

OS

Resources: OS has to make efficient use of resources and share them to multiple users. **Processors:** CPU cores, multi-socket CPU, SMT/hyper-threading. **Memory:** caches, RAM, I/O devices: GPUs, NICs, printers, **Internal Devices:** clocks, timers, interrupt controllers, **Persistent Storage:** disks, SSDs, DVDs.

Computer Architecture Overview:

I/O controller has intelligence

Disk

CD-Rom

Printer

Processor

Disk controller

CD-Rom controller

Printer controller

System bus

Sharing Resources: Fair, efficient, simultaneous, protective

Memory

System interconnection

Process: Running program, user mode (kernel through interface), some foreground (user interaction), others b-grnd (mail/print server, called *daemons*). Contain state: *addr space* (own mem, heap/static data (e.g. instrs)), *registers*, *OS resources* (e.g. open files). **A** = isolation (only kernel requests affect each other), safety, simplicity (prog-er not worried abt others), (kernel has control, concurrency. Processes created at system initialisation, user request, or a syscall by running process.

space shared. Processes have >=1 thd. Processes are slow + expensive to create/destroy, and comms are complex (pipes/sig handlers). **A** = shared data (all thds in process share addr space, e.g. R/W to same mem to comm), cheap switching (only change per-thd state - registers, SP), cheap management (basic, cheap). **D** = concurrency bugs, forks (holding locks in children never released), blocking (if 1 thd blocked, all in process blocked). **Threads** (POSIX ts): **POSIX interface for managing thds**. *int pthread_create(pthread_t *thread, const pthread_attr_t *attr, void *start_routine(void *), void *arg)* - creates new thd. Params: thread - pointer to store thd id, attr = attributes, start_routine = func to start new thd on arg, arg - arg to pass start_routine. Returns: 0 -> success, EAGAIN -> lack of resources, EINVAL -> invalid attr, EPERM -> caller lacks permissions. *void pthread_exit(void *value_ptr)* - exit thd, provide value_ptr to joins onto thd. Catch implicitly at end of start_routine. Process terminate -> terminate thds, main thd exit -> waits for children/exit call. Remember main has implicit call to exit added by compiler: *int pthread_exit(void)* - yield CPU to another thd. Returns: 0/-1. Always succeeds on Linux. *int pthread_join(pthread_t thread, void **value_ptr)* - block curr thd until other thd terminates. Params: thread - thd to wait on, value_ptr - location to put thd termination val. Returns: 0 -> success, EDEADLOCK -> deadlock (thd already joined), EINVAL -> thd param not valid, ESRCH -> couldn't find thd.

Scheduler Goals: Fairness (comparable processes get comparable CPU time), all run (avoid indefinitely postponing processes), maximise utilisation (maximise resource allocation), minimise overhead (context switches, running scheduler), correctness (scheduler correctly impl). **Fair-Share Scheduling:** each user has a ready queue used by scheduler, which RRs between them to fairly distribute CPU time. We study div time fairly for processes. **Batch Systems:** Optimise for throughput (jobs/unit time) or turnaround time (time to completion). **Interactive Systems:** Response time critical. **Real-Time Systems:** Run jobs to meet deadlines (soft deadlines e.g. slow video recoverable, hard e.g. factory arms colliding non-recoverable).

Schedulers: FCFS: No pre-emption. Assume all pcs sched eventually terminate/block (to allow others to run). **A** = all run (eventually), correctness (simple+ez impl), minimise overhead (simple queue - low overhead, low context switches as no pre-emption). **D** = turnaround time (TT) (small pcs stick behind large one). **Round Robin (RR):** Periodic timer interrupt, prev pcs at back of ready queue and next run. Larger quanta -> small overhead, higher response time. Smaller quanta -> opposite. RR correctness (ez). **D** = TT (high when runtime similar), context switches (lowers runtime response time, inc overhead). 100ms quanta by default. (*sched_get_interval*). **Shortest Job First (SJF):** Schedule ready thd with shortest time remaining, assumes we know time for each pcs, non pre-emptive. **A** = TT (optimally low, if shorter added to queue current finishes first still), correctness (simple if all short enough), context switches (very few). **D** = response time (can be high as non-pre), indefinite postponement (if many short added, long never scheduled). **Shortest Remaining Time (SRT):** Pre-emptive ver of SJF. If shorter job (than curr remaining time) pre-empt curr. Need to track time taken for jobs. **A** = TT (low as SRT always run). **D** = indef post (many short, long never sched). **Throughput** = **no. jobs / sum of time for each, Turnaround = time taken for jobs (and all before it) to finish, RR overhead = context switch / quantum + context switch.**

Synchronisation: For concurrent apps, cannot make assumptions about relative speed of exec for threads (context switch at any time), must consider op interleaving, real (true) concurrency consideration in multiprocessor systems. **A** = parallelism (independent ops in parallel reduce time taken), blocking (to avoid blocking ops e.g. I/O preventing progress, can separate them to diff threads/processes).

Race Condition: bug when result depends on non-deterministic order of ops. **Thread B**

A

Read

Write

Read

Write

Write

Fine

Race Condition!

Race Condition!

Synchronising Access: Critical Section: code accessing shared resource/data which requires mutual exclusion. Signalled at start/end of CS. **Mutual Exclusion:** 1 thd in CS at a time, no thd outside CS can prevent other thds from entering. No thd needing access to CS delayed indefinitely. No ass made about speed/scheduling. **Busy Waiting:** program continually checks on wait cond instead of blocking and being awoken (used when wait time small, wastes CPU time).

True atomicity. **Semaphores:** counter (no. threads which can enter, >1 useful for producer-consumer) and list of (of blocked threads waiting to enter). Alt to busy waiting. Can enforce ME (init sema with val 1) without caring about number of thds. *int sem_init(&sem_t *sem, int pshared, unsigned int value)* init sema with val 1. *int sem_getvalue(&sem_t *s, int *val)*, *int sem_destroy(&sem_t *s)*. Enforce ordering (init to 0, thd 0 start with down(), thd 1 start with up()) - 1 goes first as can't be downed from 0). (aka up() forces ME. 1 thd hold L at given time. If attempt to get L but already acq, block requesting thread. Only holding thd can release L. Methods: init, lock, down, Re-entrant Ls: allow L to be acquired many times by same thd. If holder attempts to lock again, will deadlock. **Spin Locks:** if expected wait low, busy waiting may be better than blocking threads. Locks busy wait (in while loop, while (TSQ==1-locked) = 0), no kernel involvement (no block/unblock syscalls), ensure check on acquiring is atomic, can run into priority inversion. **Priority Inversion:** when using priority-based scheduling algo, low-p threads can block high-p ones. LP holds resource which HP wants e.g. lock. HP tries to acq but blocked and busy waits, LP not scheduled as P is low so cannot release resource. Solution = **priority donation** -> HP donates to LP so it can finish and release resource. Aging can also let LP execute if it has been a while. Spin locks mean HP always ready to run as no blocking used, but then LP cannot release lock. **Read/Write Locks:** allow many readers to hold lock, or 1 writer. Reduced lock contention as readers don't block. Some impls allow R -> W upgrade or W -> R downgrade w/o release + re-acquire.

