





Common case: request fits in current last chunk: Primary appends data to own replica; Primary tells secondaries to do same at same byte offset in theirs; Secondaries writes with success to client.

**Consistency model:** consistent, defined (not concurrent write), defined not consistent (append)

**Atomic record append:** master guarantee at least once (unique record ID for filter), master chooses offset

**Snapshot:** copy-on-write (lease revoked thus further commit, goes via master)

**Conclusion:** Fault tolerance → constant monitoring & replicating crucial data (operation log, checkpoints, shadow master & fast and automatic recovery, backup, checksum (64-bit) to detect corruption, high aggregate throughput → sparse data transfer, master involvement is minimized → large chunk & chunk lease

Limitation: single master, bad latency, not too many files or small files, relaxed consistency

**Data model:** Renew: Master grants lease to primary (typically for 60 sec). Leases renewed using periodic heartbeat messages between master and chunkservers.

**Version number:** Chunks have version numbers. Stored on disk at master & chunkservers. Each time master grants new lease, increments version, informs all replicas.

**Master Recovery:** checkpoints & roll-forward for master in mem-metadatas.

**Data recovery:** notice via heartbeat, decrement chunk count, re-replicate chunks. Single-master. Good latency. Cannot support too many files or small files. Relaxed consistency

**Big table:** *dist* multi-level sorted map.

**Main point:** Dist storage sys for managing structured data, various access modes, flexible access patterns (e.g. random, resemble database (not fully relational, client dynamic control, data as uninterpreted strings)) → sparse dist, persistent, multidimensional sorted map

no SQL → less constrained consistency → simple design, horizontal scaling, control over availability

**Data model:** row, column, timestamp, map, cell, rows [ordered, atomic], cols [two-level, ACL], queueifier [row, lookup methods, desc sorted]; Table row range, fine-grained LB. API: change metadata, manipulate data in complex way (single row transaction, cell as counter, Sawzall script). **Access:** (passive table, active table, monitor table of table servers thru Chubby dir, LB, gc, schema migration). Tablet server [handle rw req, re-split, Chubby [dist] lock service, provide namespaces with directory and small files, lock as node registration with lock lease, meta, dir, data, a tablet], **Tablet file:** (table file, table file (highly-available and persistent), and registers itself by getting a lock in a specific directory Chubby. Chubby gives "lease" on lock, must be renewed periodically. Server loses lock if it gets disconnected.

**Metadata:** Client do not rely on the master for tablet location information, thru Chubby instead

**Tablet:** contiguous range of column families, 100-200MB, fast recovery(each server picks 1 tablet), fine-grained load balancing

**Tablet location:** chubby file → root tablet → other METADATA → usable (3 of 1 empty, 6 of failed)

**Tablet assignment:** tabletserver holds a lock in chubby (master monitors the directory), master restart using chubbylock, server directory, add root tablet

**Tablet split/change METADATA and notify master:**

**Tablet serving:** tablet log (commit) → memtable, read merge memtable and SSTable, read/write authorization from Chubby file, compaction when memtable is full (major compaction: no deleted data → timely);

**Session Lease:** each clients has a session with Chubby. The session expires if the client fails to renew the lease within its session expiration time.

**Refinement:** locality groups (Can group multiple column families into a locality group. Segregating columns families that are not typically accessed together enables more efficient reads), compression (long common string, repetition in small windows), 2-level caching (key value pair, SSTables blocks), bloom filters (Stable contains the data?), compact tablet recovery with minor compression, exploit immutability (SSTable, split share parent data)

**MapReduce:** Goal: lib to *abstract parallelization*, tolerate fault, dist data, LB. Programming model (restricted → easier to implement): high perf. on clusters of commodity PCs.

File based. M map task partitions, R reduce task partitions [hash(key)/key].

**Locality:** GFS provide redundancy, master schedule task close to map machine (important, network bandwidth is scarce).

**Performance:** detect and detect re-execute completed k in-process map (in and in-process reduce), completion committed thru master. **Master fail tolerant:** checkpoint.

