

- OS = 1. Referee
2. Illusionist
3. Giver

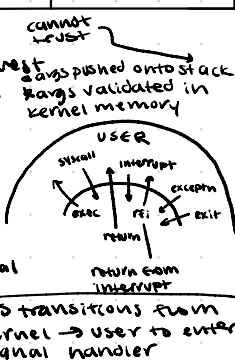
kernel	VS	User
MANAGES		MANAGES
- Process		User-level apps
- File System		Standard User-level code within restricted virtual addresses
- Memory		most shell commands
- Device		
- System calls		
- Security & Protection		

Hardware support

- Privileged instructions**
↳ unsafe instr ≠ user mode
- Memory isolation** → no access outside of process address space
- Interrupts** → kernel gain control from running process
- Safe transfers** → transfer from user mode ↔ kernel mode

kernel ↔ User

- syscalls** = user processes request certain kernel services
* **SYNCHRONOUS** event
* exec, read, write
- Traps** = software interrupts exceptions → internal
* **SYNCHRONOUS** event
* exceptions
- Hardware interrupts** = external
* **ASYNCH** event
* timer, I/O



To handle interrupt:

- Processor detects interrupt
- Suspend user prog, switch to kernel
- ID interrupt type, invoke handler
- Restore user program

Process = instance of running program

- Process control block (PCB): stores
 - process state
 - memory limits
 - PID
 - PTBR (page table ref process)
- executing program w/ restricted rights
- Illusions of: private CPU, private memory, private address per process
- resources allocated to each process
- isolate processes from all others

Process Management API

	Exit	terminate process
fork	copy current process as child & allocate PCB w/ duplicated memory content * new PID, same GID	Duplicates: - Address space - code/data - registers - stack * Same global FD as parent
exec	change program currently running on process → REPLACES it (replaces code/data) * GID = start of new program, reinit SP, FP	
wait	wait for process to finish	
Kill	send signal (interrupt-like) to another process	* can choose ANY SIGNAL to send
sigaction	set handlers for signals → can override default w/ signal handler execute signal handler when event has occurred	

File Descriptors

- POINTS TO OPEN file description in global table
- File description
 - File offset
 - file access mode
 - file status flags
 - ref to physical loc
 - # of times opened

- * dup to duplicate FD to point to same file
- * dup2 = replace existing FD
- Default FD Table
 - (1) STDIN - FILENO
 - (0) STDOUT - FILENO
 - (2) STDERR - FILENO
- * fork = duplicated descriptors point to SAME entry in global file table

two open calls & create 2 DIFFERENT entries in FD table and points to DIF file description in global table

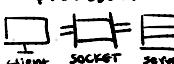
PIPES = ONE-WAY communication between two processes

pipefd[0] = read-only
pipefd[1] = write-only
Empty: blocks read
Full: blocks write

- * effectively closed after last write descriptor → reads until EOF
- * uses a kernel buffer in MEMOry

Sockets = 2 queues: 1 in each dir

* bidirectional communication between processes



- * can read/write to either and communicate between multiple processes on diff machines
- * FD obtained w/ socket()/bind()/listen()/accept()

HIGH / LOW API

LOW LEVEL ← less overhead

- direct use of syscall interface: open(), read(), write(), close()
- open() → returns file descriptor
- FD meaningful to kernel
- read() or write() causes syscall
- call read in a loop
- * write in while loop since not guaranteed to write all bytes out
- FD API = immediate I/O

HIGH LEVEL

- buffered access: fopen(), fread, fwrite, fclose
- opening returns FILE * struct
- File has memory for buffer and FD for the file
- fread / fwrite filters through buffers
- loop to get bytes w/ fread
- FILE * = buffered IO
↳ per-file user-level buffer
↳ operates on user space first

Threads = single execution sequence that can be separately schedulable

- INSIDE same process: not isolated from each other
 - share address space
 - share I/O state
- individual execution
 - individual stack & register state

Comparison category	Across processes	Across Threads
creation	fork()	pthread_create()
page table	distinct	same
registers/instruction ptr	distinct	distinct
stack	separate/inaccessible	separate but accessible
Heap/Static var	separate	shared
file descriptors	separate	shared
synchronization	wait(), waitpid()	pthread_join(), semaphore, locks
overhead	Higher	Lower
protection	Higher	Lower

- pthread_create = thread that runs start-routine
- pthread_exit = terminates calling thread & returns value - ptr to any successful join call
- pthread_join = suspends execution of calling thread until target thread terminates



