

Memory

MMU: Translate Vaddr \rightarrow paddr
 on every load/store/instruction fetch
 \rightarrow Read PTE, check valid, merge
 \rightarrow Set "accessed" bit; "dirty" on write

Base: Demand
 set registers w/base and limit

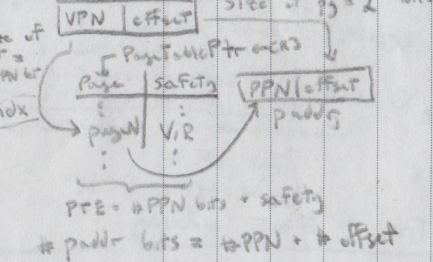
- Simple between blocks
- Internal + external fragmentation
- No spare addr; space support
 \rightarrow for inter-process sharing
 e.g. chunks for code/data/stack/heap
- Both reserved w/segmentation
- May need to swap on context switch (i.e. store process mem to disk)

Paddr bits = # PPN + # offset
 # Vaddr bits = # VPN + # offset
 $\text{VNUM}(B) = \# \text{V.pgs} \times \text{pg.size}(B)$
 $\text{pg.size}(B) = \frac{\text{PTE size}(B)}{\# \text{PTEs}} = \frac{\text{PTE size}(B)}{\# \text{V.pgs} \times \text{pg.size}(B)}$

PT size (B) = #V.pgs \times PTE (B)
 # PPN + # mem (B) = # PPN + # mem (B)

Each proc. has Page Table in p. mem.
 maps VPN \rightarrow PPN + safety bits
 (pages)
 Fixed size chunks, no external frag.
 \rightarrow Prevent internal frag. by limiting to ~1-16 KiB; use numerous pages per segment

Shareable mem: diff. VPNs map to same PPN (proc.)
 \rightarrow E.g. Kernel region has some PTE, shared exec-only user libs
 Scales w/VPN bits; sparse (external) valid



Demand Paging: Treat DRAM as cache on disk; PT tracks pages in memory not in mem \Rightarrow page fault, as loads to mem

Page Fault: Vaddr not mapped
 (1) MMU traps kernel, saving res/state
 (2) Kernel stores needed VPN in CR2
 (3) Check validity, move page to disk and load + update PTE
 (4) Restore saved state, resume exec.

Inverted Page Table

- Other "Forward" PTs have size of PT \geq vmem alloc to processes
- Pract. may be less than VMEM, i.e. on disk
- Use Hash table, mapping VPN \rightarrow PPN
- Size indep of vaddr space, good for 64B addr.
- Directly related to amount of physical mem.
- Peer cache locality
- Hash chains complexity

Multilevel: Chain PT in "tree" structure

- Much smaller: alloc PTEs when needed
- ≥ 2 lookups per reference, diminishing returns
- Unshared PT can be stored on disk (e.g. during context switch)
- Share granularity at PT/page scales
- PT must be contiguous
 \rightarrow 10-10-12 keeps table to exactly 1 pg in size

Replacement Policies

- FIFO (que). Evict oldest
 Evicts heavily used p. frames
 Does suffer from Belady's anomaly
- LRU.
 Overhead tracking + comparing (over sets) (n! compares)
- Random.
 Empirically on par w/ LRU
 Temporal locality is low
- MIN. Replace entry not used for longest time
 Non-causal; Threshold opt

Types of Misses:

- Compulsory. First access to block
 Avoid w/ prefetching (sequential data reads)
 Insignificant in large scheme (could also track working set size)
- Capacity. Too many blocks; finite size
 Increase cache size
 Reduce associativity, e.g. access pattern 0, 32, 64, ... 32N, 0... An N-line direct map cache just misses as 0; 32N conflict while fully-assoc LRU misses all
- Conflict/collision. Multiple mem. locations mapped to same cache location
 Increase cache size/assoc.
- Coherence. Other processes (I/O) update memory
 Due to parallelism

Set-Associative N-entries per cache idx

- Index refers to set of rows near
 \rightarrow Check all tags; size scales \propto #sets, n
- Reduces "pinging" of conflict misses
- Comparators scale w/n

10-10-12 nm 4KiB pages

- Space of one page mapped = $2 \times 4 \text{ KiB}$
- 2^{10} entries of {ptes/data} = 4 B, tags \Rightarrow PTE = 2^{10} B = 4 KiB = pg size w/ 10x
- Top level page references to = pg size w/ 10x
 $(1+1024) \cdot 4096 = 4096 + 2^{22}$ B

Avg. Mem. Access Time = Hit-p \times hit.time + Miss-p \times Miss.time

Demand Paging cont.

