

<p><b>THE - Dijkstra;</b> <u>Main point</u>: sequential processes(central abstraction) in layers; <u>Why</u>: handle complexity; <u>Benefits</u>: <i>verify soundness, prove correctness</i>; <u>Layers adv</u>: modular, abstraction, debug, verify; <u>disadv</u>: difficult to design(circular dependencies), separate layers, inefficient (call through layers); <u>Syn</u>c: Semaphores - mutex, waiting queues (cooperating processes); No circular wait because uniprocessor system; priority scheduler; <b>Monitors(Hoare)</b> - only one process in a monitor procedure at one time; condV; coarse grained, guaranteed cpu access to support os; deterministic; if(not invariant) wait(c)</p>	<p><b>Nucleus - Per Brinch Hansen;</b> <u>Motivation</u>: freedom of design not in existing OSes; <u>Main point</u>: <i>Small nucleus (kernel) supports &gt; 1 simultaneous OS tasks</i>; <u>Msg adv</u>: Easy to validate, explicit wrt semaphores, simpler; <u>disadv</u>: buffer management, memory is not shared, copying overhead, overflow; <u>Layers</u>: Nucleus (schedule, communicate, proc manip primitives), user; <u>Syn</u>c: msg passing; Scheduler: Round robin (slow for large n_children); general and flexible</p> <p><b>Protection - Butler Lampson;</b> control operates on domains, used for revocation; owner operates on objects</p>	<p><b>Hydra - Wulf, Cohen;</b> <u>Main point</u>: capability-based OS nucleus <u>Goal</u>: construct environments to effectively utilize HW resources Separate mechanism (protection) &amp; policy (security), rejects hierarchy (layers of privilege, not flexible); <u>Rights types</u>: kernel, auxiliary; <u>Amplification</u>: trust the procedure implementor; <u>Capability list</u>: 1 user multiple files; <u>Adv</u>: fine-grained, flexible right amplification; <u>Disadv</u>: most entries = no access allowed, wasting space; expensive search, hard to revoke, difficult to set up; <u>ACL</u>: single file, multiple users, easy mgmt; capability and type matching done for params in procedure</p>	<p><b>Tenex - Bobrow, Burchfiel;</b> Virtual “machine” = virtual memory; CLI, terminal interface Backward compatibility - Old sys call to compatibility package (Jsys) Working set principle - which processes to keep in memory based on inter-fault time (time b/w 2 page faults of a process); not good - apps should not page in and out; <u>Vm features</u> - sharing (page appears in addr space of &gt;1 proc using direct pointers in VM maps), sharing large portions by CoW; File name <u>fields</u>: device, directory, name, extension, version; FS - versioning, 5 level “tree” (symbolic name for file consists of &lt;= 5 <u>fields</u>, conceptually represents a tree of max depth 5)</p>
<p><b>MULTICS - Jerome Saltzer;</b> Right amplification by rings; Contributed protection systems, domain transitions, multilevel security policies; Permission rather than exclusion (default = not allowed), access check to each obj, principle of least privilege; no group access (as opposed to UNIX); Descriptors, segments, gates, rings - <u>still used</u> in x86; <u>ACL</u>: easy to set up, slow to check (walk through the list); <u>ACL + CL</u>: protection checked at open; FD returned with capability to R/W; HW holds, checks descriptors; restrict portability; <u>mmap</u></p>	<p><b>Medusa;</b> multi-user distributed centralized OS; Break into components (disjoint utilities); utility - distributed among Cm nodes for parallelism file system invocation is mostly remote; spin waiting (not good in uniproc, but lock can be released by other proc in DCS so efficient - short critical section); <u>pros</u>: easy to maintain, simple; <u>cons</u>: SPOF, performance, memory consumption; <u>communication</u>: msg pipes; <u>unsealing</u>: kernel obj associated w application, remote access done by first mapping it as XDL =&gt; unsealed; then manipulate; then reseal (remove from XDL) - failover; coscheduling</p>	<p><b>Plan 9;</b> specialized nodes for a particular task, more cost-effective. Everything is a file; all servers = file servers. Local namespaces - relative path names, custom NS for each user, referencing local names more familiar to users. Union combines directories from multiple NSes. <u>9P</u> - protocol to connect