Process = container where threads execute

* SIGSEGV = seg fault

* SIGKILL can NOT be overridden

- # Concurrency
- atomic operation** = always runs to completion / not interrupted
 ex: loads/stores of words
 - synchronization** = use atomic ops to ensure thread cooperation
 - mutual exclusion** = ensure only 1 thread does a thing @ a particular time
 - critical section** = piece of code that only ONE thread can execute at once
 - race condition** = when multiple threads/process access or modify shared resources at same time
 - lock() before entering critical section
 - unlock() when leaving
 - WAIT** if locked → all sync involves waiting
 - ① Busy wait = rely on atomic loads/stores → do not do consumes cycles
 - ② Enable/disable interrupts

Test & Set = simple lock that provides mutual exclusion w/ entry into the kernel

```
testset(&address) {
  result = M[address];
  M[address] = 1;
  return result;
}

int mylock = 0;
acquire(int *theLock) {
  while (testset(&mylock));
  // busy wait
}

release(int *theLock) {
  *theLock = 0;
}
```

Spinlock implementation = busy waiting
 so thread loops until lock is free
 so atomic operations

Futex = similar to testset in linux space

- ① **futex_wait** → check if condn still holds sleep until we hear futex_wake
 ↳ more busy waiting
 - ② **futex_wake** → wake up at most wait waiting threads
 ↳ more busy waiting
- lower level to maintain queues for sleeping threads
 while loop to just constantly keep trying to regain lock & avoid race condns
 for efficiency → want to minimize # of syscalls
- ```
int mylock = 0;
acquire(int *theLock) {
 while (testset(&mylock));
 futex_wait(&mylock, &futex_wait, 1);
}

release(int *theLock) {
 theLock = 0;
 futex_wake(&mylock, &futex_wake, 1);
}
```
- ↳ avoids syscalls when possible

**Semaphores** = lock w/ non-negative int value

- Down()** : wait for semaphore to be positive  
 ↳ decrement by 1 (like wait)
- Up()** : increment semaphore by 1 → wakes up waiting p
- Down on semaphore of 0** = putting thread to SLEEP
- Mutual Exclusion**: Value is 0 or 1 (like lock)
- Scheduling** : value can be >= 0

```
typedef struct semaphore {
 unsigned value;
 struct list waiters;
} semaphore;

semaphore sem;

sem_down(semaphore *s) {
 // disable interrupts
 while (semaphore_value == 0) {
 // sleep - schedule other threads
 // to do will be enable interrupts
 // decrement semaphore value by 1
 // enable interrupts
 }
}

sem_up(semaphore *s) {
 // enable interrupts
 // increment semaphore value by 1
 // enable interrupts
}
```

**locks** = prevent others from changing critical section

- ↳ can use internal semaphore, init to 1
- lock\_init(lock \*lock)**  
 initialize number of lock to null  
 sem\_init(&lock, 0, 1)
- lock\_acquire(lock \*lock)**  
 sem\_down(&lock)
- lock\_release(lock \*lock)**  
 sem\_up(&lock)

**Implementation:**  
 acquire():  
 disable interrupts  
 if (value == 0) {  
 put thread into queue  
 go to sleep  
 else: value == 1  
 enable interrupts  
 value--  
 return  
 }  
 value--  
 return

**Release():**  
 enable interrupts  
 if anyone in waitqueue  
 take thread off queue  
 place on ready queue  
 else: value == 0  
 enable interrupts  
 value++  
 return

- # conditional variables
- = queue of threads waiting for sm inside critical section  
**AVOID BUSY WAITING**
- purpose**: threads can access/modify shared data w/ certain condition
  - typedef struct cond\_var {  
 struct list waiters;  
 int cond\_var;  
 cond\_init(cond\_var \*cond)  
 initialize cond's list of waiters  
 cond\_wait(cond\_var \*cond, lock \*lock)  
 create semaphore  
 sem\_init(&waiter, sem, 0)  
 push waiter's semaphore into waiter queue  
 lock = release(lock)  
 sem\_down(&waiter, sem)  
 lock = acquire(lock)**
  - cond\_signal(cond\_var \*cond)**  
 if waiters queue not empty  
 pop waiter off waiter queue  
 execute sem\_up on waiter  
 only wakes up ONE waiting thread
  - cond\_broadcast(cond\_var \*cond)**  
 while waiters queue not empty  
 pop waiter off waiter queue  
 execute sem\_up on waiter  
 wake up ALL threads

**Monitor** = lock & zero + cond variables for managing concurrent access to shared data

- ↳ uses infinite synchronized buffers
- ↳ must hold lock when doing cond var ops