- Working set: addr space used during execution
- Resident set: subset held in mem
 \rightarrow Threshold when too small

Caching

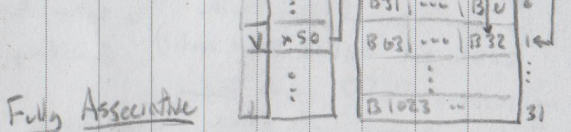
- MMU translation occur on every access, taking ~ PT height
 \rightarrow Worse if PT on disk
 \rightarrow Motivates TLB (translation cache)
 State (VPN#, PTE) + process ID
- On context switch, vaddr changes:
 Flush TLB, or (in hardware) include ProcessID
 \rightarrow Aside: Virt. Index/Tag cache would need to be flushed

- Block: minimum quantum of caching
 \rightarrow offset used to identify data inside block, sometimes not needed
- Index: Identifies set
- Tag: Identifies actual copy

- Write-Through: write to cache + mem
 Read misses + writes; I/O slow
- Write-Back: write to block; mem on evict
 Read miss may need write; I/O

Direct Single block per set/index

- 2^N byte cache has $N \log_2(32-N)$ tag bits
- 65 M bits for byte select (block size 2^N)
- Remaining bits for index (height, # blocks)
 \rightarrow Size: $2^{10} \times 2^{10} \times 2^{10}$
- No eviction policy needed



Fully Associative

- Any block can hold any line \Rightarrow no idx
 \rightarrow must compare w/all cache tags done in hardware, parallel, but difficult
- TLB of ~128-512 entries
 \rightarrow Conflict misses expensive
- TLB backup: cache access can occur in parallel if cache idx first in page offset bits (constant in p.; v. addy space)

Scheduling / Assm. Single thread, indep. programs

Goals/Criteria:

1. Min. Avg. completion time
 - Waiting + run time (response)
 - crucial for time-sensitive I/O

2. Max throughput (ops/sec.)
 - min. overhead (context switches)

3. Fair
 - Is not min. avg. response time
 - ↳ Priority schedule very unfair

Throughput \Rightarrow FCFS

Avg. response \Rightarrow SRTF approximation

I/O throughput \Rightarrow "

Fair (CPU time) \Rightarrow CFS

Fair (wait time) \Rightarrow RR

Meeting Deadlines \Rightarrow EDF

Fairness/Importance \Rightarrow Priority/Lottery

Mittdown

- Resolved by separating PT for user + kernel (eg. Flush twice per syscall, 80% overhead)
- ↳ Hardware PCID avoids flushing when changing addr space

FIFO/FCFS:

- One program runs til done
- Least overhead, simple
- Convey effect; response time depends on priority

Round Robin (RR):

- Each proc gets fixed CPU time q
- $q \geq$ longest path \equiv FCFS
- Circular convey w/ preemption max wait time is $(n-1)q$ (excluding context swaps)
- Lets of context swaps; need $q \gg$ context switch to ensure overhead isn't too great
- Small $q \sim$ SJF; lower response time higher competition "

SJF/SRTF:

- Run job w/ least amount of work to do; non-preempt
- Optimal min avg. completion
- Impossible to know proc length; starvation if many short jobs block long job

Lottery:

- Give jobs tickets; randomly choose
- Give all one ticket \Rightarrow no starvation
- On avg, CPU time \sim # tickets
- Unfair for less jobs
- Gracefully handles oddball jobs

MLFB:

- Early RR; later FCFS } Each gets set amount of time
- Approximated SRTF
- Overly provided for occasional I/O

2nd Chance: Split into Active (RR/FIFO) and SC (unified, LRU)

- Search direct-map Active first
- PF: Pop active, mark invalid push to SC
- ↳ In SC: push head of active, mark R/W
- ↳ new page in H; page out LRU

EDF:

- Given tasks w/ deadlines + compute time
- Predictable; provide guarantees
- ↳ Feasible if $\sum_{i=1}^n \frac{C_i}{D_i} \leq 1$
- by choosing earliest D_i

