

Section 2: Processes

Process: Program instance, contains code & activity.
Return: Return value from child to parent process.
Address Space: Private set of mem. addresses per process.
Stack: (IFO memory) for local thread variables.
Heap: Large memory area to allocate/deallocate vars.
fork: Fork (duplicates) process. On success, child PID returned in parent, 0 in child. Failure, -1 in parent.
wait: Wait() waits until child terminates or stops.
exec: exec() executes program in process.

Section 3: Threads

Thread: Smallest unit of seq. instructions. Mult. threads in process can share same addr. space, each has own PC.
Thread join: Waits for specific thread to terminate.
Thread create: Starts child thread in same addr. space.
Atomic: Operation atomic, i.e. executed w/o interruption.
Critical Section: Code accessing shared resource, cannot be concurrently accessed by more than one thread.
Semaphores: Semaphore increments, semaphore waits to become positive, then decrements.

Section 4: Scheduling & Synchronization

Lock: Synchronization vech. provides mutual exclusion. Only one thread holds lock. Acquiring held lock will block.
Condition Variable: Provide serialization, CV defined by associated lock, condition, wait queue. Must hold lock to access CV functions; if condition false, thread blocks.
Wait queue: Thread, signals to one/all waiting when cond. true.
Thread yield: Thread put back in ready queue (might be first).
Race Condition: State of execution causing multiple threads to access shared memory w/o mut. exc. Undefined.
Scheduler: Picks which thread to run next.
Priority Inversion: Low pri thread holds lock that high pri need.
Priority Donator: Thread waits on lock donates its priority to other thread recursively until thread has all locks and can run to completion.
Hoare/Mesa: Hoare = imm. transfers control. Mesa = while

Section 5: CV & Spin Locks

Condition Variable: Sync variable w/ lock (CV + lock = monitor).
boolean condition, wait queue. **CV wait** on lock held lock.
puts thread to sleep on wait queue, **CV signal** removes one thread from wait queue into ready state **broadcast** removes all threads from queue, put in to ready state.
When waiting thread put back into ready state, a **lock** is available.
True/False: must acquire lock, modify condition, then call signal.
Have semantics auto-wake thread when cond. true. Mesa simply puts on ready list, use while in case cond. changes.
Spin Locks: Busy wait on lock acquire until lock becomes free.

Section 6: Scheduling & Fairness

Scheduler: Decides which processes run when.

Calling Conventions: Interface of called code. Define order/atomic params allocated/passed, registers, stack. %esp, %ebp: %esp holds last thing on stack, gets moved as things get added/removed. %ebp holds begin of cur. stack.
Virtual Memory: Process believes in its own address space, each process has virtual addr. maps to physical.
Physical Memory: Actual memory addresses in hardware.
Page: Contiguous fixed size chunk allocated to process.
Page Table, Memory Management Unit (MMU): hardware that translates virtual addresses to physical. Page table stores mapping b/w virtual/physical, resides in MMU.
Stack: low-level (below malloc) & function which can increase/decrease memory allocated to program.

Section 8: Wait/Exit, Address Translation

Exit: Terminates process. OS reclaims resources (memory, files). Process becomes dead process.
wait: Process waits on another (child) to finish exec.
Zombie Process: Child process becomes zombie after termination or it still lives in system process table. If parent waits on child, will be reclaimed once actual value returns to parent. If not, reaper OS process locates zombies, retrieves exit, deallocates process.
Address Translation Structures: **Segment:** one or many addressed chunks of memory containing logically related info. i.e. program code, process heap/stack, etc. Of the form (S, i), where must be within i of base segment's. **Page Table:** stores mapping between virtual \rightarrow physical. VAD used by accessing process, PA used by the RAM.
Inverted Page Table: Uses PT that contain entry per physical frame, not logical page. Erlson IPT occupies fixed fraction of memory, proportional to physical mem. not virtual addr. space, only one per system. Can be like hash map: stores key & entry.
Translation Lookaside Buffer: TLB is cache that mem. management hardware uses to speed translation. Stores VA \rightarrow physical mappings so MMU can store recently used mappings instead of multiple lookups.