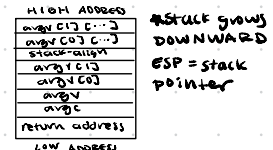
**HOARE**  
 if (is empty (b queue))  
 cond\_wait(&buff, &buff\_lock)  
 item = dequeue(b queue)

**MESA**  
 while (is empty (b queue))  
 cond\_wait(&buff, &buff\_lock)  
 item = dequeue(b queue)  
 while loop  
 ↳ more practical  
 ↳ repeatedly check for lock

**IF STATEMENT**  
 ↳ check once if code available  
**HOARE** = when thread signals cond var, the waiting thread immediately takes control of lock & starts executing  
 ✓ strong guarantee that cond is true when waiting thread wakes up  
 X dependent on OS scheduler

**calling convention**  
**CALLER** Before calling function:  
 ① save caller-saved GPRs (EAX, ECX, EDI) onto stack if access AFTER func call  
 ② push param onto stack in REVERSE order  
 Add padding BEFORE param to ensure 4 byte alignment  
 ③ CALL FUNCTION (esp is byte aligned w/ 4) and push return address & jump to func

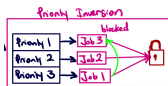
**CALLEE** new  
 ① Push ESP onto stack, set ESP = ESP  
 ② Allocate stack space for vars  
 ③ save callee saved GPRs (EBX, ESI, EDI) if used  
 ④ Perform func logic  
 EPILOGUE (before returning)  
 ① restore return value in EAX  
 ② restore callee saved GPRs (if any) from backup  
 ③ allocate local var  
 ④ restore caller's GPR from stack  
 ⑤ restore func call by popping caller return addr & jumping  
**CALLER** Epilogue after func returns  
 ① remove params from stack  
 ② restore caller-saved GPRs (if any)





**Scheduling** = deciding which threads to give access to **Starvation** = thread fails to make progress indefinitely

- ↓ **completion time** = waiting time + run time of process
  - ↳ \* time sensitive like I/O → wait time = total time on queue
- ↑ **throughput** = rate at which tasks completed (tasks/sec)
  - ↳ ↓ overhead (context switching)
- ↑ **fairness** = sharing resources in equitable manner



| Queue | FCFS/FIFO      | SRTF               | Approx.        | Linux CFS           | Round Robin             | EDF               | Priority        |
|-------|----------------|--------------------|----------------|---------------------|-------------------------|-------------------|-----------------|
| For   | CPU Throughput | Avg. Response Time | I/O Throughput | Turnover (CPU Time) | Turnover (CPU Air Time) | Meeting Deadlines | Task importance |

**Priority Inversion** = high priority task **BLOCKED** waiting on low priority task → medium runs

## POLICIES

### First Come, First Served (FCFS/FIFO)

- \* schedule in order of arrival → one program scheduled until done
- ✓ simple, good for throughput (↓ context switching)
- ✗ **CONVOY EFFECT** (short tasks stuck behind long) variable avg completion time

### Shortest Job First (SJF) NON-PREEMPTIVE

- \* schedule shortest task first
- ✓ optimal avg completion time (simultaneous arrival)
- ✗ subject to starvation, convoy effect have to know job duration

### Shortest Remaining Time First (SRTF) 2-PREEMPTIVE

- \* PREEMPT resource if task arrives & has shorter completion time
- ✓ always optimal avg completion time, <sup>NO</sup> convoy effect
- ✗ starvation, have to know job duration

### Generalized Processor Sharing

- infinitely fine grained context switching where to share processor among multiple threads

ideal service time = time that task i should receive

$$S_i(t_0, t_1) = E_i(t_0) * (t_1 - t_0)$$

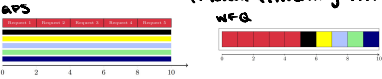
real service time = time that task i actually receives

lag = ideal - real service time

$$f_i(t) = \frac{w_i}{\sum w_j} \quad \text{proportional share of processor that thread i should receive at time } t$$

Set of active threads

weighted fair queuing = schedule threads in order of EARLIEST physical finishing time



### EEVDF = earliest eligible virtual deadline

$$\frac{dv}{dt} = \frac{1}{2w} \quad \begin{matrix} V_e = V_{d-1} + lag \\ V_d = V_e + \frac{r}{w} \end{matrix} \quad \text{used for rescheduling}$$