**Perf:** file replica, backup task mechanism (to avoid last few steps take too long of network redundancy)

**Refinement:** partition of function (same host together), ordering guarantee of files (fast random access), combiner function (output intermediate file for other tasks), custom input type (reader interface), side-effects (auxiliary files), skip bad records (master input), local execution (debug), status info, counters.

**File locations:** a. Map input: from dist FS

b. Map output: to local FS at Map node; c. Reduce input from multiple remote FS; uses local FS. Reduce output to dist FS → Local FS: Linux FS → Distrib. FS: HDFS

Evaluation: easy to use, problems expressible in MapReduce, scale well

**HayStack:** Previous: POSIX (directory and file metadata), NFS-based (gc operations are bottleneck, caching not helping due to long-latency content). **Goal:** relational database, append-only (updates are not needed), fast reads, fast fetch from H, add to cache if not from CDN and write-enabled store. 3. H store [manage multiple phy vol on dist machines, multiple photos as a single large file to reduce meta, *logical vol id=offset-photo*]

**Assumption:** long lat, recently updated most reads and likely deleted, rare modification and deletion. **HayStack:** reduce dk, IO, simplified metadata as kv., single photo severing and storage layer.

**Architecture:** 1. H dir [in-mem mapping logical vol to phy, lb, determine cache hit/miss, dramatically low cost, higher throughput, incrementally scalable

2. H dir, add to cache if not from CDN and write-enabled store. 3. H store [manage multiple phy vol on dist machines, multiple photos as a single large file to reduce meta, *logical vol id=offset-photo*]

**HFS:** superblock (updates are not needed, fast reads, fast fetch from H, add to cache if not from CDN and write-enabled store. 3. H store [manage multiple phy vol on dist machines, multiple photos as a single large file to reduce meta, *logical vol id=offset-photo*]

HFS: superblock (updates are not needed, fast reads, fast fetch from H, add to cache if not from CDN and write-enabled store. 3. H store [manage multiple phy vol on dist machines, multiple photos as a single large file to reduce meta, *logical vol id=offset-photo*]

**Optimization:** compaction (copy and swap), save more in mem (offset 0 as deleted, no cookie), batch update as possible

**Performance:** dramatically low cost, higher throughput, incrementally scalable

**Pro:** Reduced disk I/O, metadata easily cached in mem. Simplified metadata → kv instead, no directory structures/file names, results in easier lookups

Single photo severing and storage layer, direct I/O path between client and server, results in better bandwidth

**MicroK:** Traditionally, problem with monolithic K-hard to extend, heavyweight, complex and huge, hard to maintain

MicroK exports and implements only basic critical resrc abstractions: only IPC (e.g. msg), VMem (addr space), Threads Scheduling. **Pro:** Multiple OS namespaces → no need to share the same machine at the same time. Specialization → specialized imp of general functionalities (e.g. different perms); Extensibility → extend features not provided by base imp (e.g. user-level pgr/VMem, cache partitioning).

**Linux:** exports devices mapped into addr space; interrupts wrapped in msg. **Binary comp:** shared Emulation layer between Linux API and Linux server

trampoline [syscalls in statically linked unmodified Linux libraries are reflected back to emulation layer]. **Tagged TLB:** problem → TLB flush when context switch, → cache is local for proc, when proc changes, cache should be invalidated, solution → Tap proc on TLB to flush

**ExoK:** **Design:** LiBOS [user-level, OS functionally, untrusted], **K no param check**, thus faster. **Malicious behavior** only affects the single app itself that using libOS rather than entire machine. No share: ExoK hard to share FS, each app has own imp of FS and lib call, thus no hardware. No protection: FS mem errors by app can potentially corrupt on-disk data structures. **ExoK:** multiplex & protection; rescsrcs multiplexing [resrc ownership, multiplexed hw via allocation and revocation], protection → secure binding [hw: protected data sharing via VMem hw. Caching: sw TLB cache on top of hw TLB. Downloaded code, no packet filters. **ASH:** ASH packet filters: ASH

App-specific Safe Handlers) is downloaded into K, to do packet filter. A packet comes in, ExoK analyzes it, and sends it out to the right App. This is more efficient than the ExoK querying all the server processes on every incoming packet (as in LiBOS).