clients with file servers; implements RPCs for all file methods. Storage: snapshots, no backups; new level WORM (slow retrieval, not feasible today b/c large multimedia files daily). <u>UTF-8, rfork, ./n/dump</u>; Multithreading: rfork (specify shared, pvt resources, kernel thread)</p>	<p><b>UNIX;</b> unifying abstraction - file abstraction for data, I/O. No special HW support - portable Unified namespace - mount; File - linear sequence of bytes;fork optimization - CoW; set-user-id: right amplification (invoke privileged ops w/o giving privilege to user); System does buffering (reads, writes) - 2 buffers - system buffer, library buffer <u>Shell</u> - forked instead of exec b/c shell proc can crash if cmd crashes, so safer to let child execute (process isolation)</p>
<p><b>Pilot;</b> single user, single address space and language support; limited features for protection and resource allocation; no FS, hints from applications; processes cooperative not competitive; Mesa and Pilot interdependent; no protection against malicious procs, more against errors (lang based); accepts hints from apps; circularity between files and VM; kernel/manager (mechanism/policy); <u>mmap</u></p>	<p><b>Mesa;</b> notify places process on run queue, but does not switch control to it; reason for change - performance(avoid extra context switches in Hoare), remove scheduling from inside the monitor; Deadlock(circular and un-notifiable wait), priority inversion addressed by temp priority inheritance; thread.stop - threads check to see if they should shut down by polling variable</p>	<p><b>Sprite;</b> network OS, name and location transparency, caching network FS; <u>app interface</u>.- extends Unix; single uniform NS, process migration, shared memory; <u>kernel</u> - multiprocessor, RPC; (longest) prefix tables; <u>FS caching</u> - server (reduce delay caused by disk access), client(reduce no. of calls to non local disk), read reuse and client write; sequential - versioning, concurrent - caching disabled; files as backing store;double caching</p>	<p><b>Distributed V;</b> msg-based IPC for diskless WS and large file servers; no significant penalty for remote access; sync msgs (static, fixed size kernel buffers), small fixed size msgs (reduce queuing, buffering), separate data transfer facility (efficient transfer large data); <u>inefficient</u>: short msg inefficient use of large packet size, sync com prevent I/O and computation overlap, separate data and control increase nw ops; <u>remote ops</u>: implemented in kernel, raw ethernet frames, no per-pack ack</p>
<p><b>Microkernel:</b> IPC, VMem, scheduling;<u>adv</u> easy to make it reliable and secure, more stable, extensible, configurable; <u>disadv</u>: lots of system call and context switch; <u>L4 address space</u>: delegate pager responsibility, recursive; <u>Linux on L4</u>: runs as user process, syscalls from process to LServer, LServer acts as pager for user process; <u>binary compatibility</u>: shared lib(emulation layer between LAPI and LServer), trampoline (syscalls in statically linked unmodified Linux libraries reflected back to emulation layer); <u>sp ex evaluation</u>: pipes/RPC, Vm ops, cache partitioning</p>	<p><b>FFS;</b> improvement on UNIX FS, larger block size(4096), more bandwidth utilization, inodes and data closer(locality); throughput tied to free space, read as fast as write(In FFS, blocks are placed on disk in a much better layout, even with sync reads, placement roughly corresponds to a good ordering of disk requests, writes have higher overhead because they must allocate new blocks) UFS writes faster(async writes, reorder requests to minimize seek time); extent-based w/ journaling</p>	<p><b>LFS;</b> FS designed to exploit HW and workload <u>trends</u>: large memories, disk BW scaling but latency is not, smaller file access; FFS problems: sync writes, possibly related files far from each other, inodes placed separately from files; <u>challenges</u>: how to find data(indirection, inode map in log, pointers to imap in checkpoint, keep in memory), free space; <u>adv</u>: cold segments cleaned at much higher use; <u>disadv</u>: inodes appended to log, difficult to find, inode location not fixed in LFS so not as