Producer-Consumer: Constraints: **Producer:** can only deposit (produce) when space in buffer, only deposit if ME ensured. **Consumer:** only fetch (consume) if buffer not empty, only fetch if ME ensured. **Buffer:** limited capacity (0-N spaces). **P-C w/ Semaphores:** uses semas to count num items/space (if can't send down then know we have reached lim so block). **P-C w/ Monitors:** CV to model producer/consumer's constraints.

Design Decisions: Lock Overhead: Measure of extra resources required for using locks (memory alloc for lock structure, time to init/destroy, time for acq/release). **Lock Contention:** measure of num thds waiting on a lock (more contention = more thds blocked waiting on lock or waiting w/ spinlocks and less concurrency as blocked rather than running in parallel). **Lock Granularity:** measure of amt of data protected by a lock (coarse = low limit, fine = each list elem). In designing the synchronisation scheme, goals are: correctness, reduce overhead+contention.

Lock Prevention: Ostreich: ignore it. **Detection and Recovery:** Dynamically build RAG, lock for cycles (DFS), remove member in cycle if one found (e.g. revoke resource/end pcs). **Recovery:** pre-emption (take resource from owner and give to another temp), rollback (take periodic snapshot of sys and roll back+restart), killing pcs in cycle (dangerous). **Dynamic Avoidance:** for each resource request, check if it will block, e.g. Banker's Algorithm. **Prevention:** write code so 4 reqs for Lock never met. **Mutual exclusion:** remove ME and share resource. **A** = no overhead (no synch overhead), more concurrency (less lock contention). **D** = race conditions. **Hold and wait:** all pcs specify what res they need before starting. **D** = redundant requests (res which may be used req'd - waste), advance knowledge (don't know what res used in advance). **No pre-emption:** allow resources to be forcibly (sometimes temp) revoked. **D** = correctness (logic issue if revoked during operations). Allow pcs to attempt to get res, req can be revoked, release other res if reqs are revoked. **D** = complexity (makes impl even basic programs difficult. Ensure res are acq and released in same order by pcs/thds. **A** = simple (no mechanisms needing interrupts e.g. killing pcs/revoking res), efficient (less overhead due to simple point). **D** = ordering (need to order large num of res and ensure all pcs acq in order).

Fragmentation: External: Enough mem available, no hole large enough. Fix: *compaction* -> shuffle memory to put all free mem in one block (can result in I/O bottleneck - large amt of copying). **Internal:** Allocated mem larger than req'd (wasted memory).

Swapping: Num pcs limited by available main mem (only running pcs need to be in main mem). Supplement main mem with *swap space* (partition or file) on secondary storage (SSD, HDD) -> swap non-running pcs to disk, bring back when exec again. Transfer time is long.

Virtual Memory with Paging: Abstraction to separate log mem from phys. Not all of exec effs need to be in mem for exec. Log addr space can be >> phys addr space. Addr spaces can be shared between pcs. Pcs create more efficient. Fixed pg size -> ext frag impossible. **Virtual Addr Space:** Impl via paging or segmentation. **0, code, data, heap** -> **<= stack, max. Frames (F):** fixed-size blocks of phys mem, OS tracks free FS. **Page:** block the same size as F in log mem (effectively a log F). **Program Start:** Program frame in memory. 1. find n free frames, load prog into that mem, 2. create page table for pcs mapping log to phys. **Termination:** destroy prog table, freeing all associated FS. **A of VM w/ Paging:** less int (if page unused, not alloc in mem) and ext (when alloc of mem, split into pgs-dist across mem) fragmentation, and pcs alloc pgs when phys mem available (contiguous log addr can be phys non-contiguous). **Page Sizes:** Small = less int frag, potentially less mem to swap, larger pg table, best for efficient mem usage. Large = more int frag, smaller pg table, best for low addr lookup overhead. **TLB:** Translation Lookaside Buffer used to cache mappings of virtual to phys address. When virt addr used + corresponding phys addr found, added to cache to avoid re-computation. **Context Switch:** OS locate pg table for new pcs, set base reg in MMU to new pg table in mem, clear now-invalid cache addr calls from TLB.

Page Table Types: Multi-Level/Hierarchical: break up log addr space into multiple Pts. **PT size = (addr space size / pg size) * entry size**. Large addr space+small pg size = large PT - issue (if PT larger than F, added complexity, PTs uses lots of mem even if pcs doesn't access any/many pgs). Fix = use multi-level PTs (a tree). For N levels, pg num partitioned into n sections. **Pg Accessing:** 1. get outermost PT (typically in single frame), 2. for each partition of pg num: a. index from partition, b. entry at index in curr frame (a PT), c. in entry invalid/empty/dirty, page fault, d. else get frame from entry (set as curr frame), e. get phys addr (in curr frame using pg offset). D = more mem accesses for pcs to access mem (in TLB miss). **P2 ->**

2-level Paging

42 30 12

2nd pte offset outer pte inner pte

32 10 30 12

3-level Paging

2nd pte offset outer pte inner pte

32 10 30 12

Inverted PT: PT w/ entry for every pg/F of mem: virt addr of pg at location, info on owning pcs (e.g. PID). Single PT used, fast find pc+pg num assoc w/ any F. Comp to standard PT: dec mem needed to store PT, inc time to search table (filter through entries w/ other PIDs). Can use hash table to limit linear search to one/few entries.

Hashed Page Table: Hash table to map pg nums to frame nums (normal PT is hashed with id as hashing func). PT contains linked list chains of pgs hashing to same location (more complex hash func reduces PT size at expense of conflicts). Search for match of virt pg num in chain. Extract phys frame if match found.