#### scheduling

- look @ threads w/ eligible requests ( $V_e < V_{d,e}$ )
  - serve request w/ earliest virtual deadline ( $V_d$ )
  - reschedule using the  $V_e, V_d$  equations if thread finished
- $t=0$  to  $t=1$ :  $\frac{dv}{dt} = \frac{1}{w_1}$  requests made in real time not arrival time
- $t=1$  to  $t=2$ :  $\frac{dv}{dt} = \frac{1}{w_2}$  when NEW thread becomes active, set  $V_e$  = curr virtual time
- use this to calculate virtual time ( $V_d$ )
- positive lag = EARLIER next time
- ✗ if lag:  $V_{e,i+1} = V_{e,i} + \frac{U^k}{w_i}$  time thread took to complete kth request

#### FOR RR

- At end of quantum X, add prev thread to ready queue
  - Add any new threads arriving at quantum X+1 to ready list
- move context switching = more cache misses
- size of CPU time

### Round Robin (RR) = run processes in looping order for fixed quanta

- \* after q expires: task preempted & added to end of ready queue
- LARGE q = ~FCFS vs SMALL q → lots of interleavings
- ✓ no starvation, w/ n process, q time quanta, <sup>max wait</sup> = (n-1)q
- ✗ lots of context switching, HIGH completion time

### Multi-Level Feedback Queue (MLFQ): multiple queues w/ diff priorities

- \* Each queue has own scheduling policy enters at highest priority
  - \* ensure CPU doesn't ho all
  - 10 priority > CPU priority
  - ✓ Approximates SRTF
  - ✗ Omeiga strategy: request I/O to stay on higher priority
- High → RR  
Low → FCFS
- use up resource = move down  
does not use all = move UP

### Lottery Scheduling: give each task some # of lottery tickets

- \* draw random ticket at each time slice → task holding ticket is granted resource
- ✓ No starvation, no priority inversion, approximate SJF w/ more tickets to shorter job
- ✗ unfair over short run fewer tickets to longer job

### Stride Scheduling = deterministic lottery scheduler

- \* Each task given some # of tickets  $n_i$
  - \* Stride =  $\frac{W}{n_i}$  larger share of tickets = smaller stride
  - on every time slice: choose task with lowest pass
  - init pass = min (existing tasks' pass vals) at start of task
  - ↳ when task chosen: pass += stride
  - ✓ No starvation
- smaller stride → task runs more often

### STRICT PRIORITY = schedules tasks w/ highest priority

- ✓ important runs first
  - ✗ starvation, priority inversion → task given same EFFECTIVE priority as higher one
- priority donation where lower

# Deadlock = cyclic waiting for resources

deadlock → starvation, starvation ≠ deadlock

## Requirements for Deadlock

- 1) Mutual Exclusion and bounded resources  
- one thread at a time use resources
- 2) Hold and wait  
- thread holding resource waits to acquire more
- 3) No preemption  
- resource released voluntarily
- 4) Circular Wait  
- set of waiting threads waiting on each other

## Deadlock Prevention

- 1) Provide sufficient resources, VM unlimited
- 2) Abort requests, acquire atomically
- 3) Fail if waiting too long, force give up
- 4) Order resources usage in same order

Necessary but NOT sufficient

## Banker's Algorithm

add to unfinished for each thread unfinished  
if (Request [Thread] ≤ Avail)  
remove from Unfinished  
Avail = Avail + Alloc

- 1) check if request → unsafe state
- 2) state max resource needs in advance
- 3) allow thread to continue if available resources - # requested ≥ max
- 4) release resources  
Available += allocated [thread]

regular algo checks if request [thread] ≤ available

- pretend request is granted → can we enter deadlock?
- ↳ if yes: deny/hang
- ↳ if unfinished NOT empty → deadlock is possible

## Memory + Addressing

### Numbers

- 1B = 8 bits
- 1KB = 2<sup>10</sup> B
- 1MB = 2<sup>20</sup> B
- 1GB = 2<sup>30</sup> B
- 1TB = 2<sup>40</sup> B

### Time

- 1ms = 10<sup>-3</sup> s
- 1μs = 10<sup>-6</sup> s
- 1ns = 10<sup>-9</sup> s
- 1ps = 10<sup>-12</sup> s

### Virtual memory:

- protection.
- translation
- efficiency

### Base and Bound



base register = starting address  
bound register = boundary of memory segment  
✓ simple, can easily relocate  
X internal AND external fragmentation  
Outflow interprocess sharing

W/ hardware relocation:

physical address = virtual address + base  
seg MMU = hardware as memory management unit

segmentation = place segments independently  
✓ READ-ONLY can be shared, use segment map  
✓ minimizes internal fragmentation  
X external fragmentation still problem

