

Katrina Sharenin

OS Basics CPL(0)-greatest \leftrightarrow CPL(3)-lowest \star i++ \rightarrow ret original value t+ti \rightarrow ret new value \star cond? -1:1; if true \rightarrow 1

OS: layer of SW that interfaces b/w HW resources & applications running on machine, implement VA

Referee: protection, isolation, sharing + Evaluation: performance, reliability, security, portability

Process: CPU + Mem (Stack, data, Heap, code) + registers (SP, PC) + IO \rightarrow illusion of priv. CPL

\hookrightarrow PCB: PC, esp, mem limit, reg, PTBR, PID, VID, FDT, state \star Maskable = can be disabled/ignored

Kernel: lowest level of OS, full access: Applications (untrusted) \leftrightarrow Kernel \leftrightarrow HW (untrusted) \star PT base ptr \rightarrow halt processor

Dual mode: 2 bit exec. **User mode:** processor checks each instr, limited set (cannot change: priv. level, addy space, disable interrupts, I/O)

Kernel mode: OS execute w/ protection off, execute any instr. Need safe transfer & control again

Syscall: user \rightarrow kernel at specific locations by OS	Exceptions: user mode attempts priv operation User req. OS Service, transfer to kernel ① Trap instr (int 0x80) \rightarrow EV 18, exc=sys#, user ② Sys. Dispatch table ③ Ret vals in registers *iret: esp++, intr. not disabled	Interrupts: asynch signal to processor, external event needs attention Stop process \rightarrow enter kernel at interrupt handler \star Resetting timer is priv. op. ① Processor detect ② Suspend user, switch to kernel \star Save before esp = base stack Index \rightarrow INTV IRQ \rightarrow APIC ③ Interrupt type \rightarrow invoke handler ④ Restore user prog. \star disable intr. when running ⑤ Atomic Switch \star INTV in kernel \star Use interrupt stack

Process API: 1 parent, 0+ child, all start from 'init' process

\star syscalls, doesn't wait for return

exit(): terminate process \star exec("ls", args): replace code/data seg.

fork(): copy addy space (all), set esp to start of prog/reinit sp/esp

code/data seg, reg. (PC+SP), stack, FDT pointers, new PID, READY

wait(): suspend until one child exits

*arbitrary \star cpid=fork() == 0 child > parent

sigaction(): for sig. User can override signaler/default action

kill(pid): SIGKILL (interrupt 11)

\star O.1.2 (stdin, stdout, stderr)

Interprocess: Pipe: one-way, read-only and (ppipefd[0]) write only (ppipefd[1]),

Blocks if W called on full, blocks R if empty pipe, close after last W fd closed

* queue \star ret EOF \star after last R closes \rightarrow SIGPIPE signals, W fails w/ EPIPE err

Socket: 2 queues, K/W either end, fd obtained via socket/bind \star close per process

Threads: single execution sequence, separately schedulable task

* virtualizes processor, int processors, sequential code run concurrently

Kernel thread: 1 thr = 1 TCB, indiv. schedulable, switch 2 switch

pthread_create(*thread, routine, arg)

pthread_exit(val_ptr): terminate & make val_ptr available to any successful join \rightarrow signal! _exit(NULL)

pthread_yield(): yield CPU to other threads, not guaranteed

pthread_join(thread, *val_ptr): suspend until target terminates \star join(thr, NULL)

Synchronization: Multiprocessing (2+CPUs) Mult. programming (2+ processes) Multithreading (2+ thr. per process) Atomic: 1 or 0

lock_init(&mylock) or mylock = MUTEX_INIT

acquire(&mylock); wait until free thru mark | release(8 mylock) only by holder

test&set implementation:

acquire(int* lock) {

 while (test&set(lock)) {};

} \star OR !CAS(lock, 0, 1)

release(int* lock) {

 *lock = 0; \star avoid spin, like int guard

}

futex lock implementation

acq(lock){ while (test&set(lock)) {

 futex(lock, WAIT, 1); }

release(lock){ lock = 0; }

else { int r = lock; }

 lock = r; }

 futex(lock, WAKE, 1); }

 *avoid spin, like int guard

 lock = r; }

 lock = r; }

PintOS

```
struct list_elem {
    struct list_elem *prev;
    struct list_elem *next;
} list;
```

```
struct list {
    struct list_elem *head;
    struct list_elem *tail;
} list;
```

iterate & remove
or e = list.next(e)

for (e = list.begin(&list); e != list.end(&list); e = list.remove(e)) {
 word_count_t cur = list_entry(e, word_count_t, elem);
}

search for this elem embed struct type field name

Base & Bound: Context switch == change base & bound, each has VA \rightarrow PHYS_BASE
Base + offset \leq Base + Bound checked

processes: * child can exit & state freed by parent in wait(), parent doesn't need to exit * PCB malloc by kernel memory: * Kernel: 0x000 + maps to same PA, kernel always mapped, panic if deref NUL, all user: 0 \rightarrow 0xc00

* User & kernel stacks separate * Syscall can use "kernel only" PTE w/o updating PT

Threads: shared scheduler for user & kernel * Normal thr: enter kernel & save info on parent kernel stack before scheduling

Kernel only: exclusively on kernel stack

Sys call: thr blocks any time / depth \rightarrow save kernel stack on wait queue, restore state from another kernel stack linked on ready queue

Pipe cont.

```
int pipefd[2], *init_fds  
int status = pipe(pipefd);  
// for i  
if (pid == 0) { // child  
    close(pipefd[1]); // close write  
    read(pipefd[0], buf, sizeof(buf));  
    close(pipefd[0]); // close read  
} else { // parent  
    close(pipefd[0]); // close read  
    write(pipefd[1], msg);  
    close(pipefd[1]); // close write  
}
```

Sockets cont.

* all syscalls

- ① Create server socket via `socket()`
- ② Bind to specific addr via `bind()`
- ③ Listen for cons via `listen()`

* Unique conn: (srcIP, dstIP, SrcPort, dstPort, TOS)
void *serveclient(argv) {
 while ((n = read(sock, buf, bufsize)) > 0) {
 ...
 if (write(sock, buf, n) == -1) {
 close(sock)
 ... exit
 }
 }
 close(sock); thr-exit; // close up

int main(argc, argv) {

```
    server_socket = socket(family, type, protocol)  
    if (s.s == -1) ret 1;  
    if (bind(server_socket, addr, addrlen) == -1) ret 1  
    if (listen(server_socket, 1) == -1) ret 1  
    while (1) {  
        int conn_socket = accept(server_socket, NULL, NULL)  
        // etc.  
    }  
}
```

File OS cont. * guarantee they will read/write some specified amount, unlike low-level

fopen(*path, *mode): ret *FILE, open file assoc. file

fclose(FILE *stream): flush stream & close FD, 0 on success or EOF on err

fread(*ptr, size, n, *stream): read n items of data, each size bytes long, from stream, storing to ptr, $\text{ret} \#$ items

fwrite(*ptr, size, n, *stream): write " " " " $\text{ret} \#$ items written

fflush(FILE *stream): forces all data in user space buffer to stream

fprintf(FILE *stream, *format): print stream according to format

fscanf(FILE *stream, const char *format): scans stream according to format

ssize_t read(int fd, void *buf, size_t count): attempts to read up to count bytes from FD into buffer, $\text{ret} \#$ bytes

MT2 Cheat Sheet

Scheduling

: process of deciding which thread given access to resources for moments (CPU, disk access, etc)

① Min (response time / latency), user-priority time to do task ② Max (throughput): tasks/time ③ Min overhead (+ to switch) ④ Predictability

Wait time = total time in ready queue Response time = time until process first scheduled Completion E = W + runtime, total t

Policies

* Convoy effect: any non-preemptive, long blocks short

First Come First Serve (FCFS)

: FIFO arrival

Pros: simple, good throughput, minimizes overhead of context switch

Cons: avg completion time highly variable, convoy, job titanic factor

Shortest Job First (SJF): shortest task first → avoid convoy

Shortest Remaining Time First (SRTF): preemptive SJF

Pros: optimal for min avg completion time for nonpre & prem.