Section 9: IO & File Systems

Second Chance Algorithm: Modified FIFO used to approximate LRU. Each page has use bit; when evicting, if ref. set, clean & send back of list, else evict.
Clock Algorithm: More efficient second chance b/c pages don't need to go to back. Keeps head pointer @ page in circular order; hand moves round, cleans ref bit, evicts.
IO: Input/output is processes by which OS gets/puts data to/from controller which contains mem./registers for communication w/ CPU & interface for communication of hardware. Communication done via programmed IO, transferring data through registers, or Direct Memory Access, which allows controller to write directly to memory.
Interrupt IO: interrupts OS, interrupt hardware, useful for infrequent, sporadic events.
Poll: OS checks regularly for IO. Less overhead, better for regular events such as mouse input.
Response Time: Time b/w IO request & completion.
throughput: rate operations performed over time.
Asynch IO: IO returns immediately, not finished when done.
indep. File system object, communication, pointers, etc.

Section 10: File Systems & Query Theory

Simple File System: Disk treated as big array; Table of Contents (TOC) beginning followed by data. Files start contiguously, but can be unused space between. In TOC, file descriptions entries w/ name, start location, and size.

Pros: Simplicity, quick IO. Cons: External fragmentation, cannot easily get extended, no directory hierarchy, file type external fragmentation: free space (w/ allocation) ceases, still so nice, not big enough for new file.

Internal fragmentation: leftover space (end of block). File Allocation Table: FAT views disk as array. First block is boot sector, which contains bootup info. Super block w/ file system metadata. FAT after. Followed by File Allocation Table, all other data in 4kib blocks. File viewed as linked list of data blocks, pointers stored in FAT so blocks are 100% data. 1:1 correspondence b/w FAT entries & data blocks, each FAT entry stores block index. Entries: $N > 0$: No next block index. $N = 0$, end of file. $N = 1$: free block. Files can be stored non-contiguously. Maximum interval $n = 4kibytes - 1 = 4095bytes$. Directory one file store entries. Pros: No external frag, can gran file size, has directory hierarchy. Cons: No pre-allocation, not contiguous (slow IO), assume FAT fits in RAM.

UNIX/Fast File System: Disk space divided into: Boot sector, Super block, Free Block bit map (one bit per block, more efficient than table), Inode table, Data blocks. Inodes: Data structure w/ file metadata. Contains: ownership, size, mod time, permissions, ref count, data block pointers, NOT name. 12 direct, 1 indirect, 1 Doubly, 1 Triple. New Technology File System: NFS represents file as record in Master File Table. First record describes MFT, itself, followed by MFT mirror.

Log-structured FS: Data written to circular log buffer. Throughput improved bc avoids seeks, can batch writes. Writes create multiple chronologically advancing series of file data. Super-crash occurs from consistent state. Unlike Sxall to delete file. Decrements ref count, renames entry from directory. Garbage collected @ 0 ref count.

Query Theory: λ (average service rate). $T_{ser} = \frac{1}{\lambda}$, $\lambda = \text{average arrival rate (jobs/sec)}$. $U = \rho = \frac{\lambda}{\mu} = \lambda S = \text{utilization}$. $W = \frac{1}{\mu - \lambda}$, $T_0 = W = \text{average query time}$. $T_{sys} = R = T + T_{ser} = W + S = \text{response time}$. $L_q = Q = \text{average length of queue}$. Little's law: $L_q = Q = \lambda T_q$.

Section 11: FS Details & Reliability

Transaction: indivisible set of ops that fail or succeed.

ACID: Atomicity: Transaction must occur in entirety or not at all. Consistency: Data must be in consistent state.

Isolation: Transactions should not interfere.

Durability: Should persist despite crashes.

Idempotent: Operation can be repeated w/o change.

Logging FS: FS where modifications done via transactions, metadata append only log, not disk. Once committed, can be written to disk; on crash, can be recovered. 2PC: Two phase lock ensures atomic, isolated commit of concurrent transactions. Each transaction must acquire shared/exclusion lock for read/write. Grow phase: at only active, followed by Shrink phase (can only decrease). Final 2PC has all released at once. 2PL ensures acyclic dependency graph, serializable transactions.