Segmentation: Paging gives 1-D virt addr space - what about separate addr space for code/data/stack? **Segment:** indep addr space from 0 to some max, can grow/shrink indep, support diff types of protection (r/w/x), unlike pgs, programmers aware of segments. **A** = good for shared libs. **D** = memory alloc harder due to variable size, may suffer external frag.

Concurrency: Single core = time sharing, Multi core = true, both non-deterministic. **Persistence:** File system, access controls, failure protection, organise devices. **Portability:** OS written in HLL (C), then compiled down as specific interactions. **Platform specific** = re-impl of platform independent = translate down directly. **Kernel Mode:** r/w access to all memory, full access to all instrs/hardware. **User Mode:** use kernel interface, limited access, isolated. **Splitting:** **A** = protection (kernel can cause damage), reliability (u-mode restarted without entire sys), security (isolation restricts malice). **D** = Performance (interrupt overhead when switch k/u), abstraction (restricted to kernel interface)

OS Zoo: **Multiprocessor** = Windows, MacOS, Linux, **Server** = optimise throughput, Solaris, Linux, **Mainframe** = bespoke, limited workload, OS/390, **Embedded** = trusted software, low power hardware (washing machine), QNX, **Mobile Devices** = power efficient, iOS, **Real-Time** = precision and time constraints, FreeRTOS

Concurrency: Higher utilisation, kernel shares CPU by switching processes, needs fairness. **Time switching:** illusion of parallel on single core. CPU time split in quanta (time slices), at each scheduler chooses which process gets quanta.

Context Switch: When CPU switches process (context = mem map, process state). **D** = overhead (process state saved + restored), caches (some invalidated, leading to misses)

PCB: Process Control Block, data struct storing process id, pointer to it's mem map, saved state, allocated resources. **Linux:** *task_struct*, **Process end:** normal, abnormal (error/exception), aborted (p killed c, never (daemon providing service) *ptstree* -> show process hierarchy

Thread Abstraction: Thread = stream of instrs being executed, state defined by CPU, stack per thd, addr

Signals: (Inter-process comm) similar to hardware interrupts. When get signal, call appropriate signal handler. All sigs except SIGKILL/SIGSTOP can be handled + ignored (not much info in them). SIGINT = keyboard interrupt, SIGABRT, SIGFPE, SIGKILL, SIGSTOP = suspend, SIGSEGV, SIGPIPE = p closed, SIGALRM = alarm timer, SIGTERM, SIGCONT. **Pipes** (pts): 1-way comms between processes (pcs). Connect stdout with other stdin. Can buffer, if pipe full, write end blocks until read. Temp unnamed pps can comm by pcs. Named pps stored in filesys (used like file), can output creating pcs. *int pipe(int fd[2])* - create pipe. Params: fd = where to put file desc (fd[0] = Read fd[1] = W). Returns: 0/-1. *close(fd[0])*

Per process items **Per thread items**

Address space	Program counter (PC)
Global variables	Registers
Open files	Stack
Child processes	
Signals	

Pre-emption: kernel interrupts/suspends task before it completes/finishes. Done via interrupts. (Handler saves process/thread state before potential context switch to other scheduled process.

CPU Bound Process: Pcs spends maj of its CPU time in CPU. Time to run limited by CPU performance. **I/O Bound Process:** waiting for I/O, briefly use CPU to issue I/O reqs. TTR limited by len I/O wait.

involvement in thd ops e.g. creation, switching, synch), customisation (each app has own scheduler for that prog). **D** = blocking (any block syscall on 1 thd blocks all in same pcs), page faults (during interrupts for page faults all thds in pcs are blocked), pre-emptive (pre-emptive scheduling hard for u-ts). **Kernel Level threads:** kernel can manage individual thds with pcs. **A** = blocking (thds individually blocked), simpler (easier scheduler as kernel can get interrupts but must be general). **D** = performance (still cheaper than proc switching but more expensive kernel interactions, can be mitigated by thd pool areas), customisation (no app specific schedulers). **Hybrid:** take adv of thd-specific blocking, but cheaper ops. OS = kernel level thds. Program manages pool of kernel provided thds (scheduled by kernel). Can use user-level synch (faster).

Non-Pre-emptive Scheduler: Processes run until voluntarily yield. Software must be trusted. Bad for interaction. Good for batch systems. **Pre-emptive Scheduler:** Requires timer interrupts. On interrupt, kernel takes control and context switches. Good for interactive processes (smaller quanta = better interactivity).

Time Estimation: SRT and SJF need to estimate. Runtime rarely known in adv or knowable at all. Heuristics based on history (not always available/representative e.g. for I/O). User/Process provided estimates (inaccurate/malicious - can penalise incorrect estimates but adds overhead). Instead, most OSs use priority-based system. **General Purpose Scheduling:** favour short + I/O bound jobs (OS must determine nature of jobs + adapt to changes). Good resource util (after I/O uses CPU, descheduled) and short response times. Often I/O bound (get data) -> CPU to process data. **Priority Scheduling:** schedule based on priority -> always run job with highest ps, can be externally def (e.g. by user) or process-specific metrics (e.g. expected CPU burst). PS can be static/dynamic (changed during exec), SRT can be considered priority scheduling (less time = higher p).

Multilevel Feedback Queues (MLFQs): Each priority level has queue, each can sched each diff (usually RR). Highest p always running (pre-empt if higher comes along), ps recomputed periodically (uses aging / recent CPU time). Can only set niceness not priority. **A** = reactive (to changing behaviour of pcs/thds), all run (priority recomputation). **D** = inflexible (apps little control, ps no guarantees), warm-up (determining factors for estimation takes time), cheating (apps add meaningless I/O to get higher p), no donation (cannot donate priority - only based on feedback from running pcs/thds). **Lottery Scheduling:** each job gets number of tickets, each scheduling decision chooses random ticket to schedule, share of tickets = share of CPU time, tickets can be transferred between jobs, can use for resources other than CPU. **A** = meaningful (proportions tickets and CPU time), all run (no job starved, any tickets run), donation (exchange tickets if job blocked by another). **D** = proportions (meaningfulness of tickets relies on proportions - adding more affects all jobs), unpredictable (random, more freq decisions help but inc overhead).

