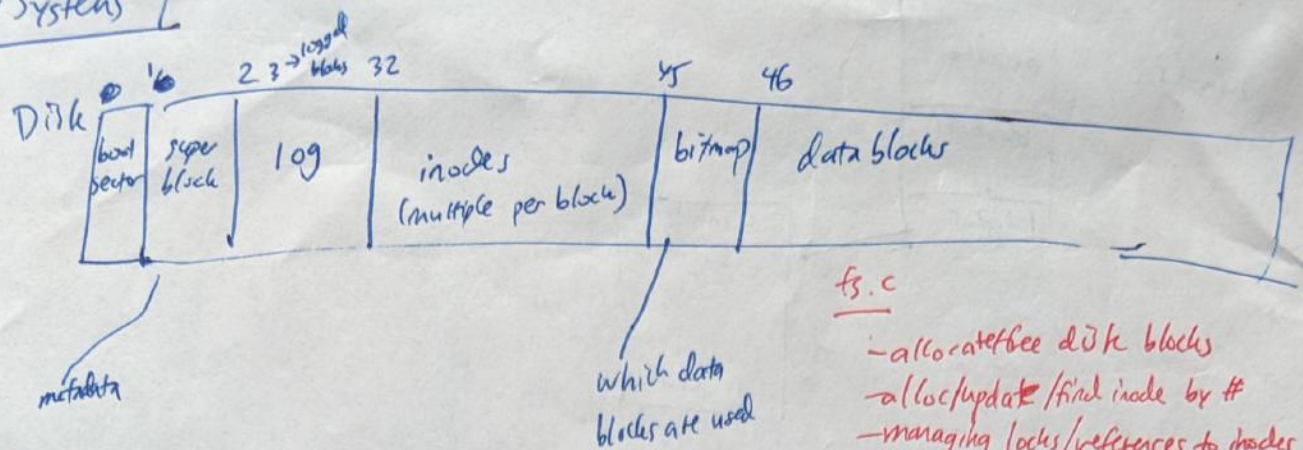


File Systems

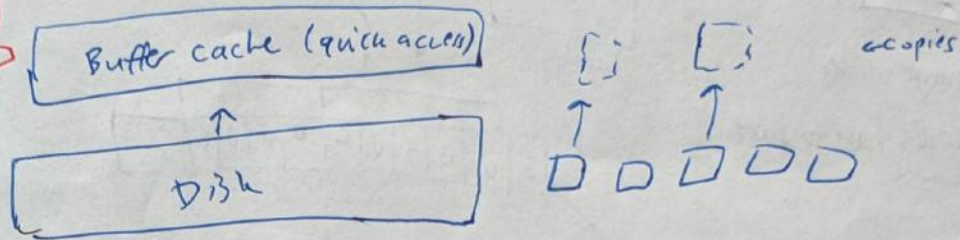


fs.c

- allocate free disk blocks
- alloc/update/fini inode by #
- managing locks/references to inodes
- map blocks to inodes

bio.c

- binit creates linked list of buffers
- each buf has a blockno, lock
- bget is called to get data
- brelse - releases a locked buffer
- bwrite: write changed data to disk before releasing buffer
- balloc - alloc new disk block



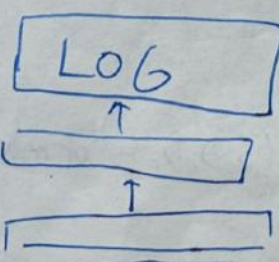
Ex →

- 1) write data to new file
 - 2) write to block 33 for allocating new inode
 - 3) write to block 46 for writing directory entry
 - 4) write to block 32 - update directory inode
 - 5) write b 45 to alloc new block
 - 6) write block 33 for adding new block to inode
- Annotations on the right side of the list:
 - alloc (next to item 2)
 - write (next to item 3)
 - alloc (next to item 4)
 - write (next to item 5)
 - write (next to item 6)

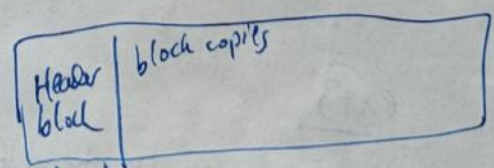
PROTECTIONS:

- * bcache.lock protects cache description
- * b->lock protects one buffer
- * b->refcnt prevents buf from being recycled

XV6 logging



description of all disk writes it wishes to make
writes commit record for each complete operation



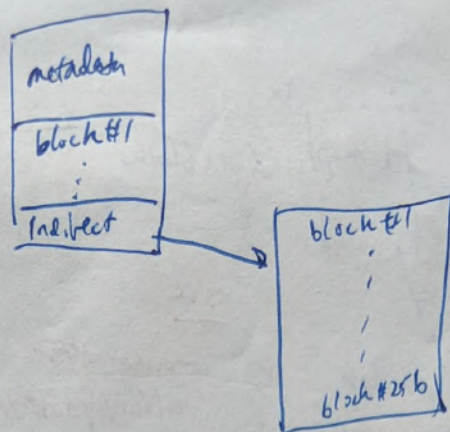
LOG/block
[sector #s for each log block] # log blocks
0 = no txn
1 = completed committed txn

log.c

- log-write writes to a buffer in the buf cache and phs it and appends block # to sector arr
- commit()
 - write modified blocks to log from buf
 - write header (commit pt)
 - install writes to actual loc
 - erase txn
- begin/end op

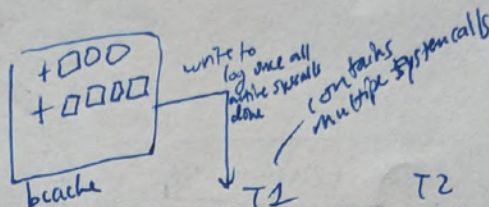
- PROS - correctness bc write ahead
- good throughput
- CONS - blocks written twice
- ops don't fit in log?

inode

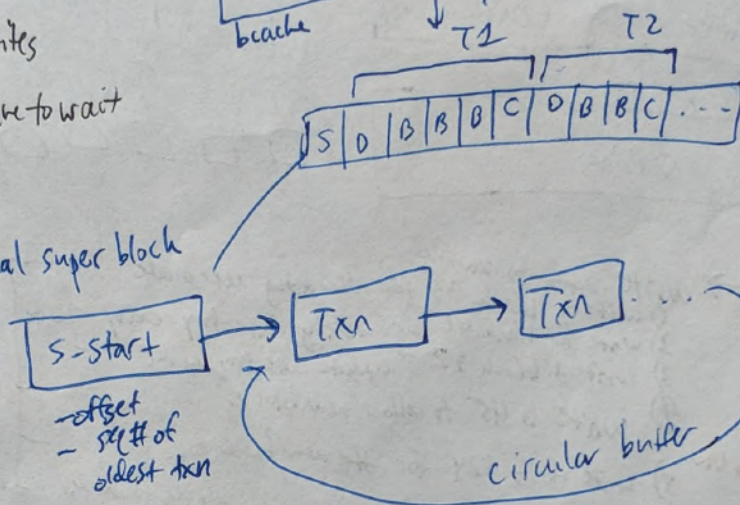


Problem w/ Logging SLOW

- 2x commit writes
- new syscalls have to wait

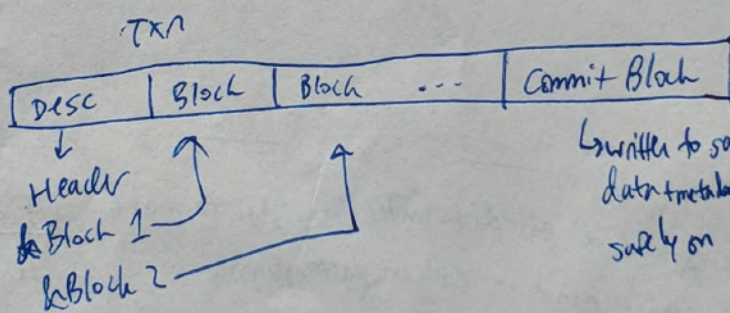


Journal super block



CRASH

- 1) Start from oldest seq #
- 2) Find end of log
- 3) Scan until missing commit or unexpected seq #
- 4) Return to last valid commit block, replay all blocks through last valid commit



