# Linux Playground and Storage Stack
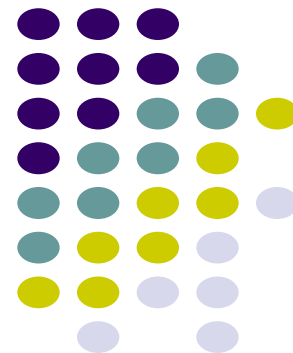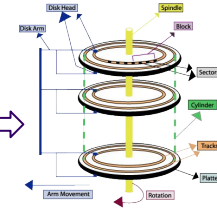
ECE 469, March 31
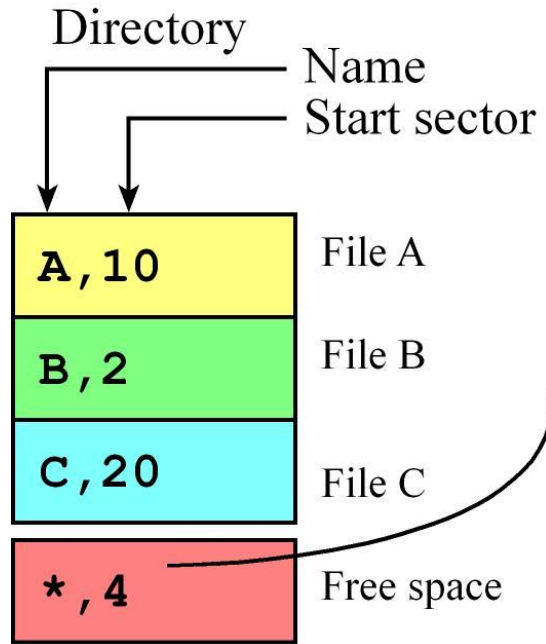
Aravind Machiry

# Preview: File System Abstraction



File System Abstraction

# Preview: File Allocation Table (FAT)

- Simple.
- Easy to implement.
- Still used in Phones and Thumb drives.

- Key data structure: File Allocation Table
  - List of all disk blocks.
  - File: Linked list of blocks.

File Allocation Table

Disk

Directory

Name
Start sector

| A,10 | File A |

| B,2 | File B |

| C,20 | File C |

| *,4 | Free space |

| 0 | x |
| 1 | x |
| 2 | 11 |
| 3 | 12 |
| 4 | 5 |
| 5 | 8 |
| 6 | 0 |
| 7 | 6 |
| 8 | 9 |
| 9 | 15 |
| 10 | 3 |
| 11 | 19 |
| 12 | 0 |
| 13 | 14 |
| 14 | 7 |
| 15 | 16 |
| 16 | 17 |
| 17 | 21 |
| 18 | 0 |
| 19 | 13 |
| 20 | 28 |
| 21 | 22 |
| 22 | 23 |
| 23 | 24 |
| 24 | 25 |
| 25 | 29 |
| 26 | 18 |
| 27 | 26 |
| 28 | 27 |
| 29 | 30 |
| 30 | 31 |
| 31 | 0 |

3

# Preview: FAT

- Pros: Simple
  - Easy to find free block.
  - Easy to append file.
  - Easy to delete a file.

- Cons:
  - Small file access is slow.
  - Random file access is very slow.
  - Fragmentation:
    - Blocks of a small file could be heavily scattered.
    - Problem becomes worse as the usage increases.

# Preview: EXT2 File System

# Preview: EXT2 File System : inode

# Preview: EXT2 File Size (Block Size: 4K)

12 File Block Pointers = 12*4 = 48K

1 Indirect block pointer (4K) = 1K direct block pointers = 1K*4K = 4MB

1 doubly indirect block pointer (4K) = 1K Indirect block pointers = 1K*4MB = 4 GB

1 triply indirect block pointer (4K) = 1K doubly Indirect block pointers = 1K*4GB = 4 TB

Total Size = 48K + 4MB + 4GB + 4TB

# Linux Kernel PlayGround

- [https://github.com/purs3lab/linux-playground](https://github.com/purs3lab/linux-playground)
- Linux kernel debugging through Docker container and VSCode