**FS in ExoK:** apps know better than OS the goal of resrc mgmt thus apps make decisions. FS done in LiBOS inside app's address space. **safety concern (ASH):** limited sharing (app development is not... OS interface the same) **Hide & Expose 1** ExoK exposes the imp., exposes the mgmt of hw resrcs to user-level OSes) to maximize their performance

2) ExoK exposed all hw resrcs to user-level processes. ExoK supported a streamlined sys call interface by which apps could request resrcs (mem) and express mgmt policy (how to map virtual pgs to physical pgs). It also supported cache (user-level to device) and cache (pg faults) and to allow processes to express preferences for revocation decisions.

3) One concrete example is virtual mem. ExoK controlled which processes were allocated which physical pgs, but those processes determined which of those physical pgs mapped to which virtual pgs they were given. On a pg fault for example, ExoK would invoke a callback to user-level asking which physical pg to map to the undefined virtual pg mapping (TLB translation). The user-level process would then be able to specify any of its physical pgs to map (i.e., it would decide its own policy). ExoK would validate the mapping (i.e., whether the process owned the physical pg) and then install the TLB entry to provide the mapping.

4) **Problem 1** that ExoK experienced is that calling back to the user-level process on each pg fault was expensive (crossing the protection boundary, deciding which pg to evict). ExoK solved this problem by implementing a software TLB → a cache of translations provided by user-level processes. This cache dramatically reduced the number of user-level callbacks. **Problem 2** was demultiplexing (dispatching) incoming network packets to their respective user-

level processes. Upon receipt of a packet, ExoK would not know which process to deliver the packet to (e.g., without asking the processes). ExoK would then send authorized use a given resource, access checks are no longer required. Therefore, the secure bindings separates binding a resource to an app from accessing that resource →

By isolating the need to understand these semantics to bind time, the K can efficiently implement access checks at access time without understanding them.

Use case. This type of operation is best when an app continuously uses a resource it has requested authorization for. If an app decides to continuously ask for different resources, more binding checks will be required which reduces the benefit of the secure binding mechanism.

**VM370:** CP. Control Program, computing machine to simulate multiple copies of the machine. CP is a general multiprogramming sys that uses virtual machines to organize input job streams. Each virtual machine is identified by an entry in the CP directory, interface to virtual machine, allocation and mgmt of resources. Domains are virtual machines. Each virtual machine is run by one person. CMS is a dk-file-oriented OS to support the personal use of a dedicated computer, data sharing (safe) is enabled via links made by CP (shared mem of all channels) **RSCS:** Remote Spooling and Comms Subsys, information transfer and management, interface with comm facilities **RCS:** is an information transfer, multi-tasking sys that uses a dedicated computer with attached comms equipment for the support of data file transfer among computers and remote work stations.

**Xen:** virtual machine hypervisor that is similar but not identical to the underlying hw, modifies OS to make it VM-aware and perf tuned). **VMMonitor** exposes hw interface [runs on raw hw, runs on host OS]. **Usage:** OS dev & debugging, sw testing, utilize hw, isolation, virtual I/O for IoT (locality, security & isolation, developer control, low cost)

**ExoK:** exoK schedule domains are virtual machines. **VM:** Virtual Machine. Fast handler (hypervisor not involved): pg fault and sys call; can register handler with VM.

**Privilege instruction:** VMm tests user's mode, emulates VM's Syscall instructions. Dom and Dom0, which is the user's OS, runs on host OS. **Protection** domains ring 0, 1, 2, 3 → Xen modifies the guest OS to run on ring 1. **Hypervisor:** extend machine interface with hypervisor calls. Control transfer: hypervisor, event (asynchronous notification from Xen to domains); Xen in top of each address space

**TLB flush:** 2 level PTE (difficulty in x86, hw manages PT, software not involved in TLB mste, TLB not tagged). Problem → x86s do not have a software managed TLB. Its TLB is not tagged, which means that the TLB must be flushed on a context switch. **Solution** → a. guest OSes are responsible for allocating and managing the hw PT, with minor involvement from Xen to ensure safety and isolation; b. Xen exists in a section at the top of every address space, thus avoiding a TLB flush of Xen when entering and leaving the address space

**Privilege:** the OS allocates mem, pg, registers with Xen, call Hypervisor for write req, validated by Xen, guest batches the update requests. **Hypervisor Mgr:** Dom0.