Differences

xv6

- every syscall is a txn
- log is an array of blocks/txn

ext3

- log is just a circular stream of blocks (S, D, B, C) of multiple txn
- less immediate availability

VM optimizations - user controlled virtual memory (faster)

3

TRAP - handle pg fault traps *sigaction*

PROT 1 - decrease accessibility of a page *mprotect*

PROTN - decrease avail of N pages (mark os vm structs as protected) *mprotect*

UNPROT - increase accessibility *mprotect*

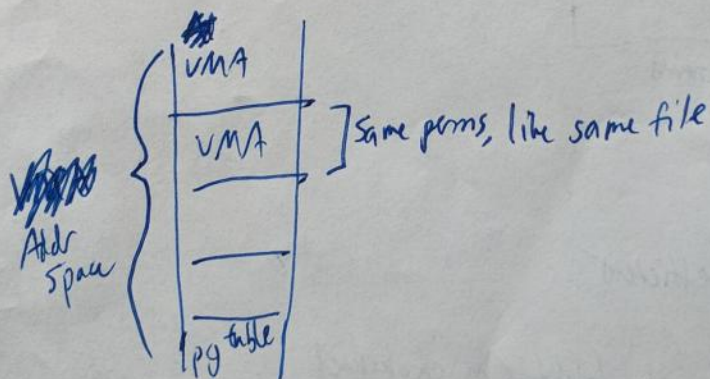
DIRTY - returns list of dirtied pages

MAP2 - map phys page at 2 diff VAs *protect levels*

mmap()

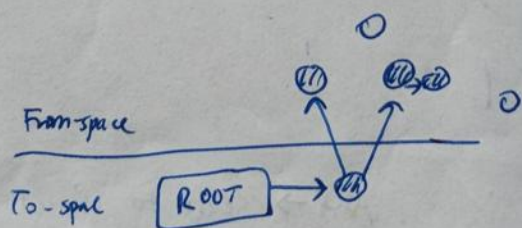
- map mem into addr space
- mprotect (PROT_READ, write)
- munmap

Process:



GC

Concurrent:



How to avoid race conditions?

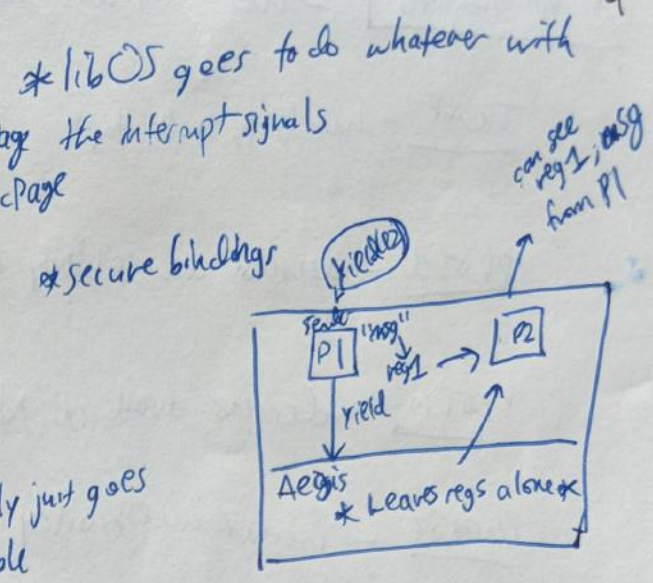
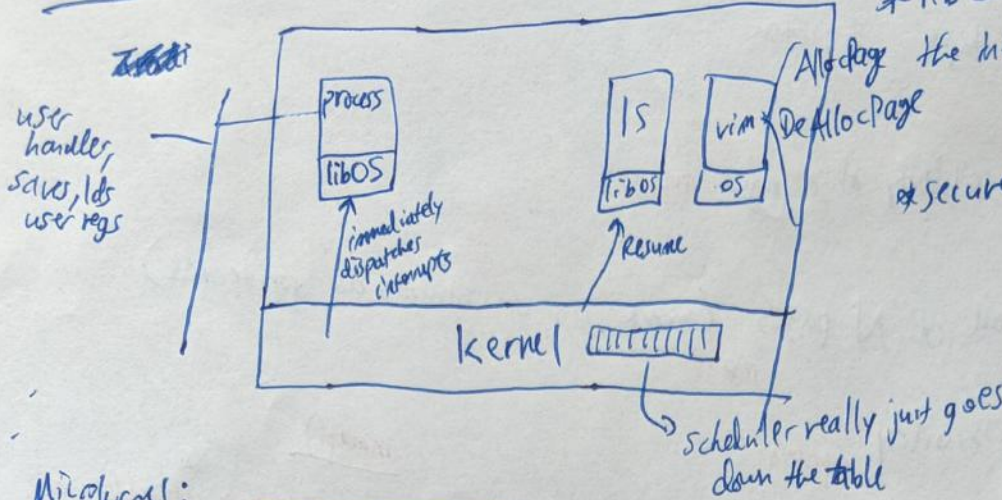
* Split To-space into Scanned R/W
Unscanned fault on access