Cons: psychic/impossible, starvation possible (SOP)

O(n) scheduler: at each switch Scan full list of processes in ready q → compute p. score → pick best

Cons: bad scalability w/ more threads (multiple CPUs)

O(1) scheduler: like MLFQ, choose next in const. time.

Ideal time/interval: 0-99 (0 highest), Vavg: 100-139 nice

Per p. intv 2 ready qs; active (not past Δt) expiring (past Δt); job finish calc. pri.

Cons: very complex to reason & tune

* Any policy that favors fixed priority for scheduling causes starvation

Stride scheduling: deterministic proportional fair sharing; stride = #W/Ni → # tickets → more = smaller stride, run more

① look for lowest priority counter ② run ③ add stride; Cons: long delays until job is rescheduled if big stride vs small

Earliest Eligible Virtual Deadline First (EEVDF): schedule tasks for Δt proportional to weight, Minimize lag

Each task has eligible time to run & deadline to finish by

* Generalized processor sharing (GPS): fluid / weighted fair grain switching $f = \text{Capacity}/(w_1 + w_2 + w_3) \approx \min(r_i, f_{\text{avg}})$

Eligible time = $t - \sum_{j \neq i} w_j \text{active} \cdot dt$ can be scheduled ⇒ ① At each t only consider eligible packets ② Pick with earliest deadline

Lag = ideal service time - real service time ⇒ ③ Lag = underallocated → move & schedule earlier vs. ④ punish

$f_{\text{avg}}(t) = w_i / \sum w_i \rightarrow S_i(t_0, t_1) \text{ideal} = f_{\text{avg}}(t_0) \cdot (t_1 - t_0)$ J assume f const in t → E,

⑤ $S(0, 1) = 0.5 * (1 - 0) = 0.5$ s, transmits 0.5 bits ($\frac{1}{2}$ packet) in $t = 0$ to 1, full access to CPU for 0.5s

Virtual time $V(t)$: rand we are cur. at real time t, rate; inversely proportional to $\sum w_i$ active tasks

Moves slower w/ more active tasks since time / round incr. Round duration = $\sum w_i$ active thr.

$V(t_1) = S_0 \cdot t_1 / \sum w_j \text{active} \cdot dt \Rightarrow$ If # active thr. const $t_1 - t \rightarrow V(t) - V(t_1) = (\sum w_j \text{active})(t - t_1)$

Virtual eligible deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Virtual deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Algorithm ① every quantum q ② get $V_e \leq V_d$ ③ select w/ earliest virtual deadline V_d

④ Compute next based on lag ⑤ New $V_e^{k+1} = V_d^{k+1} - (\text{lag})(w_i)$ ⑥ New $V_d^{k+1} = V_e^{k+1} + R/w_i$

* Next eligible t gets impacted by lag → ⑦ lag pulls closer eligible time vs. ⑧ postpares it, ← eligible

⑨ C1 C2

virtual t

physical t

VE 0 0.5 1 1.5 2

physical t 0 1 2 3 4 5 6 7

* C1: $w_1 = 2, r_1 = 2$ quanta

* C2: $w_2 = 2, r_2 = 1$ quanta

q = 1 second

* no lag → $V_e^{k+1} = V_d^k$

Weighted Fair Queuing (WFQ): approx GPS, select packet that finishes first in GPS assuming no future arrivals

* Every client receives at most q delay compared to GPS * CONSIDER finish time dynamically changes/recompute since

Solve with virtual time * Virtual finish time doesn't change w/ new arrivals, always (len packet / w_i)

* System virtual time: index of current round in weighted bit by bit RR

WFQ: F_i^{k+1} (round by which packet k+1 send) = max($V(a_i, k+1)$, F_i^k) + L_i^{k+1}/w_i

F_i^k : V finish time of packet k of client i

L_i^k : len of packet k of client i

"Takes up exactly requested service time" = no lag

Round Robin (RR): each intv with time quantum g

* Large q → higher thrnt → resembles FCFS

* Small q → lots of interleaves

Pros: fair, w/ n process & q, max waiting time $C(n-1)g$, no starvation

Cons: small q increases completion time / suffuses from low throughput

Multi-level feedback Q (MLFQ): queues for ready jobs per p. level

* Approximates SRTF; placed at highest, bump it post Δt given strvnt

Cons: games with I/O req → no matter how many times given up, reduces

Subject to starvation → after time period S, move all jobs to topmost

Proportional Fair sharing: share CPU proportionally to weights

Job share according to priority, run less but no starvation

* Y/n # users ex max-min, user reads less is share ex proportional w/ p.

Lottery scheduling: type of CFS; give ticks, prob win each Δt

* on avg. each job's CPU time proportional to # tickets it has

Cons: could choose long jobs/law prior, unfair for less jobs

Stride scheduling: deterministic proportional fair sharing; stride = #W/Ni → # ticks → more = smaller stride, run more

① look for lowest priority counter ② run ③ add stride; Cons: long delays until job is rescheduled if big stride vs small

Earliest Eligible Virtual Deadline First (EEVDF): schedule tasks for Δt proportional to weight, Minimize lag

Each task has eligible time to run & deadline to finish by

* Generalized processor sharing (GPS): fluid / weighted fair grain switching $f = \text{Capacity}/(w_1 + w_2 + w_3) \approx \min(r_i, f_{\text{avg}})$

Eligible time = $t - \sum_{j \neq i} w_j \text{active} \cdot dt$ can be scheduled ⇒ ① At each t only consider eligible packets ② Pick with earliest deadline

Lag = ideal service time - real service time ⇒ ③ Lag = underallocated → move & schedule earlier vs. ④ punish

$f_{\text{avg}}(t) = w_i / \sum w_i \rightarrow S_i(t_0, t_1) \text{ideal} = f_{\text{avg}}(t_0) \cdot (t_1 - t_0)$ J assume f const in t → E,

⑤ $S(0, 1) = 0.5 * (1 - 0) = 0.5$ s, transmits 0.5 bits ($\frac{1}{2}$ packet) in $t = 0$ to 1, full access to CPU for 0.5s

Virtual time $V(t)$: rand we are cur. at real time t, rate; inversely proportional to $\sum w_i$ active tasks

Moves slower w/ more active tasks since time / round incr. Round duration = $\sum w_i$ active thr.