Section 12: Two Phase Commit (Byzantine: need majority, not all)

2PC: Two phase commit coordinates transactions between one master and slaves. State-changing transactions must be logged and checked. TPC: master write happens across all replicas on node. Master sends message to slaves, which return COMMIT or ABORT. If all commit, master sends GLOBAL COMMIT and Slaves Ack; else, Master sends GLOBAL ABORT.

Security: Computing w/ adversers. Need reliability, robustness.

Fault tolerance: Protection Mechanism for controlling access of programs, processes, resources. Page Table Mechanism, Round Robin scheduling, data encryption. Security Protection Mech. to prevent use of resources. Requirements:

- Authentication: User is who claim to be.
- Data integrity: Data not changed for, pre-ident.
- Confidentiality: Data only read by auth. user.
- Non-repudiation: Sender/receiver can't deny did send/receive data.

Passwords: auth measure. Use cryptography. Symmetric keys: Encrypt/decrypt w/ same key. Vulnerable. XOR, use block cipher to encrypt, decrypt.

Data Encryption Standard: 56 bit. Addressed Asymmetric. Standard: 128/192/256 bits. No. Proof of strength.

Secret key: $\text{Enc}(\text{data}, K) = \text{Dec}(\text{data}, K)$. Vulnerable to man in the middle. Crypto Hash W: $H(x)$ public hash. $\text{Digest} = \text{HMAC}(K, m) = \text{H}(K \| H(m))$.

Send data to receiver. Receiver uses K to recompute $\text{HMAC} = d$. MD-5, SHA-1 broken, SHA-2 in use. Asymmetric crypto: Use e to encrypt, d to decrypt, knowing one still keeps other safe.

RSA: Use receiver's public key to encrypt, decrypt w/ private key. Digital Certificate: Binds entity w/ public/private key.

Trusted authority distributes, anyone can extract public key. HTTPS: Uses SSL (Secure socket layer) / TLS (Transport Layer Security) to distribute digital certificates. Validates, RSA.

Networking: Broadcast Medium: Shared comm. medium. Arbitration: How to negotiate shared medium. Aloha protocol. CSMA: Don't send unless idle. Collision Detect: Sender checks if packet trampled; if so, abort wait, retransmit.

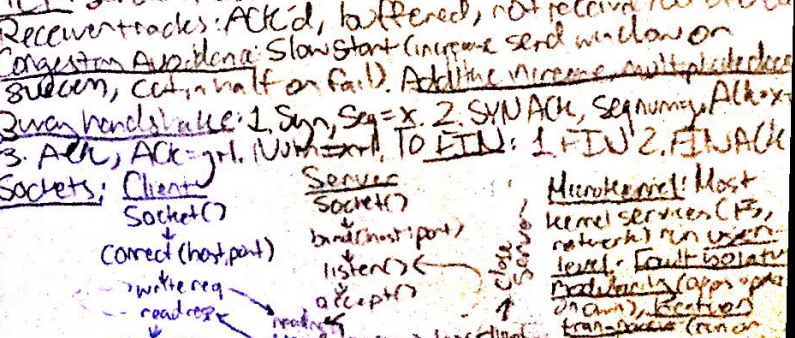
Backoff scheme: Choose wait time before trying again. Adaptive and random increases wait time after each failure. Internet Protocol (IP): 32 bit/host. Share through NAT.

Subnets: Network of hosts sharing prefix (number of bits). Routing Tables: (Addr range -> next hop). Set up dynamically.

IP sends machine -> machine, a dd ports. pre-emptive DDP. For mult pkts: window based acknowledgment. Send up to N packets w/ sequence number, receive ACKs.

TCP: Sender tracks: sent & ACK'd, sent/not ACK'd, to be sent. Receiver tracks: ACK'd, buffered, not received (w/ out of order). Congestion Avoidance: Slow start (increase send window on success, cut in half on fail). Additive increase, multiplicative decrease.

Three handshake: 1. SYN, Seq=x. 2. SYNACK, Seq=y, Ack=x. 3. ACK, Ack=y. Non-conn. To FIN: 1. FIN, 2. FINACK.