easy as FFS to calculate, manage free space (GC competes with disk read/write); crash recovery: checkpoints</p>	<p><b>Exokernel;</b> <u>main point</u>: customize OS interface / abstractions - move everything to user level. <u>Motivation</u>: performance suffers from generality, can’t customize/extend; <u>Approach</u>: OS layer exports HW resources directly, OS functionality at user level in untrusted library OS. FS part of LOS (no protection, memerr can corrupt on-disk data structures) - address by SFI; track resource ownership by tables, protection (secure bindings) techniques - HW (protected data sharing via VM HW), software TLB on top of HW TLB, packet filters for downloaded code; <u>disadv</u>: no protection within VAS; resource revocation - if LOS doesn’t cooperate, “abort protocol” forces</p>
<p><b>Grapevine:</b> distributed message(email) and registration service; <u>msg service</u>: delivery, buffering of msgs; <u>reg service</u>: naming, authentication, access control, resource location for clients; independent of each other and use IP to communicate; <u>client user package</u>: name /address transparency; client-server- ethernet, server-server - low B/w modem;<u>scaling issue</u>: distributed lists, solve by indirection; transparency:distribution and replication; <u>why not transparent?</u> Updates are not transactions, duplicate msg delivery</p>	<p><b>Xen;</b> high performance VMM supporting strict resource control among guest OSes, OS that exposes vHW interface; isolation among VMs, heterogeneous OSes simultaneously; paravirtualization - guest OS modified; x86 CPU easy to virtualize because of rings, memory difficult because HW TLB and no tagging =&gt; V-P-M, insert V-M, write protect PT; syscalls routed to guest OS via Xen w/o mods, handled so that no trap to Xen necessary, exceptions propagated to guest OS by Xen events, async notifs; guest OS modification supports event handlers</p>	<p><b>Rio;</b> write-back performance with write-through reliability; mem perf (eliminates synchronous writes) with disk reliability (protects file cache during normal operation, restore contents on warm reboot - system reset, power not lost); <u>protect file cache in kernel</u>:vm support, read-only file cache; protection domains separate from address spaces, separation of file buffer cache into separate protected module; more reliable than write-through - faults corrupting file cache pages caught with protected file cache; <u>fault injection</u>; checksum, memTest</p>	<p><b>Soft Updates:</b> order operations so that metadata is consistent in case of crash; create file: write inode safely to disk then update directory block; soft update = enable write-back caching of metadata; collapse multiple updates before going to disk; transfer multiple updates to disk in one write</p>
<p><b>MACH:</b> maintain all VM state in machine independent module; large, sparse VAS, CoW, memory mapped files, user level pagers and backing store;<u>DS</u>: address map, mem obj, resident PT, pmap (machine dependent); easy to share mem regions; different regions can have diff pagers; <u>shadow objects</u>: CoW creates a copy only if <i>extremely</i> necessary (only modified pages)</p>	<p><b>VAX:</b> OS is essentially just an extension of user address space of every process, OS can directly access user code and data, kernel stacks in user address space; command interpreter is a part of each user's address space; <u>optimizations</u>: local page replacement(FIFO),page caching,clustering; free list: clean pages evicted from a process resident set, modified list: dirty pages: second chance algo</p>	<p><b>GMS;</b> probability based (imp for DCS), clusters of computers act like tightly-coupled microprocessor than LAN, globally shared local resources; <u>goal</u>: minimize avg data reference time (go to remote MM instead of disk); maintain page age, find page in cluster, make replace and evict decisions</p>	<p><b>Cells;</b> resource constraints on smartphones, existing virtualizations require OS mods; container virtualization (OS namespace) - FS (own root partition, no files owned by other VPs), process (only procs running in own NS), NW (IPs, port mappings), device (/dev); foreground, background VP, if fg doesn’t acquire exclusive access, HW shared by bg VPs</p>