$V(t_1) = S_0 \cdot t_1 / \sum w_j \text{active} \cdot dt \Rightarrow$ If # active thr. const $t_1 - t \rightarrow V(t) - V(t_1) = (\sum w_j \text{active})(t - t_1)$

Virtual eligible deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Virtual deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Algorithm ① every quantum q ② get $V_e \leq V_d$ ③ select w/ earliest virtual deadline V_d

④ Compute next based on lag ⑤ New $V_e^{k+1} = V_d^{k+1} - (\text{lag})(w_i)$ ⑥ New $V_d^{k+1} = V_e^{k+1} + R/w_i$

* Next eligible t gets impacted by lag → ⑦ lag pulls closer eligible time vs. ⑧ postpares it, ← eligible

⑨ C1 C2

virtual t

physical t

VE 0 0.5 1 1.5 2

physical t 0 1 2 3 4 5 6 7

* C1: $w_1 = 2, r_1 = 2$ quanta

* C2: $w_2 = 2, r_2 = 1$ quanta

q = 1 second

* no lag → $V_e^{k+1} = V_d^k$

Weighted Fair Queuing (WFQ): approx GPS, select packet that finishes first in GPS assuming no future arrivals

* Every client receives at most q delay compared to GPS * CONSIDER finish time dynamically changes/recompute since

Solve with virtual time * Virtual finish time doesn't change w/ new arrivals, always (len packet / w_i)

* System virtual time: index of current round in weighted bit by bit RR

WFQ: F_i^{k+1} (round by which packet k+1 send) = max($V(a_i, k+1)$, F_i^k) + L_i^{k+1}/w_i

F_i^k : V finish time of packet k of client i

L_i^k : len of packet k of client i

"Takes up exactly requested service time" = no lag

Round Robin (RR): each intv with time quantum g

* Large q → higher thrnt → resembles FCFS

* Small q → lots of interleaves

Pros: fair, w/ n process & q, max waiting time $C(n-1)g$, no starvation

Cons: small q increases completion time / suffuses from low throughput

Multi-level feedback Q (MLFQ): queues for ready jobs per p. level

* Approximates SRTF; placed at highest, bump it post Δt given strvnt

Cons: games with I/O req → no matter how many times given up, reduces

Subject to starvation → after time period S, move all jobs to topmost

Proportional Fair sharing: share CPU proportionally to weights

Job share according to priority, run less but no starvation

* Y/n # users ex max-min, user reads less is share ex proportional w/ p.

Lottery scheduling: type of CFS; give ticks, prob win each Δt

* on avg. each job's CPU time proportional to # tickets it has

Cons: could choose long jobs/law prior, unfair for less jobs

Stride scheduling: deterministic proportional fair sharing; stride = #W/Ni → # ticks → more = smaller stride, run more

① look for lowest priority counter ② run ③ add stride; Cons: long delays until job is rescheduled if big stride vs small

Earliest Eligible Virtual Deadline First (EEVDF): schedule tasks for Δt proportional to weight, Minimize lag

Each task has eligible time to run & deadline to finish by

* Generalized processor sharing (GPS): fluid / weighted fair grain switching $f = \text{Capacity}/(w_1 + w_2 + w_3) \approx \min(r_i, f_{\text{avg}})$

Eligible time = $t - \sum_{j \neq i} w_j \text{active} \cdot dt$ can be scheduled ⇒ ① At each t only consider eligible packets ② Pick with earliest deadline

Lag = ideal service time - real service time ⇒ ③ Lag = underallocated → move & schedule earlier vs. ④ punish

$f_{\text{avg}}(t) = w_i / \sum w_i \rightarrow S_i(t_0, t_1) \text{ideal} = f_{\text{avg}}(t_0) \cdot (t_1 - t_0)$ J assume f const in t → E,

⑤ $S(0, 1) = 0.5 * (1 - 0) = 0.5$ s, transmits 0.5 bits ($\frac{1}{2}$ packet) in $t = 0$ to 1, full access to CPU for 0.5s

Virtual time $V(t)$: rand we are cur. at real time t, rate; inversely proportional to $\sum w_i$ active tasks

Moves slower w/ more active tasks since time / round incr. Round duration = $\sum w_i$ active thr.

$V(t_1) = S_0 \cdot t_1 / \sum w_j \text{active} \cdot dt \Rightarrow$ If # active thr. const $t_1 - t \rightarrow V(t) - V(t_1) = (\sum w_j \text{active})(t - t_1)$

Virtual eligible deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Virtual deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Algorithm ① every quantum q ② get $V_e \leq V_d$ ③ select w/ earliest virtual deadline V_d

④ Compute next based on lag ⑤ New $V_e^{k+1} = V_d^{k+1} - (\text{lag})(w_i)$ ⑥ New $V_d^{k+1} = V_e^{k+1} + R/w_i$

* Next eligible t gets impacted by lag → ⑦ lag pulls closer eligible time vs. ⑧ postpares it, ← eligible

⑨ C1 C2

virtual t

physical t

VE 0 0.5 1 1.5 2

physical t 0 1 2 3 4 5 6 7

* C1: $w_1 = 2, r_1 = 2$ quanta

* C2: $w_2 = 2, r_2 = 1$ quanta

q = 1 second

* no lag → $V_e^{k+1} = V_d^k$

Weighted Fair Queuing (WFQ): approx GPS, select packet that finishes first in GPS assuming no future arrivals

* Every client receives at most q delay compared to GPS * CONSIDER finish time dynamically changes/recompute since

Solve with virtual time * Virtual finish time doesn't change w/ new arrivals, always (len packet / w_i)

* System virtual time: index of current round in weighted bit by bit RR

WFQ: F_i^{k+1} (round by which packet k+1 send) = max($V(a_i, k+1)$, F_i^k) + L_i^{k+1}/w_i

F_i^k : V finish time of packet k of client i

L_i^k : len of packet k of client i

"Takes up exactly requested service time" = no lag

Round Robin (RR): each intv with time quantum g

* Large q → higher thrnt → resembles FCFS

* Small q → lots of interleaves

Pros: fair, w/ n process & q, max waiting time $C(n-1)g$, no starvation

Cons: small q increases completion time / suffuses from low throughput

Multi-level feedback Q (MLFQ): queues for ready jobs per p. level

* Approximates SRTF; placed at highest, bump it post Δt given strvnt

Cons: games with I/O req → no matter how many times given up, reduces

Subject to starvation → after time period S, move all jobs to topmost

Proportional Fair sharing: share CPU proportionally to weights

Job share according to priority, run less but no starvation

* Y/n # users ex max-min, user reads less is share ex proportional w/ p.

Lottery scheduling: type of CFS; give ticks, prob win each Δt

* on avg. each job's CPU time proportional to # tickets it has

Cons: could choose long jobs/law prior, unfair for less jobs

Stride scheduling: deterministic proportional fair sharing; stride = #W/Ni → # ticks → more = smaller stride, run more

① look for lowest priority counter ② run ③ add stride; Cons: long delays until job is rescheduled if big stride vs small

Earliest Eligible Virtual Deadline First (EEVDF): schedule tasks for Δt proportional to weight, Minimize lag

Each task has eligible time to run & deadline to finish by

* Generalized processor sharing (GPS): fluid / weighted fair grain switching $f = \text{Capacity}/(w_1 + w_2 + w_3) \approx \min(r_i, f_{\text{avg}})$

Eligible time = $t - \sum_{j \neq i} w_j \text{active} \cdot dt$ can be scheduled ⇒ ① At each t only consider eligible packets ② Pick with earliest deadline

Lag = ideal service time - real service time ⇒ ③ Lag = underallocated → move & schedule earlier vs. ④ punish

$f_{\text{avg}}(t) = w_i / \sum w_i \rightarrow S_i(t_0, t_1) \text{ideal} = f_{\text{avg}}(t_0) \cdot (t_1 - t_0)$ J assume f const in t → E,

⑤ $S(0, 1) = 0.5 * (1 - 0) = 0.5$ s, transmits 0.5 bits ($\frac{1}{2}$ packet) in $t = 0$ to 1, full access to CPU for 0.5s

Virtual time $V(t)$: rand we are cur. at real time t, rate; inversely proportional to $\sum w_i$ active tasks

Moves slower w/ more active tasks since time / round incr. Round duration = $\sum w_i$ active thr.