Atomic Operations: typically single ASM instr cannot be interr e.g. reg/mem R/W. Use special instrs to ensure atomicity of ops. **Lock prefix** can be used to ensure atomicity, and **TSX (test-and-set-lock)** allows a boolean in mem to be read and set to 1 without caring about number of thds. *int sem_init(&sem_t *sem, int pshared, unsigned int value)* init sema with val 1. *int sem_getvalue(&sem_t *s, int *val)*, *int sem_destroy(&sem_t *s)*. Enforce ordering (init to 0, thd 0 start with down(), thd 1 start with up()) - 1 goes first as can't be downed from 0). (aka up() forces ME. 1 thd hold L at given time. If attempt to get L but already acq, block requesting thread. Only holding thd can release L. Methods: init, lock, down, Re-entrant Ls: allow L to be acquired many times by same thd. If holder attempts to lock again, will deadlock. **Spin Locks:** if expected wait low, busy waiting may be better than blocking threads. Locks busy wait (in while loop, while (TSQ==1-locked) = 0), no kernel involvement (no block/unblock syscalls), ensure check on acquiring is atomic, can run into priority inversion. **Priority Inversion:** when using priority-based scheduling algo, low-p threads can block high-p ones. LP holds resource which HP wants e.g. lock. HP tries to acq but blocked and busy waits, LP not scheduled as P is low so cannot release resource. Solution = **priority donation** -> HP donates to LP so it can finish and release resource. Aging can also let LP execute if it has been a while. Spin locks mean HP always ready to run as no blocking used, but then LP cannot release lock. **Read/Write Locks:** allow many readers to hold lock, or 1 writer. Reduced lock contention as readers don't block. Some impls allow R -> W upgrade or W -> R downgrade w/o release + re-acquire.

Memory Models: Sequential consistency (all ops happen in order in source code, besides ths non-determinism) as assumed above is one mem model. Typically a **weak memory model** is used (allows compiler and hardware optimisation to improve sys performance) which can be affected by: **hardware:** reordering instrs, per-core caches on multicore processors (cache coherence) and **compiler optimisation:** instr reordering. This can affect code when conditions are changed by reordering. Typically compilers+archs have features to enforce control over ordering e.g. barriers around which instrs can't be reordered.

Resource Allocation Graphs: Directed graph to model resource allocation. Cycle = deadlock. Resource -> owned by Thd/pcs. Thd/pcs -> wants to acquire resource

Communication Deadlock: lost messages result in Dlock. Cannot use other methods to prevent Dlock here. Use timeouts (wait set period for message, if none report failure+progress to other task and retry message).

LiveLock: Pcs/Thds not blocked, but system as whole does not make progress. Common pattern: attempt to acq res, then block off+repre. No thds blocked, sys continues attempting and failing progress. Often seen with spinlocks as do not block when conds for Dlock met. Hence they are live-locked.

Memory Management: needs to provide allocation and protection. Req: no knowledge of generation/use of adrs for security (isolation to hide data/prevent mem corruption). **Mem Hierarchy:** regs (1ns), L1/L2/L3 (10ns), mem (100ns), disk/SSD (25µs), 250µs V/W, network (10ms-1s). L1+L2 per-core, L3 shared. Regs only store small amt + expensive.

Logical/Physical Address: L = gen by CPU, addr seen by comp. P = addr seen by mem unit, refers to phys mem num. Same in compile/load time addr-binding schemes (decide binding at compile/when prog loaded to mem), diff exec-time sch (dynamic decision).

Memory Management Unit (MMU): hardware device for mapping logical to physical adrs. Base reg = smallest phys addr, limit reg = highest logical addr, relocation reg = offset. Mmu ensures **base <= translated addr (base + rel reg) <= base + limit**. User pcs only deals with log adrs. Fast -> impl in hardware. MMU keeps memory of kernel (low mem) and other pcs protected.

Multiple Partition Allocation: Hole: contiguous section of unalloc'd mem. OS keeps track of alloc/free holes for when pcs are created/destroyed. **Dynamic Storage Allocation:** First-Fit: allocate first hole that is big enough (fast+simple), Best-Fit: smallest hole large enough (unless hole list sorted, search full list). Worst-Fit: largest hole (searches entire list, worse than other 2).

Address Translation: Mem addr split into page num (p) and offset (d). P = index of pg in pg table, gets base addr of frame. D = offset through pg/F, combined with base addr to get phys addr. For log addr 2^m and pg size 2ⁿ:

p	d
m-n	n

V -> P addr: V / 2ⁿ = page num, check PTE valid-dirty bits, PTE / 2ⁿ = F, add F to front of offset.

Memory Protection: Associate protection bits with each pg. **Valid/Invalid** bit: V -> associated pg in pcs' log addr space. 1 -> page not present (page fault, load in from swap), or incorrect access (by programmer?)

Page table implementation: PT kept in main mem. PTBR (PT base reg) points to start, PLTR (PT len reg) is size. Inefficient -> each data/instr access needs 2 mem accs (PT, data/instr). Use fast-lookup hardware cache as associative memory (supports parallel search). Called CAM (Content Addressable Mem), more expensive than RAM and specific to fixed mem alloc format.

Kernel Pg Access in Syscalls: when pcs makes syscall, handler runs in privileged mode so may need access to kernel pages. To ensure safe access, PT has supervisor bit (if set for translation, page can be accessed in kernel mode). Alt is to use separate PT for user/kernel.

Shared Memory: pcs can access shared mem by having pgs in 2 pcs point to same frame (imp/exp mapped to files on disk). After set up, no more kernel involvement. Useful for IPC-sharing libs (more common to use mmap). **Comparison to Pipes:** higher performance (less kernel intervention, data kept in space so less copying), bi-directional comms (would need 2 pipes), less useful for uni-directional comms (kernel has synch for pipes). **System-V API:**

shmget	Allocates a shared memory segment
shmat	Attaches a shared memory segment to the address space of a process
shmatctl	Changes the properties associated with a shared memory segment
shmdt	Detaches a shared memory segment from a process

Meltdown Attack: abuses speculative execution to access memory from cache side channel which wouldn't be accessible when running the code (e.g. to read kernel memory without correct privileges).

Effective Access Time: w = mem cycle time (ns), e = assoc lookup time, a = TLB hit ratio (related to num assoc regs), k = w/PT level + 1: **EAT = (e + w)a + (e + w)(1 - a)**. If hit, 1 mem lookup else x + 1.