real protect: *pg fault*, GC immediately scans, all its neighbors, UNPROT

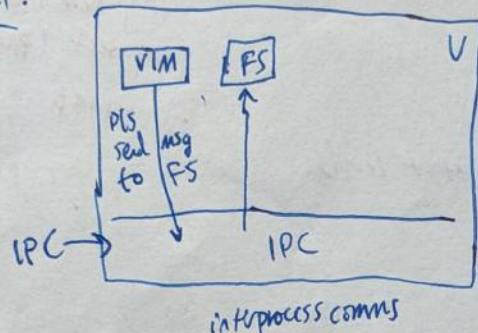
Generational: separate by how recently they were allocated
- youngest Gen is more often looked at
- also moves ptrs around in space

Exokernel - How abstracted does a CPU need to be?

4



Microkernel:



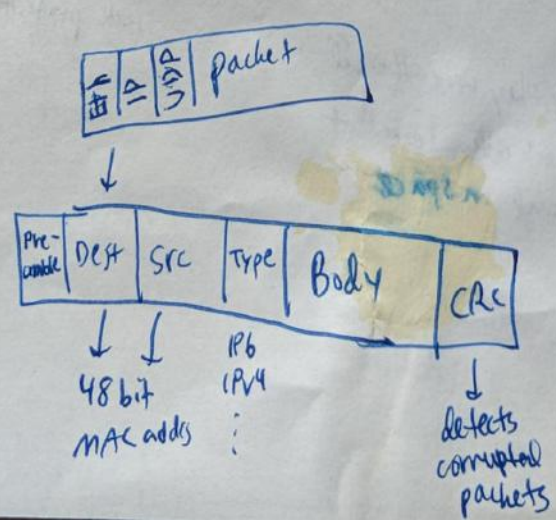
Revocation

- while revocation - libOS guides process

Secure bindings

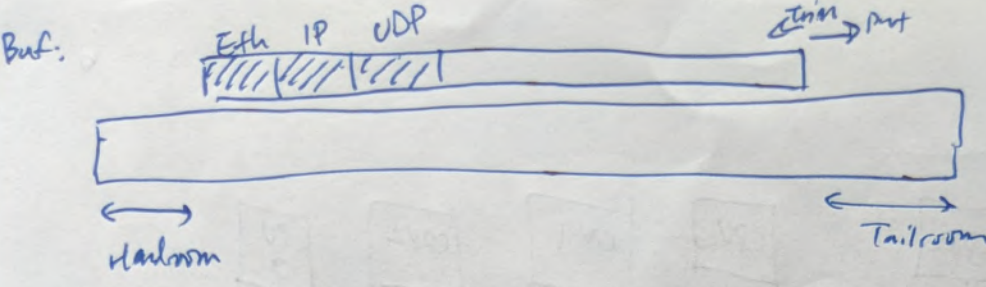
- Hardware: low level, efficient
- Software caching - cache bindings in exokernel
- downloading code into kernel