$V(t_1) = S_0 \cdot t_1 / \sum w_j \text{active} \cdot dt \Rightarrow$ If # active thr. const $t_1 - t \rightarrow V(t) - V(t_1) = (\sum w_j \text{active})(t - t_1)$

Virtual eligible deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Virtual deadline: earliest $V(t)$ thread can be fully serviced $V_d = V_e + R(\text{request interval})/w_i$ (client weight)

Algorithm ① every quantum q ② get $V_e \leq V_d$ ③ select w/ earliest virtual deadline V_d

④ Compute next based on lag ⑤ New $V_e^{k+1} = V_d^{k+1} - (\text{lag})(w_i)$ ⑥ New $V_d^{k+1} = V_e^{k+1} + R/w_i$

* Next eligible t gets impacted by lag → ⑦ lag pulls closer eligible time vs. ⑧ postpares it, ← eligible

⑨ C1 C2

virtual t

physical t

VE 0 0.5 1 1.5 2

physical t

Deadlock	cyclic waiting for resources / implies starvation (but starvation does not), can't end w/o intervention
① Mutual exclusion & bounded rsrcs: 1 thr. at a time	① Provide sufficient resources: use virtual memory
② Hold & wait: thr. holding at least one rsrc waiting for more rsrc	② Abort requests or acquire all rsrcs atomically (both)
③ No pre-empt: rsrcs only voluntarily released, after thr. finish	③ Pre-empt threads (force give up @ transaction abort) Max-con=No!
④ Circular wait: ^{vs. mul. rsrcs} $T_1 \rightarrow T_2 \rightarrow \dots \rightarrow T_n \rightarrow T_1$ waiting threads	④ Order rsrcs & always acquire in same order Total - Curvnt = Avail Demand current to avail

[Avail] = [Free Ksres] Avoidance
 Add all thr to UNFINISHED ★ IF UNFINISHED is
not empty,
Dual lock possible
 do {
 DONE = true
 for each NODE in UNFINISHED {
 [Request] = [Max Node] - [Alloc Node]
 if ([request] <= [Avail]) {
 remove NODE from UNFINISHED
 [alloc] += [request]
 }
 }
 }

$\text{[Avail]} \leftarrow \text{[Alloc NODE]}$
 $\text{DONE} = \text{false}$

Segmentation: ~~segmented-reg seg. b/w seg. selector registers - which seg. curr. memory processes~~ cont. portion of address space of particular len; place each seg indep at diff locations, base/bounds per logical seg.

Paging: divide logical space of processes into fixed size chunks (pages); PA = array of fixed size slots (page frames), contain VM pg.

#VA pages = $\lceil \frac{\#VA bits}{\#PN bits} \rceil$

Demand paging: active pages in ram, others on disk + marks PTE
 (cons: huge PT size, time overhead access internal fragmentation), sparse

Demand paging: active pages in RAM, others on disk + ^{initial}
 Page size (bytes) = 2^{# of set bits}
 Copy-on-write: copy point addy space, page sharing + marking only
 Multilevel paging:

VPN1 (base+PT1)	VPN2 (next PT)	offset
-----------------	----------------	--------

VM Size (bytes) = # VA pages × PTEs

PTE entries = # V pages

Zero fill on demand: New data pages carry no info, mark PTES invalid
Data block pts: debug, mark page R-only, fault if exec.

* If region unused → mark entire inner region invalid
 Cons: still big, 2 km address. Pros: good if spaces sparse

PT entries - # V pages	each block pte, dev, mark page R/W/R, fault if exec.	cons: still big, 2 min entries pros: good space space
PT size = # PT entries x PTE size (bytes)	processes share page, map VA \rightarrow same frame, protection bits	TLB: recent addr translations, fully assoc. \$ NOT actual data

Paged segmentation: Virtual seg # VPN # offset $\xrightarrow{\text{seg #} \rightarrow \text{base}, \text{check limit}} \text{Get PT addr} \rightarrow \text{VPN} \rightarrow \text{check permission bits} \rightarrow \text{PPN} + \text{offset}$

Inverted Page Table: each entry: physical page + pid + which VA maps to phys. \star Size proportional to physical mem, \ll virtual mem
 Anchor Hash Table: Hash(VA + RTD) \rightarrow mul physical frames \star Relies on good hash! otherwise VAs all map same PP & kickout

$$\text{Cactus AMAT} = (\text{HT} \times \text{HT}) + (\text{MR} \times \text{MT}) \quad \text{AMAT}_{\text{LI}} = (\text{HT}_{\text{LI}} \times \text{HT}_{\text{LI}}) + \text{MR}_{\text{LI}} \times \text{MT}_{\text{LI}} \quad \text{MPLI} = \text{AMATL2} = \text{HT}_2 \times \text{HT}_2 + \text{MR}_2 \times \text{MT}_2$$

Replacements: FIFO: evict longest in term, fair due to equal time. MIN: replace pages used furthest in future; optimal but needs future \$
 LRU: temporal eviction, approx MIN, expensive in HU. $\$ = \frac{\text{stack}}{\text{set pages of size } n} \leq \text{pages} \times \frac{1}{n+1}$. Belady's anomaly: $\uparrow \$ \text{ capacity} \uparrow \text{MR}$ (yes FIFO, no MIN/LRU since stack)

Clock algorithm: ^{approx LRU, old & new pages} Page empty, Use 0 → bring in A:1, rest 0. Hit: return bit=1, don't advance. Miss: 1) 0, if 0 evict & swap, advance. Demand refilling count: working set: subset in ref for space for page resident w/1: subset older space in memory - If Brooks 64 entries, our RAM

Thrashing: norm. too small for working set, const replacements @ 256KB WSS, 4KB pages $\rightarrow 2^8 \cdot 2^{10} / 2^2 \cdot 2^{10} = 2^6 = 64$

Page Fault handling: VA \leftrightarrow PA fails (invalid PTE, access violation, DNE) \rightarrow Trap \rightarrow engage OS (alloc more stack pages, make many accessible (reclaim W), bring page from 2nd-ry mem (demand paging)) ★ Need to execute SW to allow HW to proceed

* On-chip \$ \leftrightarrow SRAM \leftrightarrow DRAM (main mem, \$ for disk) \leftrightarrow Secondary Storage (SSD) \leftrightarrow TAPF tertiary

Invalid PTE (not in mem, find in disk): ① MMU traps into OS ② OS (A) chooses old page to replace
③ If old page modified, W to disk ④ Change PTE + TLB → invalid ⑤ Load new page into mem from disk

⑤ Update PTE, invalid TLB for new entry ⑥ cont. thr. from original faulting location ⑦ TLB for new page
invalid when thr cont. ⑧ while null-in. pages from disk, OS runs another process from ready Q, place off in wait