Strided

- Proportional share time w/o randomness
- Let $Stride_i = W/N_i$ for N_i tickets, const. W w/o, etc
- Pick job w/ least "push", schedule & add stride to its pos
- ↳ More tickets \Rightarrow bear stride, scheduled more often

CFS:

- Illusion of 1/n CPU time
- Track CPU time per thread, schedule lowest
- Sleeping threads don't inflate CPU time get interactivity (best when waking)
- Quanta: $= \max(\text{min-granularity}, \text{target-lateness}/N)$
- Target latency: period within all threads serviced; no starvation
- min-gran. prevents excessive overhead
- Weighted quanta: $Q_i = \text{target-latency} * (w_i / \sum w_j)$
- ↳ Less weight \Rightarrow more time
- vice versa, urgent stride

Clock: Approximates LRU/min repeats an old page

- Circle of in-min PTEs; w/ bit that's set to 1 on reference
- On PF, while hand ptr. to an in-use PTE, it zeroes and advances and advances
- then evicts candidate for replacement
- ↳ If all in use \equiv FIFO
- Partitions into new/old (2 sets)
- ↳ Extend w/ N^{th} chance: evict candidate increments counter, replaced if reaches N

Deadlock

\Rightarrow starvation; but starvation \neq deadlock

Necessary conditions:

1. Mutual Exclusion & Bounded Resources
 - Finite # threads (i.e. 1) can simultaneously use a resource
 - ↳ Infinite sessions, ex VMEM
2. Hold & Wait
 - Thread holding one resource continue to do so when attempting to acquire another held resource
 - ↳ Abort request; acquire atomically; indep. threads
3. No preemption
 - Resources released voluntarily
 - ↳ Fail if waiting too long/don't allow waiting (shows inconsistent state) (inefficient)
4. Circular Waiting
 - Set of threads waiting on each other
 - ↳ Order resource usage

Curry:

- Design specifically so deadlocks impossible, (restrictive)
- Recovery: Roll back actions of deadlocked threads; preserve same execution order
- Avoidance: Dynamically delay resource requests
- Bankers: Stay in safe state (opposite means possible ordering leading to deadlock)
- Threads state resource needs prior.
- Pretend request granted, run deadlock detection algo

Let max_i be requested resource of thread i , alloc: the $\text{unfin} \leq$ all threads; $\text{done} = \text{false}$ current possessions of i

while not done:
 done = True
 For t in unfin :
 if $\text{max}_t - \text{alloc}_t \leq \text{avail}_t$:
 $\text{unfin}.\text{rm}(t)$
 $\text{avail}_t += \text{alloc}_t$
 done = False
 $r \leftarrow \text{len}(\text{unfin}) == 0$ // True if \exists safe ordering

- N -nice \equiv LRU approx
- N loss \Rightarrow more efficient
- Can be used to manage } Dirty pgs write when added to list
- Free pgs list } Allows faster PF handling

$ms = 10^{-3} s$
 $Ms = 10^{-6} s$
 $ns = 10^{-9}$
 $ps = 10^{-12}$
 $1 \text{ KB} = 1024 \text{ B}$
 $= 2^{10} \text{ B}$
 $1 \text{ MB} = 2^{20} \text{ B}$
 $\approx 10^6$
 $1 \text{ GB} = 2^{30} \text{ B}$
 $\approx 10^9$
 $1 \text{ TB} = 2^{40} \text{ B}$
 $= (1024)^4 \text{ B}$
 $\approx 10^{12}$

0	1	7	128
1	2	8	256
2	4	9	512
3	8	10	1024
4	16	11	2048
5	32	12	4096
6	64	13	8192
		14	16384
		15	32768

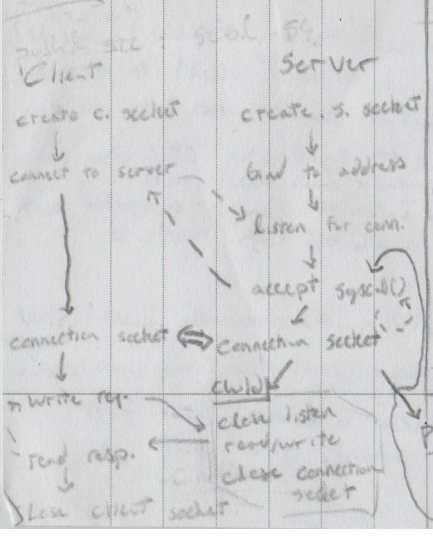
Threads;

