntu, etc.**
  - **Linux kernel version: 4.3 or earlier**

# Compile The Linux Kernel (1)

- Download the latest kernel version (4.4.2, updated on 2016-02-17) from **http://www.kernel.org**
    - Download the complete source code (not a patch or change log)
- Copy to a directory
    - **4.X.X is the kernel version number**

    **# cp linux-4.X.X.tar.gz  /usr/src**

    **# cd /usr/src**
- Extract the kernel

    **# tar –xzvf linux-4.X.X.tar.gz  or**

    **# tar –xjvf linux-4.X.X.tar.bz2**

# Compile The Linux Kernel (2)

■ Preparation

**# cd /usr/src/linux**

**# make clean**

**# make mrproper**

**Note: make clean and make mrproper can be passed with the first compilation.**

# Compile The Linux Kernel (3)

- Configure the kernel

  **# make menuconfig**

  /*creates a menu where you can browse options on what the kernel supports*/

  **# make xconfig**

  **/***same as menuconfig, except that now the configuration menu is graphics based**/**

  **# make oldconfig**

  /*minor revision on previous kernel*/

  **# make config** **/*not be recommended*/**

# Make menuconfig

# Make Xconfig



**Linux Kernel Configuration**

| | | |
|---|---|---|
| Code maturity level options | SCSI support | Console drivers |
| Loadable module support | IEEE 1394 (FireWire) support | Sound |
| Processor type and features | I2O device support | USB support |
| General setup | Network device support | Bluetooth support |
| Memory Technology Devices (MTD) | Amateur Radio support | Kernel hacking |
| Parallel port support | IrDA (infrared) support | |
| Plug and Play configuration | ISDN subsystem | |
| Block devices | Old CD-ROM drivers (not SCSI, not IDE) | |
| Multi-device support (RAID and LVM) | Input core support | Save and Exit |
| Networking options | Character devices | Quit Without Saving |
| Telephony Support | Multimedia devices | Load Configuration from File |
| ATA/IDE/MFM/RLL support | File systems | Store Configuration to File |

# Compile The Linux Kernel (4)

- Customize the Kernel
  - e.g., you may add support for NTFS file system from "File System >> DOS/FAT/NT/ >> select NTFS file system support"

# Compile The Linux Kernel (5)

○ Compile and Install the Kernel (~20 min)

**# make**

**# make modules_install**

**# make install**

○ Update GRUB (or LILO)