Distributed File System: Transparent access to files on remote
Name: Hostname/Vocabname. Mounting: Can mount
remote by giving local name e.g. /usr/sue/100, /usr/100
on server. (File has unique global name.)

API Theorem: Cannot have consistency, availability (can get result
anytime), Partition-Tolerance (system continues to work
when network part down) & save time & only two.

Remote Procedure Call (RPC): Abstracts for executing code
code on remote machines using stubgen, & run-time lib.

Stub Generator: Server sends out interface of
function calls; Client programs have client stubs w/
all more functions, client simply calls w/ args.

Internally, client stub code: - Generates Message
buffer (array of bytes). - Packs into nto buffer (includes
function identifier, args) Marshalling: packing. - Send to
RPC server (handled by run-time lib) - Wait for
reply, synchronous. - Unpack return value, unmarshalling.
- Return from caller, client stubs back to client.

Internally, server code: Unpack buffer, unmarshalling.
- Call into function. - Package (marshal) results. - Send reply
Server handles core users by keeping thread pool.

Run-Time Library: Handles performance & reliability issues.
Handles naming by building on existing protocols i.e. uses
IP address, no time port, RTL then directs. Doesn't send
over TCP b/c takes down user UDP w/ retries on failure.

Users binding to bind network name to readable name for client.
Issues include non-atomic failure (machines do different things
on failure) and performance (time + caching, memory).

Simple DFS: reads/writes forwarded to server using RPC.
Not cache coherent, only server. Server inconsistent view of
FS to clients, but slow performance, high traffic.

Use caching to improve performance, but must maintain
cache consistency (clients may have data not @ server).
If server fails, clients can lose data / be inconsistent.

Can flow state in protocol: server has no state, message
contains all info. Might lose some data on client crash.

Network File System: 3 layers: 1. User FS interface
open, read, write, close + file descriptors 2. VFS layer:
(distinguishes local files from remote) 3. NFS service layer:
bottom layer implements NFS protocol; RPC for server calls.

Write through caching. Data committed to disk as
soon as modified (slow writes). NFS server stateless.

Idempotent: Some request mult. times has same effect.
Failure model: transparent to client, can hang or error out.

NFS enforces weak cache consistency: polls for changes (30s),
Means values only updated after poll -> read new 30s after write

Pros of NFS: Portable, simple. Cons: Inconsistent, doesn't scale.

Andrew File System Callbacks: Server records who has
copy of file. On change, server alerts them w/ old copy,
write through only for client, updates to client
only when file closed. Data cached on local disk
or client + memory. Set callback on mount/write.

Less server load from NFS, but still server bottleneck.

Virtual File System Switch: Virtual abstraction similar to
local file system; virtual superblock, inodes, file etc.

Allows give API for diff file systems. FOUR princ.
object types: superblocks (mounted FS), inodes (specific file),
directories (dirs), file (open file) Microkernel (see other page)

Key-Value Store: Distributed hash table, can put, get
del key/values. Partition set of key-values across
multiple machines. Challenges: fault tolerance (handles
machine failures w/o data/performance loss),
consistency (in face of failures, message loss),
scalability, heterogeneity (1ms to 1000s latency,
32 kbps-100mbps bandwidth). Directory-based
structure: Have a node w/ dir mapping blobs to keys
and machines that store them. Recursive query:
has server w/ mapping, contact node, get value, return
iterative query has server return node to client,
client contact directly. Recursive in faster (server
closer to nodes), and easier to serialize, but has
server bottleneck. Iterative is more scalable,
but slower, hard to enforce consistency. To improve
fault tolerance, replicate values across multiple nodes.

To improve Scalability: increase storage (more nodes),
improve concurrency (replication -> multiple GETs), implement
Load Balancing (balance usage amount across nodes),
To improve consistency: make sure nodes replicate correctly.

Consistency models: atomic (RW to replican appear as if
there's a single replica; transactions), eventual (given
enough time, all updates will prop through system).