C

32-bit \Rightarrow sizeof(xptr) = 32 bits
= 4 bytes
char *a; in read-only static
char b[1]; in stack

struct point { int x; int y; }
struct point P;
P.x = 1; // OK!
struct point a
P \rightarrow x = 2; // SEGS
int x; P[1]
int x; P[0]

struct list_data {
char * names;
struct list * plist; }
struct list_node {
int value;
struct list_elem elem; }
int plist_sum(struct list_data * data)
{
struct list * l = data->plist;
for (struct list_elem * k;
k = l->next; k = k->next) {
struct list_node * temp;
list_elem * k; struct list_node * elem;
total += temp->value;
} return total; }



x86
Syntax: %register; \$const
offset (base, index, scale)
 \hookrightarrow base + index * scale + offset

Calling Conventions:
Caller Prologue:
1. Save needed GPRs to stack
2. Push args reverse-order
(first = lowest address);
padding prior for 16b-align
3. Push ret address (next instruction)
and jump to callee (call)

Epilogue:
1. Pop return value
2. Restore GPR

Callee Prologue:
1. New frame: push ebp,
set ebp to new esp
2. Allocate space for local
variables (overwrote esp)

Epilogue:
1. store return val. in eax
2. dealloc local variables.
set esp to ebp
3. Restore callers ebp
from stack
4. Jump to caller's ret address

pushl src: subd \$4, %esp
movl %src, %esp
popl dst: movl %esp, dst
addl \$4, %esp
call addr: pushl %ebp
jump addr
leave: movl %ebp, %esp
pushl %ebp
ret: popl %ebp
pop: jump to ret on stack
Parent: close connection socket
wait for child

Process		Thread
creation	Fork	P-thread, create
page table	distinct	same
registers	distinct	distinct
Stack	separate; inaccessible	separate; accessible
Heap + shared File Descriptors	separate	shared

Switching threads faster \Rightarrow sharing shared address space

Thread States:

- Running. Currently in CPU
• Spin lock (busy waiting)
- Ready. Eligible, not running
• Just finished I/O
- Waiting/Blocked. Ineligible to run
• Linked on wait queue, sleeping on CV

Secrets two-way;

Pipes one-way
Fits pts to buffer in k. mem

#4 Typed enum { UNL, LOC, CENT }
Acquire(lock *lock) {
if (cmp &swp(&lock, UNL, LOC))
return;
while (swp(&lock, CENT) != UNL) {
Futex(&lock, WAIT, CENT)
}
Release(1) {
if (swp(&lock, UNL) == CENT) {
Futex(&lock, WAKE, 1);
}
unnecessary syscalls
 \rightarrow No overhead if uncontested!

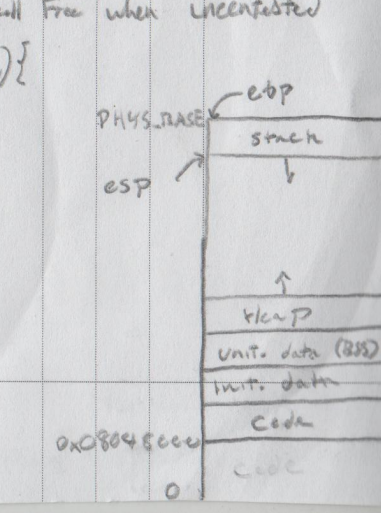
Dual Mode: Enforced in Hardware
• Swap via syscalls, interrupts, & exceptions (external, sync) (internal, sync)
• User can't change table ptr, disable interrupts, access k. mem, or interact w/hardware directly