Networking



ARP → send to IP addr... but which MAC
 if in same network:
 use MAC addr
 else:
 use MAC of gateway
 - Gateway receives, broadcasts request for src MAC/IP, dest IP
 - Target responds w/ MAC

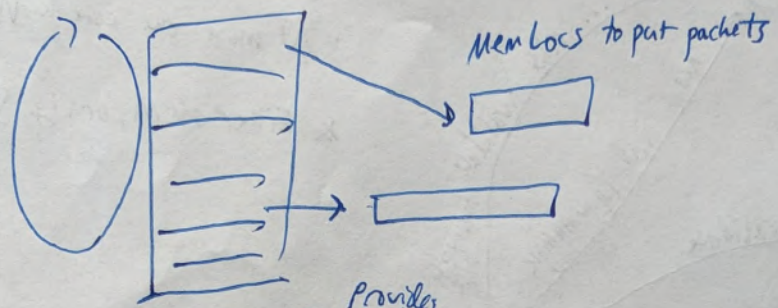
datagram (UDP) | stream (TCP)
 - note passing
 - order, data guarantee not preserved
 - faster (less overhead) | phone call
 - reliable
 - order
 - slower



Sockets

- Datagram sockets (unreliable)
- Stream sockets (bidirectional p/p's)

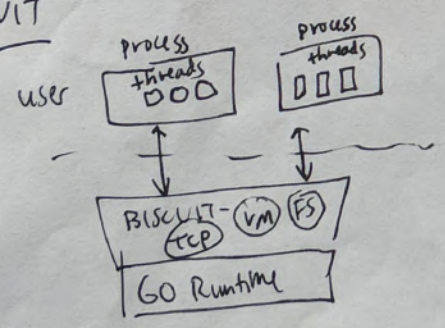
ex receive why empty? DROP
ex transmit why empty? wasted bandwidth



OS → must constantly poll or interrupt for new packets

- | | |
|--|---|
| <ul style="list-style-type: none"> - predictable tasks - unnecessary waiting | <ul style="list-style-type: none"> - unpredictable tasks - wasteful when frequent |
|--|---|

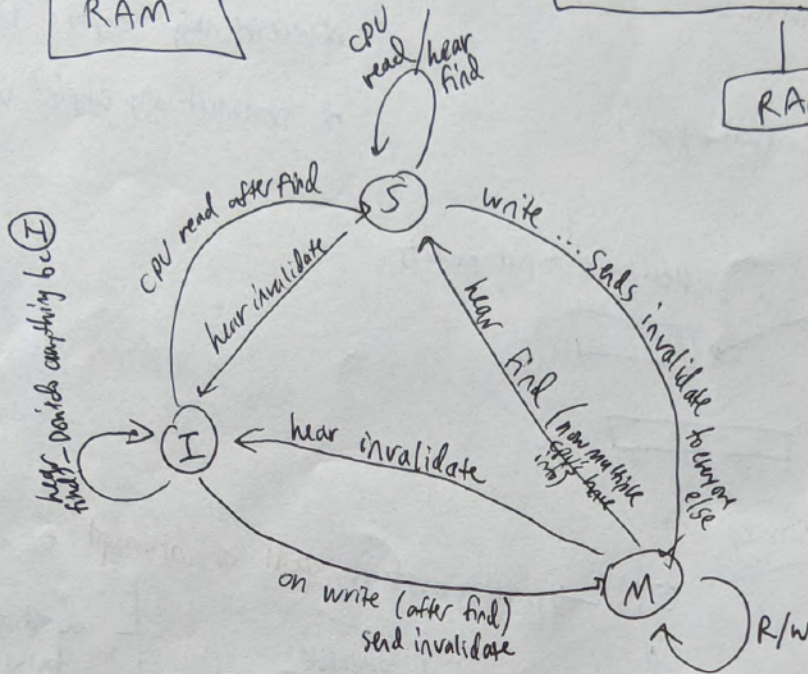
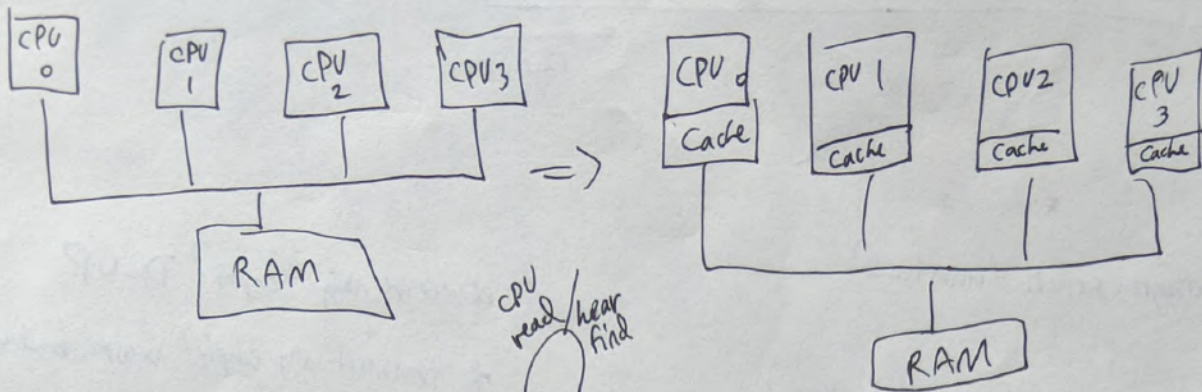
BISCUIT