## Paging (single level)

| PHYSICAL ADDR   | offset |
|-----------------|--------|
| Physical page # |        |

| VIRTUAL ADDR   | offset |
|----------------|--------|
| Virtual page # |        |

# of bits for # of pages in VA space  
20 | 12  
2<sup>10</sup> pages | 2<sup>12</sup> = 2<sup>10</sup> · 2<sup>12</sup> = 4KB

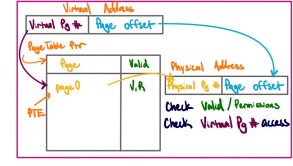
**PTBR** = page table base register = starting addy of pg table  
↳ physical address

- # virtual address bits = (# VPN bits) + (# offset bits)
- # physical address bits = (# PPN bits) + (# offset bits)
- # virtual pages = 2<sup>n</sup> (# VPN bits)
- # physical pages = 2<sup>n</sup> (# PPN bits)
- Pg size (bytes) = 2<sup>n</sup> (# offset bits)

virtual mem size (bytes) = (# virtual pgs) × (pg size B)  
physical mem size (B) = (# physical pgs) × (pg size B)

# Pg table entries = # virtual pgs  
# PTE bits = (# PPN bits) + (# metadata bits)

pte size (B) = # (pte bits) / 8  
Page table size (B) = # pte × pte size (B)



LARGER pages = fewer metadata  
SMALLER pages = minimize internal fragmentation

## Pg replacement POLICIES

- 1) FIFO: first oldest page  
Can: evict rarely used pages
- 2) RANDOM: choose random page for replacement  
Can: unpredictable
- 3) LRU: replace page not used for longest time, optimal  
Can: don't know future
- 4) LRU: replace page not used for longest time  
Can: too much overhead

NOT STACK

STACK

STACK

STACK

Belady anomaly: increasing cache can INCREASE miss rate → applies to

## NON STACK COPY ON WRITE

- 1) Mark all pgs read only
- 2) copy pg table to new process
- 3) once either attempt write  
↳ PAGE FAULT
- 4) Page fault handler → make physical copy → install new pg in pg table of triggered process

TRANSLATION LOOKASIDE BUFFER = cache for virtual-physical memory limitations

| TAG | DATA         |
|-----|--------------|
| VPN | PPN metadata |

$$AMAT = (\text{Hit rate} \times \text{hit time}) + (\text{miss rate} \times \text{miss time})$$

improve context switching

## Page walk

1st PT: Addr: Base + Index<sub>L1</sub> × Entry Size  
2nd PT: Addr: PPN<sub>L1</sub> (look in phys mem) || Index<sub>L2</sub> × Entry Size

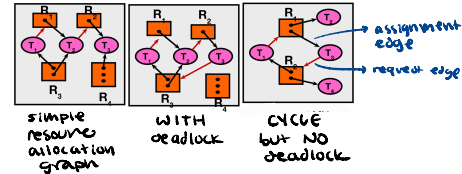
check metadata: convert to binary and check each bit  
Phys addr: PPN<sub>L2</sub> (look in phys mem) || offset

Safe space: can prevent by delaying acquisition  
unsafe space: can unavoidably lead to deadlock, with certain acquisition  
Deadlocked state: exists a deadlock

## Handle deadlock

- 1) Denial (osrith-pretend it doesn't exist)
- 2) Prevention - write code that doesn't deadlock
- 3) Recovery - let deadlock happen & recover afterwards
- 4) Avoidance - dynamically delay resource reqs so no deadlock  
cannot check for deadlock → check for unsafe state

\* Draw out current, available, needed tables



simple resource allocation graph

WITH deadlock

CYCLE but NO deadlock

Demand Paging = keep only active pgs in memory  
place others on disk, PTE = invalid

WORKING SET = subset of addy space used during execution

RESIDENT SET = subset of addy space held in memory

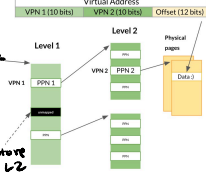
THRASHING = memory too small for working set → lots of replacements & context switches  
N:RSS > PHYS - MEM

PAGE FAULT = attempt to access virtual NOT in physical  
invalid PTE

- 1) MMU traps to OS w/ page fault
- 2) Find & replace pg w/ pg from disk w/ policy
- 3) reset pg table, restart instruction

## MULTI level page tables

- sparse addy space = large pg table w/ lots of unused entries
- TREE of pg tables where each level maps to other pg table
- BOTTOM page table = translation to physical address



don't need to store unmapped L2

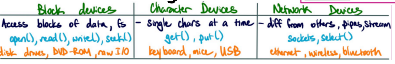
I/O

Access patterns:

- 1 Character dev = character stream  
NOT addressable
- 2 Block dev = fixed size blocks  
YES addressable
- 3 Network dev = separate interface for networking

Access Timing

- 1 Synchronous/blocking = wait until I/O request is fulfilled
- 2 Non-blocking = no wait, quickly return
- 3 Asynchronous = allow other processing w/o waiting



Device Drivers = connect I/O hardware

Standard interface OS high-level abstractions  
TOP half = start I/O operations  
BOTTOM half = service interrupts

Bus = wires for comm/connecting in devices + protocols for data transfer  
cheap + easy to connect X one at a time + arbitration logic for who gets control

PCIe I/O bus interface for high-speed, not parallel, fast serial channels  
Network cards, SSDs, graphics X modems, keyboards, USB

Device controllers = how processors talk to I/O devices  
wants to interact w/ controller

- 1 Programmed I/O (PIO) - processor involved in every byte transfer  
A Port-mapped I/O (PMIO) - separate bus  
Special privileged CPU instructions (in, out)  
Distinct memory address space from phys mem  
Good for limited physical space since its pair of physical address space

- 2 Memory-mapped I/O (MMIO) (load/store)  
Standard load/store + responsibility on controller  
No physical address space + shared w/ physical memory

- 3 Direct mapped access (DMA) device can write directly to main memory w/o CPU intervention  
CPU sets up DMA controller  
DMA request  
DMA transfer  
DMA interrupt  
DMA transfer  
DMA interrupt

Storage Devices

- 1 Hard Disk Drive (HDD) - magnetic disk storage device  
Persistent memory + stays even computer is off  
Data stored in fixed size blocks = sectors  
Sequential access + random access

Request time = queuing + controller + seek + rotational + transfer  
Software queue in device driver  
Hardware controller  
Position head/wind over track  
Rotational latency  
Transfer time

avg seek = 8 ms  
rot = 720 RPM  
data transfer = 50 MB/s  
sectors = 4 kB  
Example throughput  
seek = 8 ms  
rot time = 1/720 \* 60 = 8.33 ms  
transfer time = 4096 bytes / 50 MB/s = 81.92 μs  
throughput = 1 / (8 ms + 8.33 ms + 81.92 μs) = 12.1 MB/s

data read w/ head, platter = 2 sides  
cheaper than flash X slower than flash (non-volatile)

- 2 Flash memory = electrical circuits for persistent storage  
can erase fixed # of times  
request time = queuing + controller + transfer  
writing to cell requires ERASING it  
faster read/write w/ pages

- 3 Flash translation layer map logical flash pgs -> phys pages  
indirection + copy on write to avoid wear out  
wear leveling to spread out writes  
write to new location -> update mapping in FTL -> erase old pg in background

File System =

File system = transforms block interface of disks into files/dir  
File system = maintained by OS, file has metadata (size, owner, access)  
Directory = list of mappings from human-readable file names to specific underlying files  
Hard links = map diff names to same file  
In-entries [leaf]  
Soft links = directory entries that map name to another name  
In-entries [leaf]  
Name = file name  
File System Design

- 1 Fast File System = multi-level tree structure  
File entry -> set of indirect array  
inode = file w/ metadata, ptrs to data  
No file name/data in inode  
FO entry -> open file descriptor (application)  
ptr = block numbers  
Direct ptr = point directly to block w/ file data  
Indirect ptr = point to array of direct ptrs  
Double indirect ptr = point to array of indirect ptrs  
Number = index of inode in inode array (used as file #)  
max disk size = 2^32 blocks \* block size = 2^32 \* 4 = 4 TB  
max file size = block size \* # of direct ptrs (including from indirect)  
= 2^32 \* (12 + 512 + 512^2 + 512^3) = 16 EB

File Allocation Table

File = collection of disk blocks  
File -> linked list one on one w/ blocks  
File size encoded sequentially, no external, big file  
in 32 bits -> size limited to 4 GB  
Entry = 0 = free  
quick format = make all blocks as free  
Full format = quick format + zero out data blocks  
Entry -> 0 = index of next cluster  
Special (end of file, free cluster, bad cluster)

NTFS

Master file table, entries = file metadata + data  
variable size elements w/ variable length (not fixed size blocks like FAT)  
Big -> use extents super block: use ptrs to lists of extents  
supports journaling