<p><b>Scheduler Activations;</b> vessel to execute user thread context, mechanism to notify user of kernel events; stores user thread state in kernel. <u>Design aspects</u>: upcalls (add proc, proc preempted, SA blocked/unblocked), downcalls (hints: add more procs, proc idle), critical sections (on preemption/unblock: recovery - run thread till end of CS; deadlock free); SA not reused to notify b/c contains exec state of thread; goals: other ULTs can run during blocking, no proc idle when ULT ready, no priority inv</p>	<p><b>Lottery Scheduling;</b> probabilistic mechanism for proportional CPU share; relative rate; scheduling; <u>tickets</u> - abstract, relative, starvation free. Transferable (inheritance via tickets) - solves priority inversion, inflation/deflation, currencies(hierarchical allocation + keep allocation isolated along trust boundaries), compensation tickets - temp inflated value; proportional share is goal; can't express response time differently from share</p>	<p><b>Taintdroid;</b> android;identify sensitive data, taint and track data flow (vars - local, args, static, classes, arrays; msgs - upper bound of tainted vars in msg; methods; files - same as msg), real-time monitor behavior of running apps, identify misuse of pvt data; works even if data encrypted; significant performance overhead <b>VM370:</b> Control program(CP)-VMM, CMS (Conversational Monitor System) - guestOS; FS - disk; single dir, no hierarchy, sharing(CP to mount), I/O(record); mem model single addr space, RSCS(store/fwd nw router)</p>	<p><b>MapReduce:</b> Reporter: provided for MR to report progress,indicate that they are alive; inefficient for multipass algorithms, no efficient primitives for data sharing; intermediate written to local disk, output files written to DFS; utilize data localization by running map tasks on machines with data, one becomes master and assigns tasks to idle. Tasks scheduled based on location of data, map fail - rerun; Workers are periodically pinged by master, writes periodic checkpoints, avoid straggler</p>
<p><b>GFS;</b> Failures are norm, append writes, co-design apps and FS API, automatic failure detection, tolerance, recovery, <u>concurrent</u> appends by multiple clients, 64MB chunk servers: intra-rack BW &gt; inter-rack BW, garbage collects orphaned chunks, migrates chunks b/w chunkservers, file data not cached on chunkserver or client, append at least once, separate data and control flow: client interacts directly w chunkserver, master server state logs stored on disk, atomic metadata changes by single master: <u>consistency</u>, client deletes file: master records deletion in log, renames file to hidden name, removes every 3 days <u>Disadv</u>: single master, bad latency, multiple/small file support weak, relaxed consistency</p>	<p><b>BigTable;</b> distributed multi level map, fault tolerant, persistent, scalable, self managing; <u>NoSQL benefits</u>: auto sharding(splits a single dataset into partitions), replication, integrated caching; tablet - row range, unit of distribution and loadBal, split as table increases; single row transactions(atomic read, modify, write sequence), no support for general transaction across keys, <u>components</u> - GFS (persistent data storage, SSTable), Chubby(distributed lock manager, master election), master(load balancing, garbage collection), tablet servers, lib linked to client (read/write directly with tablet server); 3 level (B+ tree); empty cache - 3RTT, stale cache - 6 RTT; locality groups (Segregating columns families that are not typically accessed together enables more efficient reads) and compression</p>	<p><b>HayStack;</b> goals to serve photos: high throughput, low latency, fault-tolerant, cost-effective, simple; <u>adv</u>: reduced disk I/O, simplified metadata, single photo serving and storage layer; haystack cache uses DHT, add a photo to cache if request directly from browser not CDN, photo fetched from write-enabled store machine; index file allows quick loading of needle metadata without traversing larger Haystack store file; cookies insufficient protection; using column family / tagging helps for albums</p>	