∨ VARIABLES
∨ Locals
  err: <optimized out>
  tsk: <optimized out>
> stack: <optimized out>
> stack_vm_area: <optimized out>
  node: -1
> orig: 0xffff888003520000
∨ Registers
> CPU
> Segs
> FPU
> Other Registers
> SSE

∨ WATCH

∨ CALL STACK          PAUSED ON BREAKPOINT
dup_task_struct(int node, struct task_struct * orig)     fork.c   881:1
copy_process(struct pid * pid, int trace, int node, struct kernel_clone_
kernel_clone(struct kernel_clone_args * args)            fork.c   2565:1
__do_sys_clone(unsigned long clone_flags, unsigned long newsp, int * par
do_syscall_x64(int nr, struct pt_regs * regs)            common.c   50:1
do_syscall_64(struct pt_regs * regs, int nr)             common.c   80:1
entry_SYSCALL_64()                                       entry_64.S  113:1
[Unknown/Just-In-Time compiled code]                     Unknown Source  0

∨ BREAKPOINTS
☐ All C++ Exceptions
☑ fork.c  kernel                                         874
☑ fork.c  kernel                                         874

kernel > C fork.c > ⊕ dup_task_struct(task_struct *, int)

```c
868        unsigned long *stackend;
869
870        stackend = end_of_stack(tsk);
871        *stackend = STACK_END_MAGIC;     /* for overflow detection */
872 }
873
874 static struct task_struct *dup_task_struct(struct task_struct *orig, int node)
875 {
876        struct task_struct *tsk;
877        unsigned long *stack;
878        struct vm_struct *stack_vm_area __maybe_unused;
879        int err;
880
881        if (node == NUMA_NO_NODE)
882                node = tsk_fork_get_node(orig);
883        tsk = alloc_task_struct_node(node);
884        if (!tsk)
885                return NULL;
886
887        stack = alloc_thread_stack_node(tsk, node);
888        if (!stack)
889                goto free_tsk;
890
891        if (memcg_charge_kernel_stack(tsk))
892                goto free_stack;
893
894        stack_vm_area = task_stack_vm_area(tsk);
```

PROBLEMS 250    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS 2

```
[    2.588740] PM:    Magic number: 14:562:333
[    2.589157] tty tty54: hash matches
[    2.589863] printk: console [netcon0] enabled
[    2.590067] netconsole: network logging started
[    2.593440] cfg80211: Loading compiled-in X.509 certificates for regulatory database
[    2.627515] tsc: Refined TSC clocksource calibration: 1607.892 MHz
[    2.627963] clocksource: tsc: mask: 0xffffffffffffffff max_cycles: 0x172d4430b0c, max_idle_ns: 440795221797 ns
[    2.628509] clocksource: Switched to clocksource tsc
[    2.671384] modprobe (64) used greatest stack depth: 14488 bytes left
[    2.761915] cfg80211: Loaded X.509 cert 'sforshee: 00b28ddf47aef9cea7'
[    2.764809] platform regulatory.0: Direct firmware load for regulatory.db failed with error -2
[    2.765642] cfg80211: failed to load regulatory.db
[    2.767459] ALSA device list:
[    2.767658]   No soundcards found.
[    2.784558] Freeing unused kernel image (initmem) memory: 1252K
[    2.785062] Write protecting the kernel read-only data: 22528k
[    2.789961] Freeing unused kernel image (text/rodata gap) memory: 2032K
[    2.792105] Freeing unused kernel image (rodata/data gap) memory: 1176K
[    2.793031] Run /init as init process
Boot took 2.79 seconds


$$$$$$$$\  $$$$$$\  $$$$$$$$\        $$\   $$\  $$$$$$\   $$$$$$\
$$  _____|$$  __$$\ $$  _____|       $$ |  $$ |$$  __$$\ $$  __$$\
$$ |      $$ /  \__|$$ |             $$ |  $$ |$$ /  \__|$$ /  $$ |
$$$$$\    $$ |      $$$$$\           $$$$$$$$ |\$$$$$$\  \$$$$$$$ |
$$  __|   $$ |      $$  __|          \_____$$ | \____$$\  \____$$ |
$$ |      $$ |  $$\ $$ |                   $$ |$$\   $$ |$$\   $$ |
$$$$$$$$\ \$$$$$$  |$$$$$$$$\               $$ | \$$$$$$  |\$$$$$$  |
_____| _____/ _____|             \__|  _____/  _____/


                 Operating Systems Engineering


Welcome to the Linux Kernel Playground
Type CTRL-A x to exit QEMU

/bin/sh: can't access tty; job control turned off
/ # [    3.158838] input: ImExPS/2 Generic Explorer Mouse as /devices/platform/i8042/serio1/input/input3
ls
bin     etc     linuxrc  root     sys
dev     init    proc     sbin     usr
/ # ls
```