**# cd /boot**

**# mkinitrd –o initrd.img-<kernel version>**

○ Reboot the system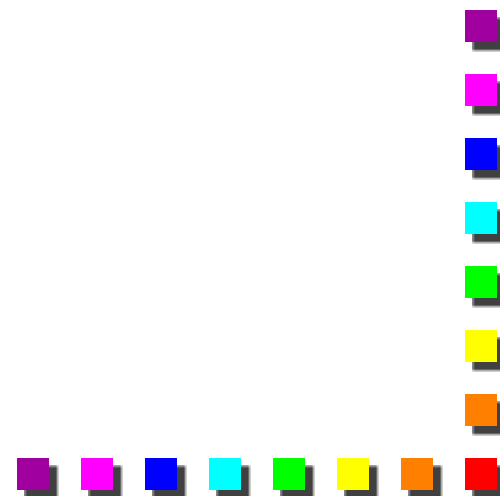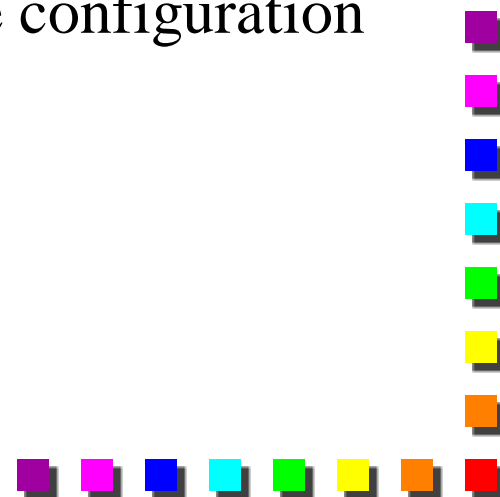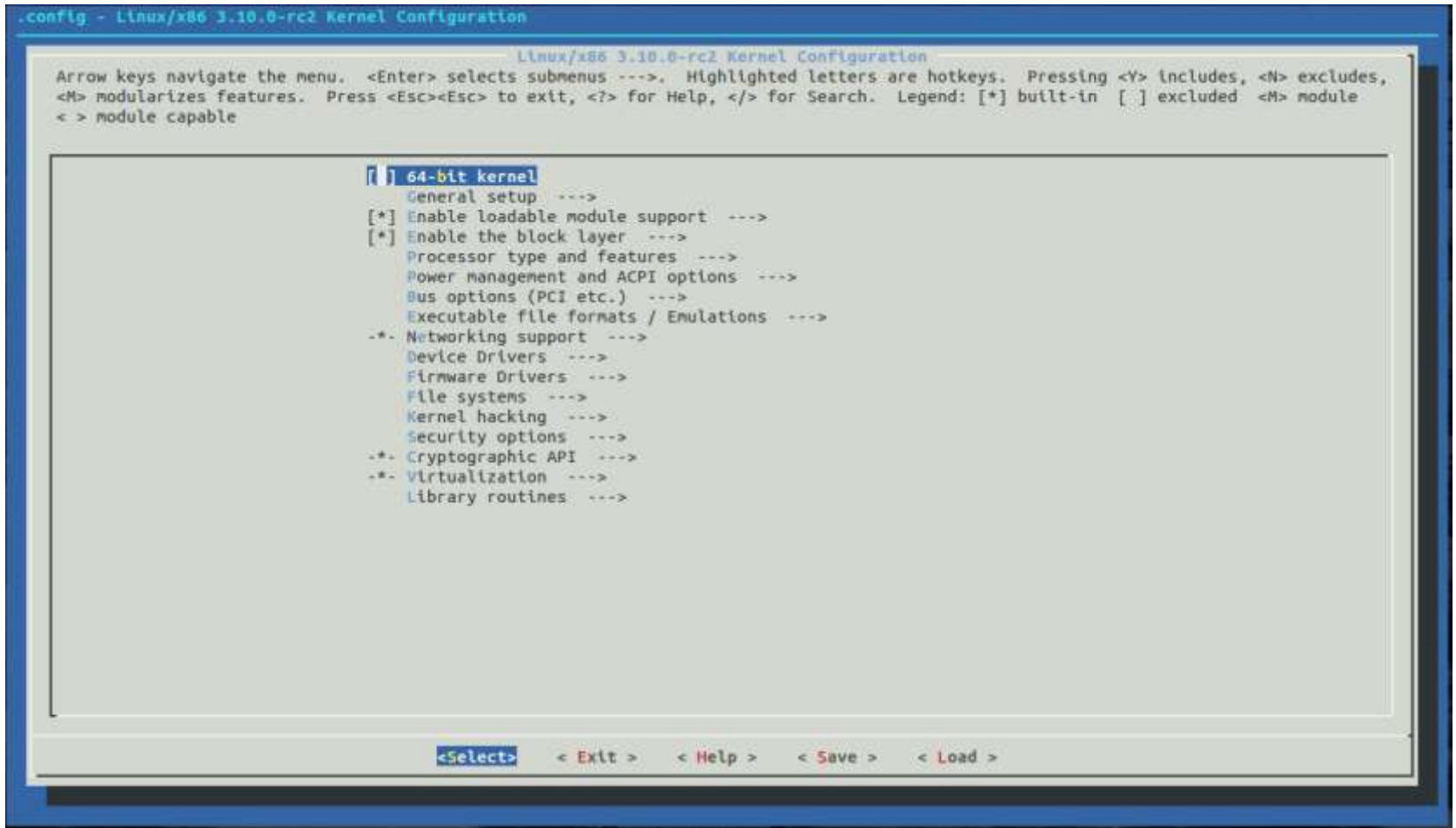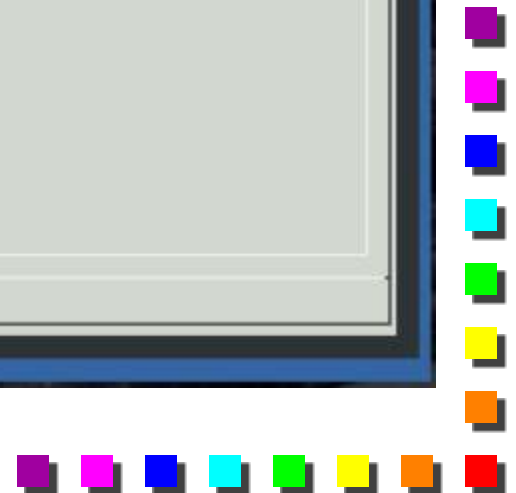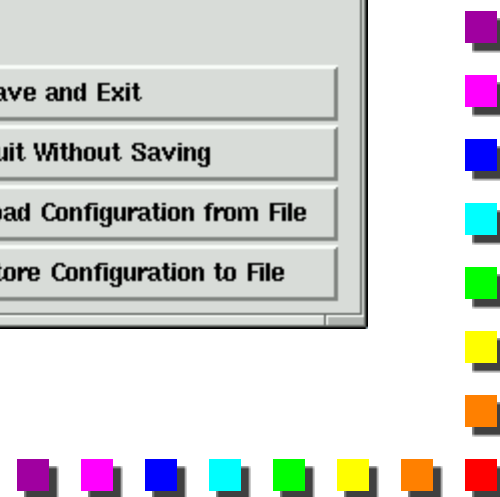