HLL	C
- automatic mem mgmt	- no implicit code
- GC	- flexibility
- less use-after-free bugs	- control of mem/alloc/free
- poor performance	- hard to write
	- RCU expensive
GC makes this way easier	

SCALABLE LOCKS

6



messages —
invalidate
find

* At most one core in M
* Either one M, or I + S

MCS Locks

cpu requesters: **q node** {
 locked: T/F (have lock or nah)
 * next — next person to have the queue.
 }

1) **Lock** → empty node
 NULL

2) **Lock** → **My qnode** ^{!!} running!
 - locked = **T**
 - next = NULL

3) **Lock** → **My qnode**
 - locked = **T**
 - next → **Other qnode**
 - locked = **T**
 - next → NULL
 waiting, spinning on locked

RCU Usage

7

Reader - getval

prevent
interrupts
(disable)

```
rcu_read_lock()
p = rcu_dereference(itm) // read
    : ops
rcu_read_unlock()
```

* no
concurrent
writes!

* Read doesn't
lock.

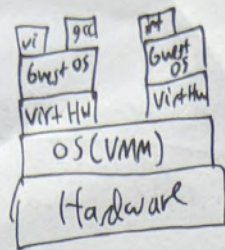
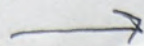
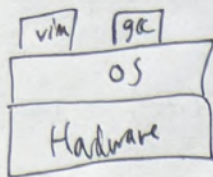
Writer - setval

```
spinlock
oldp = rcu_dereference_protected(itm) - mem barrier
```

```
ops, alloc new p, copy, set val
```

```
rcu_assign_pointer(itm, newp) - mem barrier
spinunlock
rcu_sync()
kfree(oldp)
```


VMS



Before: Host VA $\xrightarrow{\text{MMU}}$ Host PA

maps guest PAs \rightarrow Host PAs
 \rightarrow this + guest PT ~~page~~ helps VMM post to shadow pg table

Now: Guest VA $\xrightarrow[\text{Guest PT}]{} \text{Guest PA} \xrightarrow[\text{Map}]{\text{VMM}} \text{Host PA}$

* on reads, VMM will translate guest PAs to host PAs.

Actual: Guest VA $\xrightarrow{\text{Shadow pg table}} \text{Host PA}$

* on writes, guest OS is fine, but real MMU traps bc no mapping. VMM handles this and puts a mapping in shadow page table, or the actual table.

* 1 shadow page table / VM

memory tracing \rightarrow

* Guest PTs are read-only. Writes fault. Traps to VMM