(ex) PA: 8GiB, Page 8KiB, PTE 4 bytes, needs PT to map 46 bit VA if every PTE fits in 1 page? (x) w/o TLB: 1 PA access = 3 PT lookups + 1 data access

$$\# \text{bytes mapped} = \# \text{pages pointed by table} \times \text{page size} = (\text{size PT}/\text{size RTE}) \times 2^{\text{size PT}} = (2^{13} / 2^2) \times 2^{13} = 2^{24} \text{ bytes}$$

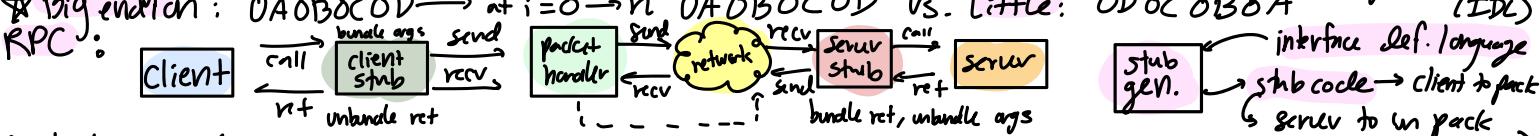
$\hookrightarrow 2^n$ entries per table \rightarrow need 2^{46} mapped; find n where $(2^{11 \times n}) \times 2^{13} = 2^{46} \rightarrow 46-13 = 33 = 11 \times n \rightarrow n=3$ levels

Digits	2^2	4^2	2^8	16	Ki 10	$1 \text{ ms} = 10^{-3}$	Type	tos	cons	Little endian (LSB smallest)
--------	-------	-------	-------	---------------	-------	--------------------------	------	-----	------	------------------------------

Digits	2^2	4	2^8	16	Ki 10	$ms = 10^{-3}$	Type	Pros	Cons	★ Little endian (LSB at smallest)
0000 0 0101 S 1010 A	2^3	8	2^9	512	Mi 20	$1\mu s = 10^{-6}$	Simple Segmentation	Fast context switching (segment map maintained by CPU)	External fragmentation	<u>FS F4</u> bottom 4 permission
0001 1 0110 C 1011 B	2^4	16	2^{10}	1024	Gi 30	$1\mu s = 10^{-9}$	Paging (Single-Level)	No external fragmentation Fast and easy allocation	Large table size (~ virtual memory) Internal fragmentation	$20 \rightarrow 12 \rightarrow 2000$ (X20AH) $\times 52 \rightarrow 16 \times 2 \text{ bytes} \rightarrow 0x4A \text{ VPN}$
0010 2 0111 7 1100 C	2^5	32	2^{11}	2048	Ti 40	$1ps = 10^{-12}$	Paged Segmentation	Table size ~ # of pages in virtual memory	Multiple memory references per page access	★ PA: Base + VPN × (PTE size bytes)
0011 3 1000 8 1101 D	2^6	64	2^{12}	4096	Pi 50		Multi-Level Paging	Fast and easy allocation		★ Renation: keep highest until release
0100 4 1001 9 1111 E	2^7	128	2^{13}	8192	Ei 60		Inverted Page Table	Table size ~ # of pages in physical memory	Hash function more complex No cache locality of page table	★ Next renamer is writer w/ highest off pri Keep memory free space for future allocation

$WSS \leq RSS \leq PHYS-MEM$ (minimizing phus= RSS fits PHYS-MEM)

Distributed Systems Client/server: It server, client makes remote procedure call, server requests
 Peer-to-Peer: no hierarchy, gossiping ★ Availability (proportion of time sys in func. cond, use backups)
 ★ Fault-tolerance (well defined behavior if fault) ★ Scalability (can add more servers) ★ Transparency (easy UI, hide location, migration, replication, concurrency (how many users), parallelism (split job), fault tolerance
Protocol: agreement (syntax: how structured) (Semantics: what a com means) ★ Marshalling = serializ.
 ★ Big endian: DAOBOD → at i=0 → n DAOBOD vs. Little: ODOBODA (IDL)



★ which mbox/dest to send to? Binding: convert user-visible name ↔ network endpoint (dynamic w/ service)
 ★ CONS: RPC not performance transparent (cost high bc marshalling, stubs etc. Same-machine RPC vs. RPC)

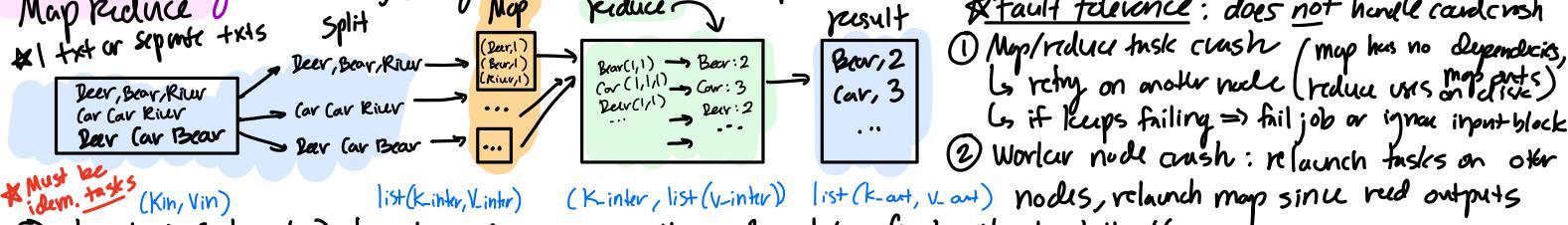
Dist File Systems: Mount files (local can refer to remote), global namespace or (host name, local name)
 Virtual File System Switch (VFS): allow same syscall API to be used on diff file systems in same hierarchy
 ↳ Superblock (mounted file sys) inode (specific file) devtory (dir) file (open f, assoc w/ proc) might take time
 ★ CONS: performance (Network > local mem, server bottleneck) ⇒ caching ⇒ consistency problems
 ↳ Dealing w/ failures? ↳ remove file w/o server ack, changes in buffer cache, shared state across RPCs + poll
 ↳ Stateless protocol: all info to service is included Idle timeout: 3x = 1x = nx execution traffic

Network File System (NFS): status, idempotency, RPC for file ops write-thru \$ (modify disk before ret to client)
 ↳ Transactional failure model to client ↳ Weak \$ consistency: client polls server for changes, may use old ver, arbitrary
 ↳ Transient failure model to client ↳ Weak \$ consistency: client polls server for changes, may use old ver, arbitrary

Internet: IP (can use any network) ⇒ Layering (cons) (may duplicate N-1 functionality, may need same info, hurt performance (hide details)), not cleanly separated (dependencies) e2e: hosts implement themselves, check by seq. order

★ E2E == increase network complexity, delay/overhead, host sometimes reduces complexity, conservative (all) moderate (lower)

Data processing: Bad: msg passing b/w nodes ⇒ Breakup tasks & distribute to workers (9 coord)



③ slow task (straggler): launch 2nd copy on another node, take first output, kill other

★ CONS: MR stores intermediate state on disk ⇒ skew

Spark: efficient handling of intermediate datasets, high level ops (SQL, Stream etc), reliable virt. mem
 ↳ Resilient Distributed Datasets (RDDs): immutable collection of objects stored mem/disk, auto rebuild on failure
 ↳ Transformations: apply f(RDD) → new RDD Actions: return val to phy or W to disk

Librairie: layered seq. of transformations to create RDD ★ Failures: recomputate RDD partition w/ lineage

Hadoop: single namespace for whole cluster, replicate data 3x, MR framework, Colocated w/ file sys, append only