Local Memory Layout: 1:1 Mapping → turn log addr to phys addr for kernel pgs by subtracting 3GB. V efficient for kernel mem access, no change of TL (no change in context when pcs switching so PT not flushed when user pcs makes syscall), on-demand mapping contains temp mapping for use of ~896MB of mem in remaining 128MB of virt mem. IA32: 4KB pg size, 4GB virt address space, 2-level PT (up to 3 w/ physical addr extension), offset bits contain page status (dirty, R only etc). AMD64/x86_64: larger pg sizes (e.g. 4MB), 48-bit adds, up to 4-level PT, offset bits can contain can-excc (prevent malicious code taking over process).	Demand Paging: avoid swap pages from swap when user attempts to access them. A = lower I/O load (unused pgs never loaded), less mem required (fewer pgs resident in mem), faster response time (no need to wait for all pgs to be loaded), supports more users (lower mem usage allows this). Valid/Invalid Bit: 1 = in mem, 0 = not, all pg entries init = 0, if pg on 0 accessed, page fault and trap to kernel. K uses table to determine if ref invalid (abort) or valid but pg not in mem. Valid req: 1. get empty block, 2. swap pg to F, 3. reset tables (valid bit = 1), 4. restart last instr.	Performance (Demand Paging): Pg Fault Rate $0 < p < 1$ (0 = no PFS, 1 = every ref is PF). Effective Access Time (EAT) = $(1 - p) * \text{memory access} + p * (\text{pg fault overhead} + \text{swap pg out} + \text{swap pg in} + \text{restart overhead})$. Virt Mem Tricks: Copy-On-Write (COW): pcs accessing identical pgs use same frame, only copy when one wants to write/modify it. Parent and child init share same pgs in mem. Efficient pgs creation (copy only modified pgs). Free pages alloc from pool of 0'd out pgs. Using fork(): C's PT points to S's pgs (marked read-only in both PTs). Protection fault causes trap by kernel. Kernel alloc new copy of pg to process altering, replaces old one in PT. P and C PT sets page to R-W. Memory Mapped Files: Map files into virt addr space using paging. Only load parts of file when accessed. Simplified programming model for I/O (easy access to stdin/out).	Page Replacement: When out of free mem+new pg needs creating, find unused pg to swap out. Victim does not have to be pg of same pgs. 1. Access pg not loaded in mem (PT bit shows this), 2. update PT (ensure no RCS when involving multiple pcs on multiple cores), 3. write victim back to disk, 4. read req pg from disk, 5. restart op. Goals: reduce num pg faults (In general, more frames = fewer PFS), prevent over-alloc of mem (PF service routine should incl pg replacement), reduce redundant I/O (use dirty bit to only load modified pgs back to disk). Replacement Algos: FIFO: replace oldest pg. A = simple to impl. D = may replace heavily used pg. Belady's Anomaly: more frames → more page faults (see graph). Optimal Algorithm: Replace pg which will not be used for longest time period. D = impossible in practice (can be used as benchmark to compare other algos tho). LRU (Least Recently Used): each pg has counter, every time referenced counter clock to ctr. When replacing, choose pg with lowest ctr. Proper LRU is expensive (search for lowest ctr-storing ctrs) so use approximations instead: Reference Bit: each pg has ref bit = 0 (init), when pg referenced $r = 1$ (done by MMU), periodically reset bits, when evicting choose pg with $r = 0$. Second Chance: Evict pgs in RR, but if $r = 1$ give second chance to stay in mem. Checks circular queue of pgs to evict (like RR). If pg to be replaced (in order where $r = 1$, set $r = 0$ and try another pg). Counting Algos: keep counter of num of references to each pg. LFU (Least Frequently Used): replace pg with smallest counter, may replace very new pg just brought in, never forgets heavy pg usage (reset counters/age aging). MFLU (Most Freq Used): rep pg w/ largest cnt. Newly accessed pgs have low count → prioritised, pgs barely used hog mem+heavily used pgs evicted quickly. Locally vs Global Pg Repl: Local: when repl pgs, pcs can choose a pg belonging to same pcs. Reqs keeping track of changes in working set size + some partitioning to determine how many pgs each pcs owns/can be alloc'd - 1. fixed partitioning (each pcs gets fixed size), 2. balanced set algos. As pcs manage pg faults indep, more scalable than global. Global: Can choose pg from any pcs. Mem div shared between pcs. Init alloc mem prop to pcs size. Use pg fault freq to tune alloc (more PFS → more alloc). Use all pgs to determine repl. More efficient than local. No best sol - Linux = global, Windows = local. Linux Page Repl: variation of clock algo to approx LRU pg repl. Active list → contains active pgs, MRU near head of list. Inactive list → approx LRU near tail. Only repl pgs in inactive list. Ref bit determines when pgs moved between A+1 lists. kswapd (swap daemon) → pgs in inactive list reclaimed when mem low, uses dedicated swap partition/file, must handle loaded/shared pgs. pdflush (kernel thread) → periodically flushes dirty pgs to disk.
Locality of References: progs tend to req same pgs across space+time. Design for this to ensure good performance. If we don't, could res in thrashing. Thrashing: excessive paging causing low processor util. Prog repeatedly reqs pgs from secondary storage, destroying performance. I/O Device Management Objectives: Fair access to shared devices (prevent pcs hogging res, limit dedicated devices), exploit parallelism (can use devs in par - e.g. send packets using NIC while writing to disk - allows for multiprogramming), provide uniform-simple I/O view (abstract devs from pcs, uniform naming+err handling, hide complexity of dev handling). Device Independence: dev indep from it's type (e.g. terminal, disk, DVD drive) and instance (disk 1, 2, or 3). Device Variations: unit of data transfer (char (bytes) or block), supported ops (R, W, seek), synchron or async (s = disk, as = NIC), speed differences (NVMe SSD vs tape drive), types of err concs (disk errs vs GPU °C warning).	Working Set Model: working set of pgs $W(t, w)$ = set of pgs ref'd by a pcs while running from time interval $t - w$ to t . Can use this in clock algo, keeping track of pgs in working set by adding "time of last use". At each pag fault, $t = 1$ → set $r = 0$ and move to next pg, $t = 0$: a.c. alg, e.g. if age < w (working set age) → move to next pg, c. age > w → if pg clean, replace, otherwise start write-back+continue to find another pg (once w-b of frame complete, pg marked as free+clean - don't want to halt while waiting for pg to write back to disk). Effectively, only repl pgs if haven't been ref'd in w time, to take adv of temporal locality. Working Set Size (w): processes transition between different working sets (transition has higher pg frames alloc to pcs, then const again in new WS). Do not dealloc pgs ref'd within w time (OS keeps some outside WS hanging around). In transitions num alloc pgs tend to be higher as pgs from 2 sets referenced within w time. WS too large → too many allocs, pcs use too much mem. WS too small → too few allocs, many pg faults, slow. Don't want to misalloc (OS alloc for pgs prog doesn't want anymore). If many faults → alloc more pg frames (interrupt time inc). When alloc all pgs in pcs, no pg faults so IFT = excc time. Graph (x = num pg frames alloc to pcs, y = IFT) follows S curve (top stops at (num pgs in pcs, total excc time)).	