bash sources
(Advanced) Use vscode to debug  Task
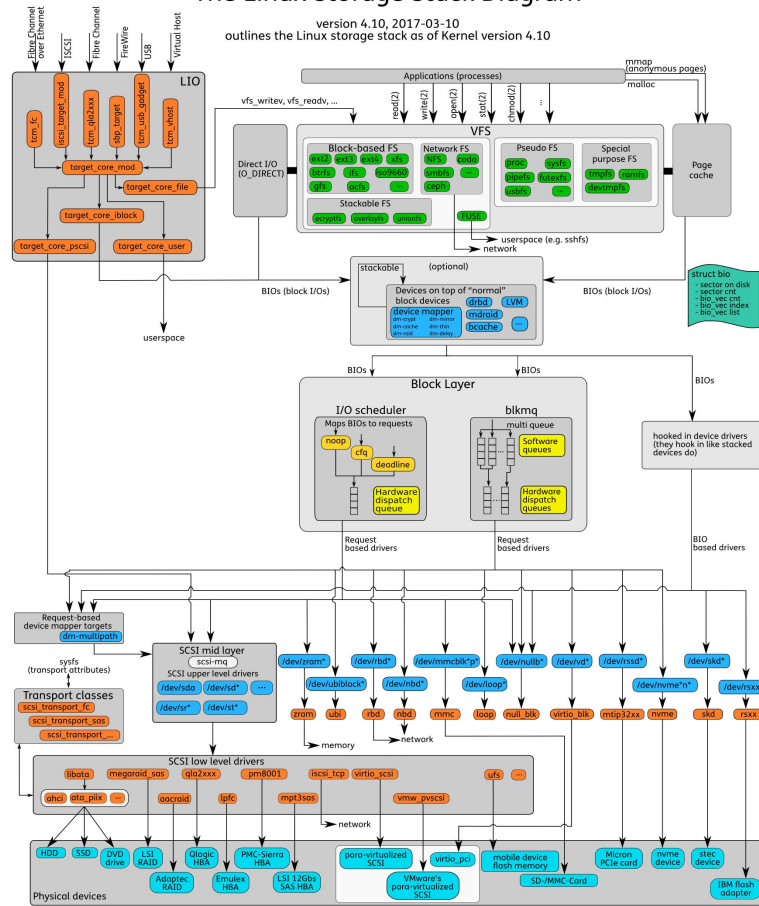(debug) Start Kernel in QEMU  Task
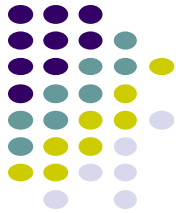
# Linux Storage Stack

- Exhaustive and Modular

# The Linux Storage Stack Diagram

version 4.10, 2017-03-10
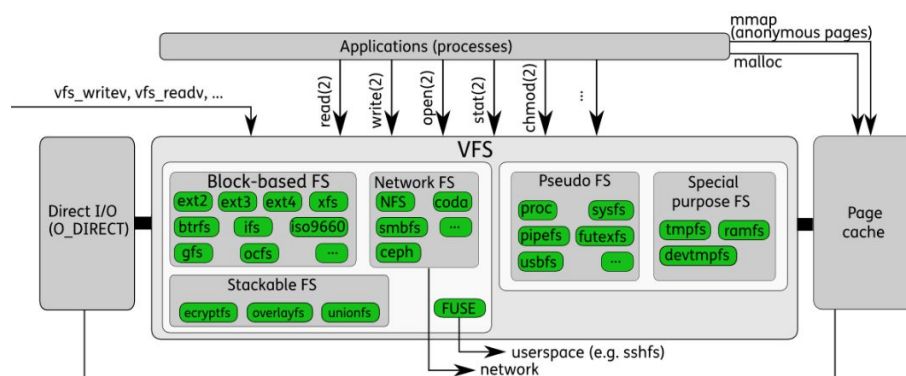outlines the Linux storage stack as of Kernel version 4.10

The Linux Storage Stack Diagram

version 4.10, 2017-03-10
outlines the Linux storage stack as of Kernel version 4.10

Syscall

VFS

FileSystems

Page Cache

Block Scheduler

Device Drivers

Hardware devices

THOMAS KRENN®

12

# VFS



- Virtual File System (~22K SLOC).