Buffer Cache

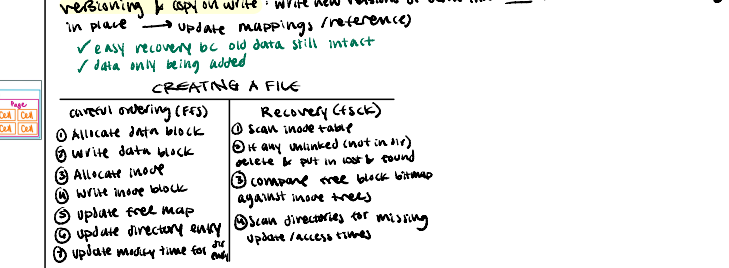
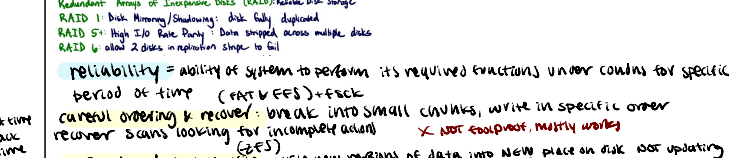
memory to cache disk blocks/name translations  
implemented in OS software w/ CPU replacement policy  
Read ahead prefetching = fetch sequential blocks early  
Delayed writes = flush to disk -> either on failure or periodically  
Reliability = probability that system can accept & process requests  
"Nines" of probability: 99.97 = 3 nines  
Durability = ability of system to recover data despite faults / fault tolerance  
Small defects = Reed-Solomon error correcting codes (ex. single disk failure)  
Short-term data preservation -> battery backed RAM (NVRAM) (ex. or better cache)  
Long-term data preservation -> redundant array of inexpensive disks (RAID)

Reliability

Reliability = ability of system to perform its required functions under conditions for specific period of time (FAT + FFS) -> check current ordering & recover: break into small chunks, write in specific order  
reducer scans looking for incomplete data  
versioning & copy on write -> write new regions of data into new place on disk not updating in place -> update mappings (reference)  
easy recovery bc old data still intact  
data only being added

CREATING A FILE

- 1 Allocate data block
- 2 Write data block
- 3 Allocate inode
- 4 Write inode block
- 5 Update free map
- 6 Update directory entry
- 7 Update mtime for dir
- 1 Scan inode table
- 2 Scan unlinked (not in dir) delete & put in lost & found
- 3 Compare empty block bitmap against inode trees
- 4 Scan directories for missing update/access times



# Transactions = atomic reversion of operations

**consistent state** → **consistent state** updated disk w/ transactions  
 \* if any fail → roll back  
 ↳ otherwise, commit!  
 Atomicity  
 Consistency (1 valid state to another)  
 Isolation (concurrent is ok)  
 Durability (once committed won't be erased)

**Methods:** (BOTH USE LOG = BUFFER ON DISK)  
 \* **Journaling Filesystem** → **multistage in SUBSYSTEM before committing**

- Preparation phase = append all updates to log
- Commit phase = log commit record
- Write back = async apply committed update
- garbage collect any **completed transaction** in background

\* write to log instead of disk directly  
 \* once **WHOLE** transaction written to log - disk applies changes  
 \* remove transaction from log ✓ can roll back

\* if crash when writing → incomplete, NOT applied  
 \* if crash when **applying** → re-apply transaction after recovery

- log structured File system (LSFS):** log is storage  
 write everything sequentially \* all data in log  
 ✓ writing random sectors better than normal

# Distributed Decision Making

Generalis **Paradox** = impossible to achieve simultaneous acknowledgement over **UNRELIABLE** network

**Two Phase Commit** = decide if all processes commit/abort transaction eventually → 1 coordinator, rest participants

- Coordinator asks all processes to vote **VOTE-REQ**
- Participants vote **VOTE-COMMIT** / **VOTE-ABORT**, log
- If ALL vote commit → **GLOBAL-COMMIT** otherwise → **GLOBAL-ABORT**, log
- Participant commit or abort on receive, log  
 respond w/ ACK

**Failure:**  
 \* any participant error → coord votes **ABORT**  
 \* if all voted commit → wait on coord to recover

when **coordinator logs** → becomes **SOURCE OF TRUTH**  
 persistent decision

\* DOES NOT solve **General Paradox** → still **Uncertainty**

# Internet = allow apps to function on all networks

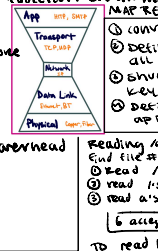
**Packet** = basic unit of communication  
 ↳ chunk of data, some metadata, source destination IP