Quorum Consensus: Improve put/get performance & Define
replica size N, put waits for A client w/
replicas, get waits for B client R; W + R > N.

To scale up, use consistent hashing: give node ID
#s, hash (key), store key pair in node w/ id right above

MapReduce: Map: (K_i, V_i) -> list (K_{inter}, V_{inter}). Reduce:
(K_{inter}, list (V_{inter})) -> list (K_{out}, V_{out}). Restricted KV model
deterministic, idempotent, some fine-grained op on input

Pros: Distributed in transparent, automate fault-tolerant
(from failed tasks), scaling, load-balancing.

Cons: Restricted programming model, high latency, read-
only requirements for Deadlock: (1) Mutual exclusion
(anyone thread @ a time can use resource) (2) Hold
and Wait (thread w/ resource waits for rest) (3) No
preemption (thread only voluntarily releases resource
after finished) (4) Circular Wait (wait on each other)

Banker's Algorithm: Keep system in SAFE state,
only grant resources if some thread has enough to
proceed after allocation

Address Translation: Base and Bound. Thread gives
VA blocks can only access if above base, below bound

Cons: Fragmentation, bad fit - sparse space can't show
Segmentation: VA made of segments - offset maps
to table w/ base and bound, PA = base + offset + offset

Page Table: VA made of VPN + offset, VPN maps to
PPN in PT, PA = PPN + offset, PT has address per
Multi-Level: VA = VPN1 + VPN2 + offset, VPN1 maps
to second level PT, DN = DPN from PT2 using VPN2's
Invented, VPN hashes to PPN - offset

Calculations for single level: offset bits: log₂(pages)
VPN = 32 - offset - bits U = RPN, PTE = PPN + mem
PI = (PTE bits) * (num pages), num pages = 2^{VPN}

Caching: Cache miss (complicated first access) Capacity
conflict (mapped to same), coherence (invalidated on
Direct map: Addr = Tag + Index or Byte, index - row,
store tag, select w/ byte. Set associated. Index
maps to rows, store tag in any of sectors

Fully associative Tag + byte, store tag anywhere
Cache coherence translation for PT's, over the wire

Demand Paging uses memory on cache & disk
maps map. disk to memory, improves IO

MTJ: disabling interrupts or CPU does nothing, or on the
Simultaneous multithreading = hyperthreading = splits
threads every cycle. For HTTP sockets, child & parent
must both close. Kernel mode happens @ well-
defined entry points. Thread blocked only on
one CV @ once. Disks don't use float
point b/c it doesn't save in registers.

CVs use queues of semaphores so need to
implement scheduling for CV, not semaphore.
Printf buffer, doesn't need kernel.

Disabling interrupts to create crit. section: (1)
user level code cannot (2) lock out HW, important
events can be missed (3) Slow performance.

Only work when: interrupts would cause incem. below

Hoare hands control to waiting threads immediately

Mesa puts on queue, no guarantee on waking
Context switch: Threads: save registers, program
counter, condition registers, thread execution state

Processes: additional save PT pointer, seg. registers

User \rightarrow kernel: Syscall, exception, interrupt

fork(): Create child w/ dup. addr. space.

exec(): throw away addr. space, runs executable

Overuse of threads bad: (1) waste cycles on syscall

(2) waste memory for stacks & TCB

Spinlocks only efficient when expected wait time
is less than time to switch threads

MTJ + mmap maps VA to pages in buffer cache
which hold contents of file. FFS does not

include name. MMIO maps device memory
to physical addr. space, threads can access.

That means not all physical addresses in DRAM.

(ZPL guarantees consistent serializable, bc
conflicts turn into deadlocks which are aborted.

Combo of policy & interrupts to handle events.

Network Addr. Translation allows more than 2^{32} comp.

IP router needs preexist + next hop to them.

Cachon effect when workset fits in them.

Process exception is when everything before exception
guaranteed to have happened; helps restart.

RAIDS: XOR blocks to recover failed one

Copy on write: copy addr. space read only, make copy
of page or page write change in PT.

syscall # & syscall args passed from user to kernel
by pushing on stack in reverse.

Direct Memory Access: DMA controller handles
data transfer to memory