Coordination: General's Paradox: if network is unreliable, impossible to guarantee 2 entities do smthg simultaneously

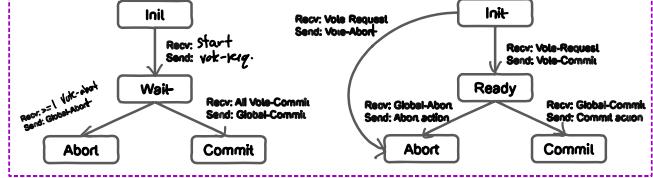
★ If malicious, impossible if less than 3f+1 parties (f misbehav) eventually choose to commit or abort

2 Phase Commit (2PC): automatically agree, no constraints on time, either commit or abort transaction

★ Agreement (same decision) ★ Finality (cannot reverse decision) ★ Consistency (if no fails + yrs, commit) ★ Termination (no faults / repair)

① Spray VOTE-REQ ② Send VOTE-COMMIT or VOTE-ABORT (if any ABORT, abort) ③ collect votes, if all commit then GLOBAL-COMMIT else GLOBAL-ABORT ④ GLOBAL dictates local commit or abort

Coordinator (FSM) Worker



★ Stable storage (assume logs & state restored) ① VOTE-REQ but no decision → Abort, voice & no vote → VOT ABORT

★ 2PC not subject to paradox bc all nodes eventually converge to same decision (not at same time)

★ 2PC CONS: blocking VS. 3PC (allow progress if block), PAXOS ↓

Consensus: Termination (eventually decide) Agreement (some val) Validity (proposed by some process)

↳ Impossible on asynch system ⇒ we still do it b)

★ Coordinator can self abort if waiting too long → GLOBAL ABORT
 ★ Worker waiting after voting → timeout & terminate
 ↳ Coordinator crash if waiting too long → termination protocol:
 ① Block & wait until card receives (get GLOBAL-*)
 ② Ask other worker "p" (optimization) ③ p voted COMMIT, p abo no stuck
 ④ p decided Commit/abort → final decision ⑤ p not decided, VOT ABORT → help decide

★ Stable storage (assume logs & state restored) ① VOTE-REQ but no decision → Abort, voice & no vote → VOT ABORT

★ 2PC not subject to paradox bc all nodes eventually converge to same decision (not at same time)

★ 2PC CONS: blocking VS. 3PC (allow progress if block), PAXOS ↓

Consensus: Termination (eventually decide) Agreement (some val) Validity (proposed by some process)

↳ Impossible on asynch system ⇒ we still do it b)

PAXOS: Safety (consensus not violated) & Eventual liveness (high chance of consensus, no guarantee)

* NO COORDINATOR, NO upper bound on msg delay (eventually will arrive) 2f+1 nodes, f fail
 Rounds: unique ballot id, asynch, j vs. j+1 round → take j+1 → 3 phases / round
Diagram: shows a timeline with three phases: prepare, propose, and commit.br/>
Election: majority on ballot id → once elect cannot pick smaller id
Bill: if participant already got b>ballot_id, ignore propose, else accept
Law: majority ACCEPT (ballot_id, v), send COMMIT (ballot_id, v),

* Version 2 election: G sent ACCEPT (old-ballot_id, v)
Version 2: Prepare(id) → if b>id seen → send PROMISE(id) BUT if already
 If majority promise(id) → elect, if majority promise(id, old_id, old_val) → pick v, w/ highest old-ballot_id val

Core safety theorem: If 2f+1 participants accepted v in round r, then for all rounds r'>r, proposer gets at least one PROPOSE(v', r') {0 0 0 0 0} * Always 1 overlay * Cannot prove liveness

General I/O: Device types: Block (green box)
 Block (raw I/O ex disk/read), Character (single char / func, get(), keyb)
 mouse (USB) Network (unique I/O, socket interface)
 Single bus: O/I layer: device driver. Bus: high BW but slow
 Single bus: O(n^2) relationships w/ 1 set of wires, but blocks all except 1 Parallel: hard to sync.

Device controller: registers (R/W) ⇒ programmed I/O, CPU involved in every data transfer, heavy polling
 ① Port-mapped I/O: separate addy space (in/out instr to comm. w/ device registers)
 Mem-mapped I/O: looks like mem access, consumes local addy space (id/str instr)

Direct Mem Access (DMA): CPU gives DMA req. to device controller
 When done DMA controller informs CPU DMA has direct control of mem. bus
 If uses DMA, device driver (device specific code in kernel that interacts w/ HW) can split into 2 pts
 Top half: accessed in call paths via syscalls, kernel interface, Bottom: interrupt routine, get in/transfer, wake sleeping I/O
 Blocking: wait until data ready Non-blocking: immediate return nothing Asynch: ptr to buf, immediate ret, having fill

Storage Devices HDDs (Magnetic): rarely corrupted, large capacity, block-level, slow for random, good sequential
 Latency = Q time + controller + seek + rotational + transfer
 Seek: position head on track vs. if same cylinder
 Rot: wait for desired sector (here if same track + sequential)
 Transfer: R/W min time, transfer sector under head
 ex 4096 sectors, 3×10^6 sectors/tracks, 100 tracks/platter, 2 platters, 5400 rpm, 5.6 ms seek, 1 ms cont + Q
 1 Size? $(3 \times 10^6 \times 4096) \times 100 \times 2$ 2 Throughput 64 KB read? $1 \text{ ms} + 5.6 \text{ ms} + (\frac{1}{5400 \text{ min}} = 11.1 \times 2 = 5.5 \text{ sec}$
 + 64 KB / 100 MB/s = 0.457 ⇒ total t = 12.6 ⇒ throughput = 64 KB / 12.6 **Block size ≥ Sector size, BOTH use block API**

Disk scheduling: FIFO (fair but random), SSTF (shortest seek) (may starve), SCAN (closest req. in dir of travel, no starvation) \geq C-SCAN (one-dir, skips otw blocks, fairer than SCAN since not biased towards middle)

SSDs (storage state): no moving parts, limited W, low power **R/W pages, erase blocks (page < block)**
 Block access in parallel > Blocks (64-256 pages) > pages fw 1KB > Cells (1-4 bits) > 2-3 (MLC) wears faster

Flash: Invalid → erase → valid JEE Latency = Q + controller + transfer
 No seeks, expensive / smaller Flash translation layer (FTL) = translate logical block to flash blocks, reduce amplification (control / grp W traffic), avoidular, indirection (DMA mappings) copy on W (w new page tag update ptr)

File Systems: Most files small **Most bytes in large files** **Fixed size blocks** $\xrightarrow{\text{ptr size}}$

FFS: File # → inode array → inode (metadata + DP + IP) → DP or array of pointers $\xrightarrow{\text{ptr size}}$
 ex 12 dir + 1 Dbl + 1 Indir (4 KB block) = 12 × 4 KB + 1024 × 4 KB + 1024² × 4 KB $\xrightarrow{\text{1024} = 2^{12}}$

New file → find free inode + free DB **Disk optimizations:** dist inodes on same cylinder (reduce seek), continuous alloc **Dir + files in same block groups + free map** 10% free inode + dir entry / 2+ diren.