Character vs Block Dev: C = minimise latency (keyboard/terminal - mem = dev file for phys mem of OS, pty = pseudoterminal - bidirectional comms channel, USB = USB devices). B = maximise throughput (disks/NIC - randisk = access RAM disk in raw mode, fd = floppy disk, loop = loop device - maps data blocks to file in filesegs, or other block dev).	I/O Layers: User-level I/O Software, Device-independent OS software, Device drivers, Interrupt Handlers, Hardware. Interrupt Handler: Interrupt is sig from dev to CPU that dev needs attending to (dev connecting, finish reading, error etc). Drivers register handlers to deal w/ diff interr types. Block: on transfer completion, sig dev handler. Character: when char transferred, process next. Modern Systems - Hybrid: (e.g. nvme storage) polling (queue of reqs+resps) very fast but uses CPU time, wait for interr (better for longer waits as can sched other pcs). Device Driver: handles type of dev, can control mult of same type. Impl block/R/W, access dev reqs (write control info to file), init ops (e.g. start dev at boot), sched reqs (if dev shared), handle errs. Dev-Indep OS Layers: standard interfaces for drivers of dev types → simps OS des, interface to write new drivers to, no OS changes to supp new drivers (just new interfaces). Provides dev indep: map log to phys devs (naming+switching, map 1:m any l:p or p:l), req validation against dev (check dev driver working correctly), alloc (which pcs can access which devs), buffering (performance+block size indep), err reporting.
Buffered I/O: Output → use data transferred to OS output buffer. Pcs control (only suspend when buffer full). Input → OS reads ahead, reads taken from buff, pgs blocks when buff empty. Smooths I/O traffic peaks (limited load balancing), avoid diff data transfer unit sizes between devs (buff contains blocks). Unbuffered I/O: data transferred directly between dev+user space, each R/W causes phys I/O (dev does smth, not just buffer), dev handler used for every transfer, high switching overhead (every R req driver to take over+do phys action).	Device Drivers: Mem Mapped I/O: dev addr as mem location where can use virt mem setup to restrict acc (supervisor bit). Ways to do I/O: Programmed: simple but ineff, wait for device (spin) then continue excc. Interr Driver: large overhead, good for long expected waits. Hardware sends interr when op complete, do other work while waiting. Direct Mem Acc (DMA): requires hardware, reduces CPU intervention. DMA controller (often in dev) waits for dev to respond, then once full res avail, places direct into mem.	Device Alloc: Dedicated Device: (DVD writer, terminal, printer) - 1 pcs gets exclusive access, if another tries to acc it fails (can keep queue of open reqs), typically alloc for large periods of time, only alloc to authorised pcs (avoid malicious pcs blocking access to res). Shared Device: (disks, window terminals) - OS provide sys for sharing e.g. files/share to use disk. Spooled Device: blocking user access to alloc (non-shareable) devs causes delays+bottlenecks - spool to intermediate media (disk file). Daemon pcs has dedicated access, pcs send reqs/jobs to daemon which provides sharing on non-shareable res, and red I/O time (greater throughput).	User-Level I/O Interface: syscall interface to allow user progs to interact, often with 3 rd party libs. Basic I/O ops (close, R, W, seek), set up params (dev indep), can be sync (blocking) or async (non-blocking). UNIX files: accesses virt devs as files (standard I/O calls). FD Name → 0:stdin, 1:stdout, 2:stderr. /dev/kdb for keyboard access.
Loadable Kernel Module (LKM): Dev drivers loadable modules, loaded+linked dyn with running kernel. Required binary compatibility (mod specific to kernel ver). Kmod: kernel subsystem manages mods w/o user intervention, determines symbod devs, loads mods on demand. init_module call code, cleanup_module clean shutdown. Kernel can open file, look for symb table (gen by compiler) and call corresponding funcs. sudo insmod same, sudo rmmod same. I/O Management: Device classes: group sim types of dev. ID numbering: major = det which driver is controlling, Minor = distinguishes dev of same class. Special files: mod dev rep by dev/special files in /dev dir. Device Access: accessed via virtual file system (VFS). Most drivers impl R/W/seek but all have other ops too. ioctl/cdrom_CROMEJECT_0 syscall supports special tasks. Character Dev I/O: transmits data as stream of bytes. Represented by char_device_struct (driver name, registered maj+min nums, ptr to driver's file_operations struct via cdev). All registered drivers ref'd by chrdevs vector, file_operations ops supported by dev driver, stores funcs called by VFS when using syscall access dev special file. Block Dev I/O: Block I/O Subsys: several layers, modularised ops (common code in each layer). 2 main strats to time access block devs: 1. caching data, 2. clustering I/O ops (store up tasks, excc many at once). Direct I/O: bypass kernel cache when accessing driver, good for databases (+others) where caching can reduce performance/consistency (e.g. vals rarely accessed more than once). Linux API: I/O Classes: character (unstructured, files+devs), block (structured, devs), pipes (message, interprocess comms), socket (message, network interface). Sockets: can be local or across network (pipes use machine-specific fds so cannot go over network). Types: TCP, UDP. fd = create(filename, permission), fd = open(filename, mode) mode = 0 (R), 1 (W), 2 (RW), close(fd), numbytesread/write = read/write(fd, buffer, numbytes), pipe(&fd[0], newfd = dup[oldfd]), ioctl(fd, &termios), fd = mkfifo(filename, perm, dev). File Descriptors: each pcs has fd table, 3 fds when created (0:stdin, 1:out, 2:err) all refer to terminal from which prog was started.	RAMR: (Heat Assisted Magnetic Recording) to ensure high precision, platter can only be written to when heated (laser heats platter, cold sections ignore R/W ops).	Blocking: process suspended until op complete, I/O call only returns once complete (appears instantaneous from prog perspec). Pcs making call could be blocked for a while, threads get around this (complex). Non-Blocking: I/O call ret as much as is avail. Constantly R until get all data (could return none sometimes). App-level polling for I/O (constantly req R from stdin, or loop if only data at moment). Turn on using fcntl syscall for sending comms to manage fds. Async: pcs runs in parallel w/ I/O ops (non-blocking). When ops complete pcs notified, pcs can check/wait for I/O completion. Flexible+efficient, more complex code, potentially less secure if buffs mismanaged.	