

# Section 11: File Systems, Journaling

CS 162

November 13, 2020

## Contents

<b>1</b>	<b>Vocabulary</b>	<b>2</b>
<b>2</b>	<b>Extending an inode (from last week)</b>	<b>4</b>
<b>3</b>	<b>Comparison of File Allocation Strategies</b>	<b>6</b>
<b>4</b>	<b>Logs and Journaling</b>	<b>7</b>

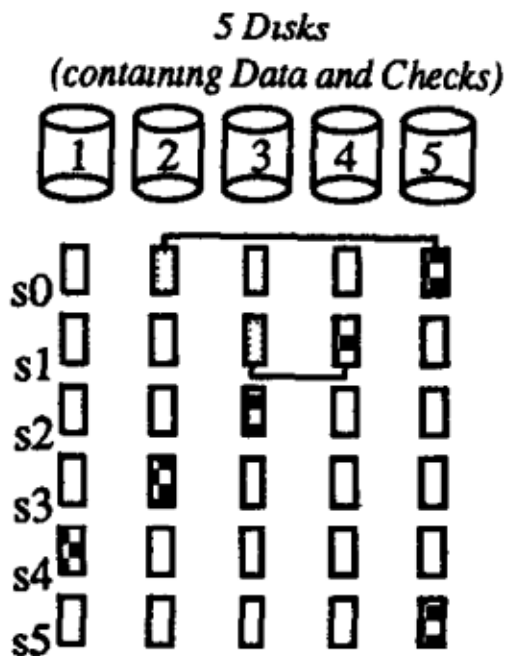
# 1 Vocabulary

- **Fault Tolerance** The ability to preserve certain properties of a system in the face of failure of a component, machine, or data center. Typical properties include consistencies, availability, and persistence.
- **Transaction** - A transaction is a unit of work within a database management system. Each transaction is treated as an indivisible unit which executes independently from other transactions. The ACID properties are usually used to describe reliable transactions.
- **ACID** - An acronym standing for the four key properties of a reliable transaction.
  - *Atomicity* - The transaction must either occur in its entirety, or not at all.
  - *Consistency* - Transactions must take data from one consistent state to another, and cannot compromise data integrity or leave data in an intermediate state.
  - *Isolation* - Concurrent transactions should not interfere with each other; it should appear as if all transactions are serialized.
  - *Durability* - The effect of a committed transaction should persist despite crashes.
- **Idempotent** - An idempotent operation can be repeated without an effect after the first iteration.
- **Log** - An append only, sequential data structure.
- **Checkpoint** - Aka a snapshot. An operation which involves marshaling the system's state. A checkpoint should encapsulate all information about the state of the system without looking at previous updates.
- **Write Ahead Logging (WAL)** - A common design pattern for fault tolerance involves writing updates to a system's state to a log, followed by a commit message. When the system is started it loads an initial state (or snapshot), then applies the updates in the log which are followed by a commit message.
- **Serializable** - A property of transactions which requires that there exists an order in which multiple transactions can be run sequentially to produce the same result. Serializability implies isolation.
- **ARIES** - A logging/recovery algorithm which stands for: Algorithms for Recovery and Isolation Exploiting Semantics. ARIES is characterized by a 3 step algorithm: Analysis, Redo, then Undo. Upon recovery from failure, ARIES guarantees a system will remain in a consistent state.
- **Logging File System** - A logging file system (or journaling file system) is a file system in which all updates are performed via a transaction log ("journal") to ensure consistency in case the system crashes or loses power. Each file system transaction is written to an append-only redo log. Then, the transaction can be committed to disk. In the event of a crash, a file system recovery program can scan the journal and re-apply any transactions that may not have completed successfully. Each transaction must be idempotent, so the recovery program can safely re-apply them.
- **Metadata Logging** - A technique in which only metadata is written to the log rather than writing the entire update to the log. Modern file systems use this technique to avoid duplicating all file system updates.
- **EXT4** - A modern file system primarily used with Linux. It features an FFS style inode structure and metadata journaling.
- **Log Structured File System** - A file system backed entirely by a log.

- **RAID** - A system consisting of a Redundant Array of Inexpensive Disks invented by Patterson, Gibson, and Katz.

The fundamental thesis of RAID is that in most common use cases, it is cheaper and more effective to redundantly store data on cheap disks, than to use/engineer high performance/durable disks.

- **RAID I** - Full disk replication. With RAID I two identical copies of all data is stored. If disk heads are not fully synchronized, this can decrease write performance, but increase read performance.
- **RAID V+** - Striping with error correction. In RAID V, 4 sequential block writes are placed on separate disks, then a 5th parity block is written by XORing the data blocks on the same stripe. RAID VI uses the EVENODD scheme to encode error correction. In general, Reed Solomon coding can be used for an arbitrary number of error correcting disks.



Note: Due to the large size of disks in practice, RAID V is no longer used in practice, because it is too likely that a second disk will fail while the first is recovering. RAID VI is usually combined with other error recovery techniques in practice.