Pros: seq., random, no external frag., big files, locality for file + meta, ok for small

Cons: tiny (read inode + whole block) → internal frag., but DP do ok, bad encoding, 10% free

FAT (File alloc table): file # → index into table → LL of blocks, free list (LL), append to grow Dir: < file name : file # > file attributes stored in dir not file itself → linear search

Pros: sequential, no external frag., big files **Cons:** random (LL traverse), internal frag. (if 2 many small) stored sequentially

NTFS (new tech): file # → index MFT → Record: file meta + data OR ptr to extents (variable len chunks of data) OR if massive, ptr to lists of extents (like dbl); MFT always has file name

Hard link: < name: file # >, link(), part of 2+ dirs, cannot be removed until all file gone **NOT supported by FAT**

Soft link: < name: path >, symbolic(), no refcnt, if path DNE, fail **Since no inode, cannot make mult directories?**

Buffer cache: temporal locality, LRU valid unless scanning which kicks delayed W/write back periodically **Paging uses clock**
 Cons: interrupted → loss of shared data → ① careful order ② copy on W: make new file, then pt ptr to there

Durability: being able to recover from disk failure. **Reliability:** after recovery want consistent state

Transactions: unit, atomic seq of R/W that goes from consistent state 1 → state 2

Atomicity (valid state to another) **Consistency** (valid state to another) **Isolation** (concurrent exec. ok) **Durability** (once committed, not erased)

Log (data stays in log form), write commit to log ↓ **Journaling** (log for recovery) **Journaling** (log for recovery) Write data to log; once transaction written, disk apply changes, then remove transaction

↳ Case 1: sys crashes before Wr to log → incomplete transaction → no changes applied ***Idempotent**

↳ Case 2: sys crashes applying changes to disk → re-apply transaction by looking at log **(ex ext3 FFS+ Jlogs)**

Discussion Problems

Disc 7

Corrupted sectors: spread backup evenly, checksum, incr. perceived wr speed by signaling from buffer, platter

Thruput ↗ seek times, 7200 RPM, transfer rate 50 MiB/s

① random access, 4 KiB throughput? ② 4 KiB same track?

G on avg 1/2 revolution G No seek time

G $\frac{1}{2} \times 1/7200 = 4.17 \text{ ms rot.t}$ G $4.17 + 0.078 \dots$

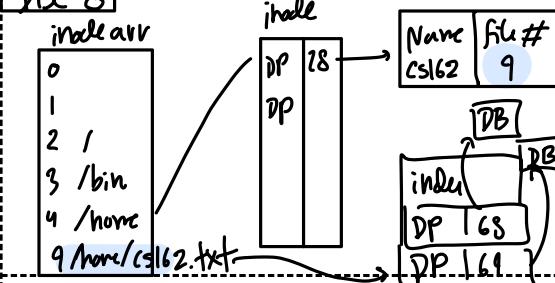
Transfer t: $4 \text{ KiB} \times 1/50 \text{ MiB} = 0.078125$ ③ Next sector?

$t = 8 + 4.17 + 0.07 = 12.248$ G No rot/seek

Total throughput = $4 \text{ KiB}/12.2 = 326.6 \text{ KiB/s}$ G just transfer rate

in background ↗
not
buff[28]
if (id->indir == 0)
 if (id->dir[i] == 0) **Shrink**
 if (size <= S12 * i + 8 & id->dir[i] == 0)
 free (id->dir[i])
 else if (id->dir[i] == block->alloc)
 id->dir[i] = block->alloc();
 else
 block->read (id->indir, buffer);
else if (id->indir == 0 & size <= 12 * S12)
 id->len = size;
 ret true;

Disc 8



Quick format: mark all blocks free **Full:** quick + zero out DB
FAT vol size / block size = 2¹⁹ entries × 4 bytes size = 2²¹ bytes
***FAT32 → 4 GiB size limit**

FAT steps D:\MyFiles\video.mp4

- ① Read in blocks of / until MyFiles found (name: file#)
- ② Read in blocks of /MyFiles until video.mp4 found (name: file#)
- ③ File # = block 0 → read all blocks until EOF

```
for (int i=0; i<110, i++)  
  if (size <= (12+i) * S12 && buffer[i] == 0)  
    block_free (buff[i]) Shrink  
    buff[i] = 0  
  else if (size > (12+i) * S12 && buffer[i] == 0)  
    buff[i] = block_alloc(); grow
```

```
if (size <= 12 * S12)  
  block_free (id->indir);  
  id->indir = 0 Shrink indir grow  
else  
  block_ur (id->indir, buff);  
id->len = size;
```

```
Sector alloc failures  
sector = alloc_block()  
if (sector == 0)  
  inode_resize (id, id->len)  
  ret false
```

Disc 9 Availability: probability sys. can accept & process req (99.9% = 3 nines of availability)

Durability ≠ Availability → Reed-Solomon ECC (short term) RAID (long term)

***commit if commit → Nothing applied OR intermediate state OR all updates written to disk**

Linux uses -f vfat to access files in flash drives formatted w/ FAT32

NFS: may have inconsistent views due to pulling/caching vs. local file sys. sas truth

Disc 10 completed map tasks re-execute when worker crashes, coord. finds locations of map results to reduce workers → reduce worker includes RPC to read data
Lazy eval. transformation → avoid unnecessary comp., optimize reading knowing computation graph

```
map(key, val){  
  reduce (key, list<val>) {  
    larest = INIT_MAX  
    for (val : List<val>)  
      larest = min(larest, val)  
    }  
    Emit (problem name, larest)  
}
```

2PC: optimize by ignoring about decisions → reduce logging
***Blocked worker back, may hang vsyncs**

worker	send/recv t	log	latency
w1	600	10ms	3ms
w2	300	20ms	5ms
w3	200ms	30ms	parallel comp

① Min time for 2PC to complete successfully
freq host = worker recvs + log + send back + above logs
= $\max(400+10+400, 300+70+300, 200+30+200) + 5$

Commit/abort phase = worker recvs + log result + send back (no need log)
max(400+10+400, 300+70+300, 200+30+200) \Rightarrow Total = Prep + Commit

② All 3 commits, broadcast comit, w2 crashes & recvs after timeout of C → transaction comit/abort? latency?
Still commits

Prep host = 815 (idle lost)
Commit: freq 1: 300ms (timeout)
freq 2: 300+70+300
total = 815 + 3000 + 670

- Deadlock**: cyclic waiting for resources, implies starvation (but starvation does not), can't end w/o intervention
- Mutual exclusion & bounded vsrcs: 1 thr. at a time
 - Hold & wait: thr. holding at least one rsrc waiting for more rsrcs
 - No pre-empt: rsrcs only voluntarily released, after thr. finish
 - Circular wait: {T1...Tn} waiting threads T1 → T2 → ... → Tn → T1

$[Avail] = [\text{Free Rsrcs}]$

Add all thr to UNFINISHED
done
 $DONE = \text{true}$
for each NODE in UNFINISHED:
 [Request] = [Max NODE] - [Alloc NODE]
 if [request <= [Avail]]:
 remove NODE from UNFINISHED
 [Avail] += [Alloc NODE]
 DONE = false

3
until (DONE)
Safe state vs.
unsafe (will
invariably
lead to)

Avoidance
★ If UNFINISHED
not empty,
deadlock possible

(1) Deadlock Init A=0, B=4,
void threadA() { 12 max
 request(3) } executes
 request(6) ↳ 12-4-3 = 5
 free() blocks left

(2) Deadlock void threadB() {
 if rand(1000) {
 request(6) } blocks
 free() ↳ free in A
 or B!