**Lazy:** An example of lazy evaluation in the Xen operating system is GuestOS (Xen) and Dom0, which is the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

Xen tries to make the exception registration process as fast as possible. Therefore, exception handling validation does not do safety checks such as verifying that the handler code is present in mem (which could be expensive). **Dom0** and **Dom1** are the user's OS, runs on host OS. **Protection** domains ring 0 (highest instruction privilege). Guest OSes which Xen virtualizes operate on ring 1, but their code has been written to function on the highest privilege; as such, Xen has to simulate some registers which are only available on ring 0. It has to handle exceptions such as segmentation faults and allow Guest OSes to register exception handlers with the hw.

**File-level:** Persistent storage tracked at the file level. The design stores data tag per file. The taint tag is updated on file write and propagated to data structures. **Reliability:** Reliability is ensured by the use of redundant, low-bandwidth Sensors → place hooks in Android's LocationManager and SensorManager apps. High-bandwidth Sensors → place hooks for both data buffer and file linting for tracking microphone and camera information. **Privacy:** Privacy hook is implemented by the use of a unique phone number, SIM card identifiers (IMSI, ICC-ID), and device identifier (IMEI). Network Taint Sink → identifies when tainted information transmits out the network interface.

**Scheduler Activations:** **User vs. K threads:** *perf* on context switch, flexible customized. **Pro:** *negation:* K threads block, resume and are preempted by user-level threads. **Cons:** K threads scheduled obliviously to user-level state; when user threads block, K threads also blocks. **Solution:** → a new K interface (straightforward modification) and user-level thread package (programmer sees no difference) → functionality, perf., flexibility → Share necessary control and scheduling information between K and app address space

**Scheduler activation:** 1) Vessel for executing a user thread context. 2) Mechanism for notifying user-level of K events (upcall, downcall). 3) Data structure for storing user thread state in the K. **Virtual Processor:** a. Each process supplied with virtual multiprocessor; b. K allocates processors to address spaces; c. User level threads sys (ULTS) has complete control over scheduling; d. K *upcalls* ULTS whenever it changes the number of processors, or user thread blocks and another thread. 2. IO completion, K preempts processor, takes the preempted SA to form a new fresh SA as the context of upcall.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

**Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL. **Example:** 1. IO threads block in K; K upcall in the context of a fresh SA. UL.

Reading data is where LFS does not perform well, an extra layer of caching could only be beneficial to the overall performance.

**Soft updates:** Reliability, can waste some inode or blocks, but data pointed cannot be corrupted or uninitialized.

**Soft updates:** write-back caching of metadata, avoid reliability-induced write → *perf* (metadata write overhead), integrity (handle possible corruption), fast recovery (instant recovery, no log replay) → app view is up-to-date. **Soft updates:** Reliability is consistent (dependency is resolved). **Operands:** direct, indirect. **Operands:** 1. Never point to a structure before it has been initialized (e.g., an inode must be initialized before a directory entry references it).

**2. Never reuse a resrc before nullifying all previous pointers to it (e.g., an inode must be nullified before it is allocated to a new block, or a block must be reallocated for a new inode).** 3. Never reset the last pointer to a live resrc before a new pointer has been set (e.g., when renaming a file, do not remove the old name for an inode until after the new name has been written). To leave an inode in a consistent state [-] point to uninitialized inode, okay for dangling inodes.

**Undo/reredo operations:** soft updates tracks dependencies on a per-pointer basis and allows blocks to be written in any order. To solve circular dependency, any still-pending updates in a metadata block are rolled-back before