Disk Scheduling: FCFS: reqs completed in order received. A = best for lightly loaded disk (time between reqs larger than time to fulfill any req), fair (no bias). D = low performance for heavy load. SSTF (Shortest Seek Time First): order reqs on shortest seek dist from curr head pos. D = biased against inner/outermost tracks (middle on avg closer), unpredictable performance, pcs can use dummy reqs to keep control of head, can delay reqs indefinitely. SCAN: elevator - select reqs w/ shortest seek time in preferred direction. Only change dir when at inner/outermost cylinder (no dir reqs in dir). Base for most common algos used. D = same delay issue as SSTF (but reduced - only in 1 dir), long delays for reqs not in dir of algo/on extremes. -SCAN: SCAN in 1 dir, jumping to start (innermost) when at the end (outermost). A = lower variance of reqs on inner tracks, largely red indef wait arising from SSTF. N-Step SCAN: only services reqs waiting when sweep began (for each sweep). 1 sweep = entire → outer → inner. A = reqs arriving during sweep serviced before end of sweep (no long waits), no indef waits poss. RAID (Redundant Array of Inexpensive Disks): Disk performance not kept pace with CPU performance → RAID inc disk based sys performance by using many disks in parallel. Arr of phys drives appears as single virt drive, distribute stores data over disks to allow parallel operation (imp performance). Use redundant disk capacity to respond to disk failure (more disks → lower mean time to failure). Levels: 0: (diag 1) spread blocks in RR on disks, concurrent seek/transfer of data (for blocks on diff phys disks), sometimes balance load across disks, no redundancy, 1: level 0 but all disks duplicated. Rs serviced by either disk. WS must update both disks in parallel, failure recovery easy, low space efficiency, high cost. 5: (diag 2) most commonly used, distributes parity so potential for concurrency. Some potential for W concurrency as parities on diff disks, good storage/efficiency tradeoff, reconstruction of disk non-trivial+slow.
Disk Caching: can cache sectors of disk in main mem to reduce access times. Buffer in main mem contains copies of disk sectors, OS manages disk in terms of blocks (likely much larger than sectors so loads multiple). Must ensure contents saved in case of failure (lazy writing comp). Cache has finite space → need replacement policy. LRU (Least Recently Used): replace block in cache longest with no refs. Cache = stack of pointers to blocks in mem, when ref'd pushed to top of stack, algo evicts bottom of stack. D = doesn't track num of accesses, only ref time. LFU (Least Frequently Used): replace block w/ fewest refs - each block has counter. Some blocks ref'd many times in short period → misleading ref count → use freq-based repl. Frequency-Based Replacement: Divide LRU stack into 2 secs - new+old. Block ref → move to top of stack, only inc ref count if not already in new. D = blocks access out too quickly → use 3 secs, only replace blocks from old. File Attrs: Basic: name, type, organisation, creator. Addr info: volume, start addr, size, used, size alloc'd. Access Control Info: owner, authentication, permitted actions. Usage Info: creation timestamp, last modified/read/archived, expiry date, access activity. File Syscalls: fd = open(file, how...), s = close(fd), s = read(fd, buffer, nbytes), n = write(fd, buffer, nbytes), pos = lseek(fd, offset...) - move fp, s = stat(pathname, &buff) - get file metadata, s = fcntl(fd, cmd...) - file locking/other ops.	Filesys Organisation: Space Alloc: File size variable (can inc/dec) and is allocated on disk in blocks (512-8192 bytes). Block size determined by filesegs. BS Too Small: high overhead for managing large files (many blocks to keep track of), high file transfer time (more blocks = more seeking back+forth). BS Too Large: internal frag (small files wasteful), caching based on blocks → large cache space/unable to cache many blocks. Contiguous File Alloc: Files at contiguous addrs on storage device. A = successive logical records typically phys adj. D = external frag, poor performance if files grow/shrink, file grows beyond init size specified and no contig free block avail → transferred to new area of adequate size → many additional I/O ops. Block Linking/Chaining: Each file contains linked list of blocks. When locating block, chain searched from beginning (O(n)), if blocks dispersed through disk searching is slow. Waste pointer space in each block. Insertion/Deletion by modifying pointer in prev block. Block Allocation Table: Uses dir mapping files to first block. Table maps blocks to next block in file, indicating free spots (block with no next is EOF). A = file allocation table (FAT) can be cached in mem for quick lookup, no lengthy seeks to traverse block nums for file. D = files can become fragmented reducing R/W speed → needs periodic de-fragmenting (expensive), FAT can become impractically large using lots of mem (>1 block of storage). Index Blocks: Each file has >=1 index blocks containing list of pointers to file data blocks (effectively page table for files). File's dir entry points to index block. Chaining: reserve last free entries in index block to store ptrs to more index blocks. A (over simple linked-list indep) = searching in index blocks themselves, index blocks placed near corresponding data blocks → quick data access, cache index blocks in mem.	Directory Syscalls: s = mkdir(path, mode), s = rmdir(path), s = link(alldpath, newpath) new hard link, s = unlink(path), s = chdir(path), s = opendir(path), s = closedir(dir), dirent = readdir(dir) R 1 entry from dir, rewinddir(dir) rewind dir to re-read. ext2fs: the second Extented File System is high-performance and robust (formerly standard, remains robust tho). Uses block sizes typically 2^a 10, 11, 12, 13. 5% block reserved for root (safety mech to ensure root pcs can always run even after malicious ones use all disk space). ext2 inode used to rep files+dirs, uses 12 direct pointers, 13 th indirect, 14 th double, 15 th triple (fast access to small files, also allowing for large files). Block Groups: clusters of contiguous blocks. FS attempts to store related data in same BG. Reduces seek time for accessing groups of related data. BG structure: [superblock group descs block alloc bitmap inode alloc bitmap inode table data blocks]. Buses: Bus 1: High bandwidth, low latency devices. Devs that can function at same speed as RAM. Bus 2: High bandwidth, medium latency devs. Devs have high throughput+req low latency but slower than mem (e.g. graphics cards, network). Bus 3: Other. Slower devs/reqs with lower bandwidth (e.g. block devices - disk).	