- Provide sufficient resources: use virtual memory
- Abort requests or acquire all rsrcs atomically (both)
- Pre-empt threads (force give up ex transaction abort)
- Order rsrcs & always acquire in same order

Total - Current = Avail
Denote current to avail

(3) Uniprogramming: app always in same PA, 1 at a time, any PA/no protect
Mem transl. Thru Relocation: use loader/loader to adjust addr, no protect+
Base & Bound: registers, Base + addr offset < Bound (call, data/
stack) / stack
★ Kernel executes w/o Base/Bound reg., Loader rewrites addr to control
offset
Limits (1) No expandable mem: static w/ relocation: loaded at addr &
 w/ compiled as if loaded at addr &
(2) No mem sharing b/w processes: PA = VA + base & translated by MMU in CPU
(3) Non-relative mem addr.: location of code/data at run time (1) relocate
(4) External Frag.: cannot relocate programs (2) Relative to adjustable reg., can
 relocate process
(5) Internal Frag.: addr space must be cont (4) More processes around &
 consolidate

Segmentation: shared-mem seg. btw processes * seg. selector registers - which seg. curr. running
cont. portion of addr space of particular len/ place each seg. indep at diff locations, base/bounds per logical seg.

* minimize internal frag. * Seg map in processor, Seg # ↔ Base/limit pair; Base + offset = PA (Ex 0x8020 → 100,000 0010 0000)

seg # | offset @ 4 seg → log2(4) = 2 bits @ 16 bit addr → 16-2 = 14 bit offset → Seg size = 2¹⁴ bytes Seg 2 offset 0x20 → 0xF000 + 0x20 = 0xF020 (curr. segment)

* HW adds protection bits + limit Pros: support sparse oddy spaces, avoid internal frag, min. HW req/efficient transl. Cons: external frag, many processes

Paging: divide logical space of processes into fixed size chunks (pages); PA = array of fixed size slots (page frames), contain VM pg.

VA | offset + PTE → VPN index into PT → PPN | offset * PT resides in mem. * each process has PTBR (base), saved/reload on context switch

VA pages = 2^{VA bits} # VPN bits # of offset bits Demand paging: active pages in mem, others on disk + mark invalid

Page size (bytes) = 2^{VPN size} Copy-on-write: copy parent addrspace, page sharing + marking R only

VM size (bytes) = # VA pages × Pgs size Zero fill on demand: New data pages carry no info, mark pages invalid on page fault

PT entries = # V pages Data blocks pts: debug, mark page R only, fault if exec.

PT size = # PT entries × PTE size (bytes) * Process share page, map VA ↔ same frame, protection bits

PTE size (bytes) = # PTE bits / 8 * Kernel region of each proc. has same PTE, running same binary → same seg.

Paged Segmentation: Virtual seg # | VPN # | offset * seg # → Base, check limit → Get PT addr → VPN → check permission bits → PPN + offset

Inverted Page Table: each entry: physical page + phys. & which VA maps to phys. * Size proportional to physical mem, << virtual mem

Anchor Hash Table: Hash(VA + PID) → mul! physical frames. * Relies on good hash! Otherwise VAs all map to same PB & kickout

Cactus AMAT = (HT × HT) + (MR × MT) * AMAT.L1 = (HT_{L1}) × (MR_{L1}) + MR_{L1} × MP_{L1} * MP_{L1} = AMAT.L2 = HT_{L2} × R_{L2} + MR_{L2} × P_{L2}

Replacements: FIFO: evict largest in mem, fair due to equal time MIN: replace pages used farthest in future, optimal but needs future stack = set pages of \$ size mem in stack since stack

LRU: temporal eviction, approx MIN, expensive in HW * Belady's anomaly: ↑ \$ capacity ↑ MR (yes FIFO, no MIN/LRU since stack)

Clock algorithm: approx LRU, old & new pages → bring in A:1, next. Hit: set turn bit = 1, don't advance MISS: → 0, if 0 evict & swap, advance

Demand paging cont: working set: subset in addr space for exec. resident set: subset addr space in memory TLB reads 64 entries

Thrashing: mem. too small for working set, const. replacements (Ex 256 KB WSS, 4 KB pages → 2⁸ × 2¹⁰ / 2 × 2¹⁰ = 2⁶ = 64)

Page Fault handling: VA ↔ PA fails (invalid PTE, access violation, DNE) → Trap → engage OS (alloc more stack pages, make page accessible (copy on write), bring page from 2nd-level mem (demand paging)) * Need to execute SW to allow HW to proceed

* On-chip \$ ↔ SRAM ↔ DRAM (main mem, \$ for disk) ↔ Secondary Storage (SSD) ↔ TAPR tertiary

Invalid PTE (not in mem, find in disk): (1) MMU traps into OS (2) OS (A) chooses old page to replace

(B) If old page modified, W to disk (C) Change PTE + TLB → invalid (D) Load new page into mem from disk

(E) Update PTE, invalid TLB for new entry (F) cont. thr. from original faulting location (3) TLB for new page located when thr cont. (4) while pulling pages from disk, OS runs another process from ready Q, place old on wait

(Ex) PA: 8GiB, Page 8KiB, PTE 4 bytes, levels PT + map 46bit VA if using PT¹¹ * (A) w/o TLB: 1 PA access = 3 PT lookups + 1 data access

bytes mapped using PT = # pages pointed by table × page size = (Size PT / size PTE)¹¹ * (B) w/ TLB: Best: 1 data access Worst: 3 PT + 1 data

= (2¹³ (since each PT = size of page) / 2²) × 2¹³ = (2¹¹) × 2¹³ = 2⁴⁶ bytes * AMAT TLB: ((0+S0) × 0.S + (1-S0) × (10+S0+S0))

↳ 2¹¹ entries per table → need 2⁴⁶ mapped; find n where (Z¹¹ × n) × 2¹³ = 2⁴⁶ → 46-13 = 33 = 2¹¹ × n → n = 3 levels * AMAT NOTLB: S0 + S0 = 100 * S0 = main mem, TLB = 10

fetch PTE * read data

Digits 2² 4 2⁸ 16 Ki 10 | ms = 10⁻³ **Type** **Pros** **Cons** * Little endian (LSB smallest)
0000 0 0101 5 1010 A 2³ 8 2⁹ 512 Mi 20 | ms = 10⁻⁶ PSF → bottom 4 permissions
0001 1 0110 6 1011 B 2⁴ 16 2¹⁰ 1024 Gi 30 | ms = 10⁻⁹ 20 07 → top 12 → x2000 (0x200)
0010 2 0111 7 1100 C 2⁵ 32 2¹¹ 2048 Ti 40 | ps = 10⁻¹² 0x200 → 16 × 2 bytes → 0x AU VPN
0011 3 1000 8 1101 D 2⁶ 64 2¹² 4096 Pi 50 | ps = 10⁻¹⁵ * PTA: Base + VPN × (PTE size bytes)
0100 4 1001 9 1110 E 2⁷ 128 2¹³ 8192 Ei 60 | ps = 10⁻¹⁸ * Denomination: keep highest until release
 | ms = 10⁻¹⁸ * Next runner is water w/ highest off pri, then promoted by any donor's mem cont
 | ms = 10⁻¹⁸ * B C B D BA * C = donor 100 > B50 > A1
 | ms = 10⁻¹⁸ * D = donor 50 > B50 > A1

WSS ≤ RSS ≤ PHYS-MEM (simplifying phys = KSS fits PHYS-MEM)