- Everything is a File!!
  - E.g., Network file system! sshfs!?


- ~42 File Systems supported in Linux!!

# **VFS to Applications**



- Common interface for accessing files irrespective of file systems.

- File systems no need to worry about interface to user.

# VFS to File System Implementers

- Exposes common optimization logic. E.g., Page cache, Path lookup.

- Define functions to be implemented by the filesystems.

# What does File System Implementers do?

# Page Cache



- Reduce Disk IO

- Memory pages maintained by the kernel for storing contents to/from disks.

- Disk block  <-> Page

# File IO with Page Cache

- *read()*: Serviced by Page Cache!
  - Optimization: Read ahead!

- *write()*: Dirty pages; will be written to disk later!
  - Can loose data!?

- *sync()*: Flush all writes to files.
  - Synchronous

# File IO with Page Cache

char buf[n]

**USER**

read()

**KERNEL**

page cache

disk blocks

# File IO with Page Cache

USER

char buf[n]

read()

KERNEL

page cache

disk blocks

# Page Cache Implementation

- For each file (inode):

  - Has addr space.

  - File offset -> Page cache.

- For each page:
  - A reference to the file/process.

  - The offset with in the file.

# The mmap system call

- Bind virtual memory to file blocks.

```
fd = open("hello.txt", O_RDWR);

// map 4k from offset 0 into virtual address space of the
process.
char *data = mmap(..,fd, 0);

// read 7th character from file.
char c = data[6];

// write 101th character into file.
data[100] = 'a'
```

# Flushing mmap region to file



```
MSYNC(2)

NAME
        msync - synchronize a file with a memory map

SYNOPSIS
        #include <sys/mman.h>

        int msync(void *addr, size_t length, int flags);

DESCRIPTION
        msync() flushes changes made to the in-core copy of a file that was mapped i
        part of the file that corresponds to the memory area starting at addr and having
```

# Memory RW with Page Cache

**USER**

**KERNEL**
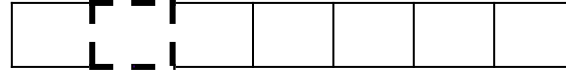
mmap

page cache

disk blocks

# Memory RW with Page Cache

USER

mmap

KERNEL

page cache
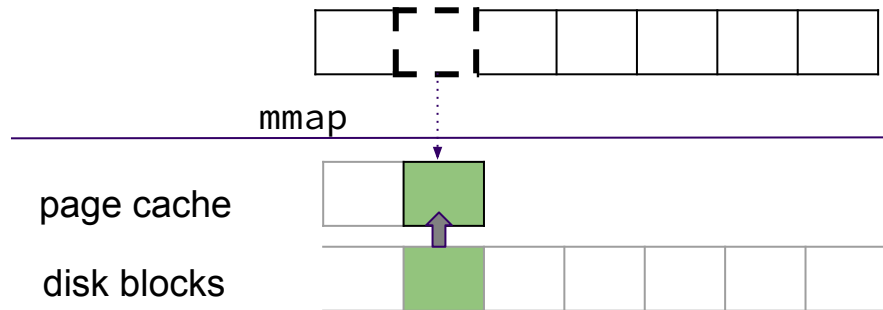
disk blocks

# Mmap v/s Explicit IO

- Mmap:
  - **No syscalls on each access.**
  - **Page cache <-> Disk.**
  - **Dynamic paging.**
  - Extra PTEs.
  - Mapping large files? IO Errors?

mmap

page cache

disk blocks

char buf[n]

read()

- File IO
  - **Universal.**
  - app buffer <-> page cache <-> Disk.

page cache

disk blocks

26