SSD: Sched: no sched algo as many mem modules can be R/W in parallel, and R/W speed approx constant. Drivers need to overcome issues with limited WS, tracking virt to phys blocks-assignment of free blocks. Comp to HDDs: high bandwidth (1GB/SSD vs 100MB/HDD), lower latency, high parallelism. File Systems: Organise info. Main obj: non-volatile+long-term storage, sharing info, concurrent access, convenient organisation, easy management of data, security. File Types: hard links (file aliases data of another file), soft links (aliases path to a file by another - in command), regular (bat, sh, lib, zip, docx...), directory, char special (access to char I/O dev), block special (same but for block devs). Filesystem Support Functions: name translation (conv paths to disks-blocks for log use by driver), management of disk space (alloc/dealloc storage for files), file locking for exclusive access (important when concurrency - fcntl syscall), performance optimisation (caching/buffering), protection against sys failure, security. File User Funcs: truncate (erase contents but keep all other attrs), reposition/seek (set curr pos to given val), R attrs (creation date, size, archive flag), W attrs (protection, immutable flag). Filesys Directory Organisation: map symbolic names (text.txt) to logical disk locations (Disk 0, block 234). Helps with file organisation+prevents naming collisions (all names in given dir unique). Hierarchical FS → root = where root dir begins, root dir points to various dirs each having own entries for it's files. Link: ref to dir/file in other part of FS (allows alt names/def locations in tree). Hard Link: ref to dir or file (not dirs). Soft Link: ref full pathname of file/dir, created as dir entry. Problems: file deletion (search for links+remove, leave links and cause exception when used (symbolic links), keep link count in file+delete when cnt=0, looping → if file points at self/loop of others, they never get deleted). Directory Operations: open/close, search (pattern match on strings), create/delete files, link/unlink files, change dir, list files, read/write attrs, mount. Mounting: combine (make links to) multiple FSs into one namespace. Allows ref from single root dir, support soft links to files in mounted FS (not hard as they are dependent on FS impl). Mount point = dir in native FS assigned to root of mounted FS. FSs manage data. With mount tables (info abt location of mount points/devices - when native FS encounters mount point, use MT to determine dev and type of mounted FS). Directory Representation: dirs map symbo names to inodes (can be other dirs - typically only single block or files).
Free Space Management: When alloc new block for file, need to quickly determine which is available. Free List: linked list of blocks containing locations of free blocks. Bs allocated from start of free list, newly freed blocks appended to end of list. A = freeing+alloc fast (O(1)). D = files unlikely to be contiguous alloc → when reading files have to seek across many random locations on disk which is slow. Bitmap: 1 bit in mem for each B, 1 bit = 1B. A = uses little mem (single bit/entry), quickly determine contiguous blocks at locations, highly optimised bit ops as standard on most CPUs. D = may need to search entire bitmap to find free spot (O(n)).	Filesys Layout: boot block, super block, free inode bitmap, free block (zone) bitmap, inodes+data. Super Block: contains crucial info about filesegs: num inodes, num data blocks, start of inode-free space bitmaps, location of 1 st data block, block size, max file size.	Multi-Level/Hierarchical Page Table: Allocating non-pg sized mem can be hard so nice if ind ptrs of pg table fit in 1 frame each (partially what multi-level/PTs achieve). Single level wastes lots of mem.	Access Control: Principals = users/ops. Authentication: Personal characteristics: (key based on user, usually biometrics/signature). A = hard to forge, convenience. D = special hardware (expensive), false +ves/-ves. Possessions: RFID cards, implants, secure keys, ZFA. A = convenience. D = lost/stolen (impersonation), sometimes expensive (complex/locks/card scanners). Knowledge: secret known to user (password, with freq turnover). A = cheap, secure if secret. D = password reuse diminishes security, only as secure as password storage, dictionary attacks. Authorisation: who can access what and how. Principle of Least Privilege: user gets min reqs to carry out task. Sometimes more given by default/for convenience. Access rights defined as objects with operations permitted on them. Principal exccuting in domain D has access rights specified by D. Access control matrix (col = obj, row = principal). Access control list → 2D arr too large to store, so store each col as list + associate with object. Capability list = rows of matrix as lists assoc with principals.
UNIX/Linux File Access: users are principals, each has UID. Can belong to +1 groups. UID = 0 = root (all access rights). All files are objects, can belong to at most one group. -rw-r--r- 1 prp lds 2517522 Mar 9 09:14 test.txt → file type (- = reg, d = dir, b = block, c = char, l = symbo link, s = pipe, s = socket), owner rights (group rights, e = rights, owner, group name, permission bits) = none, r = read, w = write, x = excc, s = setuid/did, t = sticky bit. File excc w/ privs of user exccuting. SUID set → use owner's privs. Each user has effectively 3 UIDs: 1. Real UID = ID of user that started pcs, 2. Effective UID = ID of exccuting process, 3. Saved UID = UID the effective UID can be switched to.	Capability Lists: capabilities are protected pointers to objects, specifying permitted operations. Kernel owns capability info, user uses fd as a key of sorts for indirect access to info.	Security: Goals: data confidentiality (theft of data), data integrity (destruction/alteration), system availability (DoS). Policy = Mechanism: P = what security if provided (what is protected, who has access, what access permitted), M = how to impl P. Same M can support diff P's. People Security: attacks by insiders (abuse privileges), social engineering (phishing, blagging, shouldering), convenience (not changing pwds often enough), ignorance/lack of knowledge. Hardware Security: phys access (damage machine, read/alter contents of disk, snop on f/wire network traffic), hardware flaws (side channel attacks e.g. cache attack, timing attack, data remanence, or badly impl access control). Software security: buffer overflow, integer overflow, string formatting. Can gain root privileges/crash app/steal data/deny access to sys.	DAC vs MAC: DAC = Discretionary Access Control = principals determine access rights to their objects (default). MAC = Mandatory Access Control = system rules/policies determine ARs. Bell - La Padula Model: MAC policy where info doesn't travel down security hierarchy. Pcs only read at it's security level or lower, only write at same/higher. Biba Model: opp.
Design Principles: Least privilege = pcs run with lowest priv poss, default no access. Simple-uniform mechanism, psychologically acceptable (if too complex/burdensome ppl will not impl), public (can be critiqued for flaws).			