- **Eventual Consistency** - A weaker form of a consistency guarantee. If a system is eventually consistent, it will converge to a consistent state over time.

## 2 Extending an inode (from last week)

Consider the following `inode_disk` struct, which is used on a disk with a 512 byte block size.

```
/* Definition of block_sector_t */
typedef uint32_t block_sector_t;
/* Contents of on-disk inode. Must be exactly 512 bytes long. */
struct inode_disk
{
    off_t length;                /* File size in bytes. */
    block_sector_t direct[12];   /* 12 direct pointers */
    block_sector_t indirect;     /* a singly indirect pointer */
    uint32_t unused[114];        /* Not used. */
};
```

Why isn't the file name stored inside the `inode_disk` struct?

What is the maximum file size supported by this inode design?

How would you design the in-memory representation of the indirect block? (e.g. the disk sector that corresponds to an inode's `indirect` member)

Implement the following function, which changes the size of an inode. If the resize operation fails, the inode should be unchanged and the function should return `false`. Use the value 0 for unallocated block pointers. You do not need to write the inode itself back to disk. You can use these functions:

- “`block_sector_t block_allocate()`” – Allocates a disk block and returns the sector number. If the disk is full, then returns 0. For this question, assume unlimited disk space.
- “`void block_free(block_sector_t n)`” – Free a disk block.
- “`void block_read(block_sector_t n, uint8_t buffer[512])`” – Reads the contents of a disk sector into a buffer.
- “`void block_write(block_sector_t n, uint8_t buffer[512])`” – Writes the contents of a buffer into a disk sector.

1. Handle resizing the inode for direct pointers first

```
bool inode_resize(struct inode_disk *id, off_t size) {
    block_sector_t sector; // A variable that may be useful.
    for (int i = 0; i < 12; i++) { // Handle direct pointers
        if (size <= 512 * i && id->direct[i] != 0) { // Shrink inode
            -----
            -----
        }
        if (size > 512 * i && id->direct[i] == 0) { // Grow inode
            -----
        }
    }
}
```

2. Next, try to extend your implementation to the indirect pointers. Note we have to check if the indirect pointer page is even allocated yet, and do so if necessary. (we are still in `inode_resize`)

```
if (id->indirect == 0 && size <= 12 * 512) { // Check to handle indirects
    id->length = size;
    return true;
}
block_sector_t buffer[128];
memset(buffer, 0, 512);
if (id->indirect == 0) { // Allocate indirect page if does not exist
    -----
} else {
    -----
}
for (int i = 0; i < 128; i++) { // Handle indirect pointers
    if (size <= (12 + i) * 512 && buffer[i] != 0) { // Shrink inode
        -----
        -----
    }
    if (size > (12 + i) * 512 && buffer[i] == 0) { // Grow inode
        -----
    }
}
if (id->indirect != 0 && size <= 12 * 512) {
    -----
    -----
} else {
    -----
}
id->length = size;
return true;
}
```

### 3 Comparison of File Allocation Strategies

In lecture three file allocation strategies were discussed: (a) Indexed files. (b) Linked files. (c) Contiguous (extent-based) allocation.

Each of these strategies has advantages and disadvantages, which depend on the goals of the file system and the expected file access patterns. For each of the following situations, rank the three designs in order of best to worst. Give a reason for your ranking.

1. You have a file system where the most important criteria is the performance of sequential access to very large files.

2. You have a file system where the most important criteria is the performance of random access to very large files.

3. You have a file system where the most important criteria is the utilization of the disk capacity (i.e. getting the most file bytes on the disk).

## 4 Logs and Journaling

You create two new files,  $F_1$  and  $F_2$ , right before your laptop's battery dies. You plug in and reboot your computer, and the operating system finds the following sequence of log entries in the file system's journal.

1. Find free blocks  $x_1, x_2, \dots, x_n$  to store the contents of  $F_1$ , and update the free map to mark these blocks as used.
2. Allocate a new inode for the file  $F_1$ , pointing to its data blocks.
3. Add a directory entry to  $F_1$ 's parent directory referring to this inode.
4. *Commit*
5. Find free blocks  $y_1, y_2, \dots, y_n$  to store the contents of  $F_2$ , and update the free map to mark these blocks as used.
6. Allocate a new inode for the file  $F_2$ , pointing to its data blocks.

What are the possible states of files  $F_1$  and  $F_2$  *on disk* at boot time?

Say the following entries are also found at the end of the log:

7. Add a directory entry to  $F_2$ 's parent directory referring to  $F_2$ 's inode.
8. *Commit*

How does this change the possible states of file  $F_2$  on disk at boot time?

Say the log contained only entries (5) through (8) shown above. What are the possible states of file  $F_1$  on disk at the time of the reboot?

What is the purpose of the *Commit* entries in the log?

When recovering from a system crash and applying the updates recorded in the journal, does the OS need to check if these updates were partially applied before the failure?