**Transport Layer** on top of packets

**UDP** = best effort delivery

**TCP** = reliable in-order, but **MORE** overhead

**end-to-end:** functions should be implemented at endpoints of communication system = **NOT** in middle  
 \* implement or lower layer only as **enhancement**



# Distributed Systems = coordinating multiple computers

**Scalability:** add resources to system to support more work

**Transparency:** mask complexity behind simple interface

**Protocols:** agreement on how to communicate (syntax/semantics)



# Remote Procedure Call = execute code on remote machine in simple way

**Serialization:** express object in byte sequence  
 big/little endian: first bit in array = most/least sig bit  
**Marshalling:** converting values to canonical form (serializing)  
**Unmarshalling:** converting user-visible name to network endpoint  
 ↳ dynamic binding allows server flexibility

**STUB** = glue on client/server

- Client stub = marshalling/serializing args  
 unmarshalling/deserializing return value
- Server stub = unmarshalling/deserializing args  
 marshalling/serializing return values

# Distributed File Systems

**Virtual Filesystem Switch (VFS)** = interface between kernel's view of files & underlying file system (-ATFS, NFS...)

\* **SAME** system call interface API for diff types of file system

\* **VFS** like translator for file systems

- Object types**
  - superblock object = specific mounted filesystem
  - inode object = specific file
  - directory object = directory entry
  - File object = open file associated w/ process

\* distinguish between **remote** & **local** files

**Network File System (NFS)**

Translates read/write VFS calls → **RPCs** for remote file system

\* **STATELESS:** RPC has all info to ID a file

readable (number, position) vs read (offset)

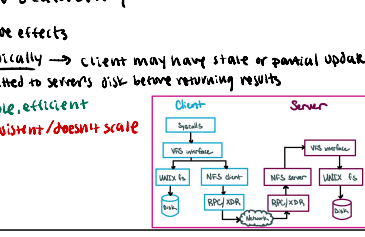
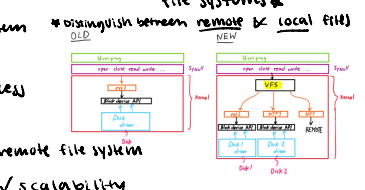
\* **IDEMPOTENT:** re-read / re-write blocks w/o side effects

\* **WEAK consistency:** client polls server periodically → client may have state or partial updates

\* **write through caching:** modified data committed to server's disk before returning results

\* **wait sequential ordering** ✓ **simple, portable, efficient**

✗ **sometimes inconsistent / doesn't scale**

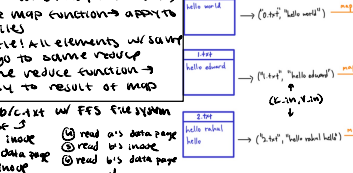


# MapReduce & Spark

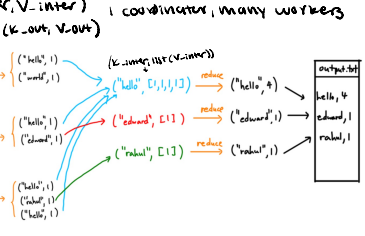
**MR** = model that distributes tasks over large clusters

**map:**  $(K-in, V-in) \rightarrow list(K-inter, V-inter)$

**reduce:**  $(K-inter, list(V-inter)) \rightarrow list(K-out, V-out)$



1 coordinator, many workers



**Reading w/ b/c/txt w/ FFS file system**

- Read 1/3 inode
- Read 1/3 data page
- Read 1/3 inode
- Read 1/3 data page
- Read 1/3 inode
- Read 1/3 data page

**6 accesses** directly has inode of c/txt  
 To read 165 Kib give c/txt files  
 165 Kib = 165 Kib  
 Blocksize = 4 Kib → round(41.25) = 42  
 165 / 42 = 3.9285714285714284  
 3.9285714285714284 \* 42 = 165 Kib

**Fault Tolerant** = large cluster w a lot of failures

- Coordinator fails → **ABORT**
- Reassign in progress or completed map tasks if mapper fails
- Reassign in progress reduce tasks if reducer fails

↳ check or complete map

**Apache Spark** = how to reuse intermediate results  
 ↳ cache instead of recomputing

\* **Resilient Distributed Dataset (RDDs)** = immutable collections of objects in memory

- Transformations: apply func on RDD → make another RDD  
 ↳ lazily evaluated w/ **iterator** → replay transformation to regenerate

- Actions → return value to program / write data to disk  
 ↳ **lineage**: logged sequence of transformations to create RDD

\* **Failures** Node fails → lost RDD partition → recompute w/ lineage