Lock #2 w/o int. ds. short busy wait
int guard = 0;
int lock = FREE;
Acquire(int *lock) {
while (test &swt(&guard));
if (lock == BUSY) {
sleep & guard = 0;
// guard = 0 on wake race
} else {
*lock = BUSY;
guard = 0;
}
Release(int *lock) {
while (test &swt(&guard));
if (non-empty wait queue) {
wake up
} else {
*lock = FREE;
}; guard = 0;
}

bool maybe = false; \rightarrow syscall free when uncontested
int lock = 0;
Acquire(int *lock, bool *maybe) {
while (test &swt(&lock)) {
*maybe = true;
Futex(&lock, WAIT, 1);
*maybe = true;
}
Release(1, 1) {
*lock = 0;
if (*maybe) {
maybe = false;
Futex(&lock, WAKE, 1);
}
}

Sync. | Lock-less
leave note A
while (note B) { pass }
if (no Milk) { buy Milk }
rm note A
Thread A busy-wait
leave note B
if (no note A) {
if (no Milk) { buy Milk }
}
rm note B
Thread 2

Lock #1
Non-busy: doesn't hog interrupts
Acquire() {
disable interrupts
if (value == BUSY) { sleep(); }
else { value = BUSY; }
enable interrupts
}
Release() {
add to wait queue and enable int.
will re-enable on return
disable interrupts
if (non-empty wait queue) {
pop: give to ready queue
} else { value = FREE; }
enable interrupts
• Can't be used by users
• Not great w/multiprocessors; disable int scales $O(n)$



$$\begin{aligned}
 \text{AMAT} &= P_H \cdot t_H + (1 - P_H)(t_H + t_m) \\
 &= P_H t_H + (P_H t_H + P_H t_m) \\
 &= (P_H + P_H^2) t_H + P_H t_m \\
 &= t_H + (1 - P_H) t_m
 \end{aligned}$$

asm working out fits in DRAM
e.g. no disk ops. required
check in parallel

$$\begin{aligned}
 &\{ \text{Max}(t_1, t_{TLB}) + P_H(t_2 + P_H t_m) \} \\
 &+ \{ P_{TLB}(t_{TLB} + 2(t_1 + (1))) \} \\
 &\text{Get VPN=PPN for TLB} \\
 &\text{Then retrieve actual data (possibly cached)}
 \end{aligned}$$

c/w,

$$\begin{aligned}
 &\{ t_{TLB} + (1 - P_{TLB} P_F) [t_1 + P_H(t_2 + P_H t_m)] \} \\
 &+ \{ P_{TLB} (2[t_1] + P_F(t_1 + t_2 + t_m + t_{TLB})) \}
 \end{aligned}$$

in TLB \Rightarrow no fault possible

FS

- inode-disk; stored in buffer cache (not inode)
 \hookrightarrow holds ptr to data
- inode; in memory, holds sys. primitives

Distributed Sys As opposed to central server for many clients, we desire no central coord (all-way communication)

- Scalability - Add resources dynamically
- Transparency - mask complexity behind interface

General's Paradox: Impossible to confirm lost msg went through (confirm act. simultaneously)

\hookrightarrow Motivates distributed transaction:
2+ machines agree to (not) do smthg in some manner

Two-Phase Commit: Delegated coordinator decides if all procs should commit or abort a transaction
 \hookrightarrow Log ensures guarantee of atomicity and allows recovery (log on disk)
Blocking recovery coord, motivates 3PC & Paxos

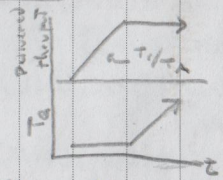
Byzantine GP: One general, $n-1$ lieutenants where f may be malicious

- Integrity Constraints:
- All loyal lieutenants obey same order
 - If general is loyal all follow his order
- Can't solve $n=3$, need $n > 3f$

Queue

$\mu = 1/T_s$: service
 $\lambda = 1/T_A$: arrival
 $u = \lambda/\mu$: Utilization
 $= T_s/T_A$
 $C = \sigma^2/T_s^2$
 $= \text{Var}(X)/E[X]^2$
 $= 0$ if periodic deterministic
 $= 1$ if memoryless, poisson process
 $= 1/2$ in FIFO sys
 $= 1/2$ in M/M/1

$T_q = T_s \times \frac{C+1}{2} \times \frac{u}{1-u}$
 \hookrightarrow Avg time in queue
 $L_q = \lambda T_q$ Little's Law: Avg jobs in sys
 $L_q = \lambda \times T_q$
 $= \text{arrival rate} \times \text{response time}$
 \hookrightarrow Assn stable state
3 memoryless



Storage

HDD: Magnetic Disk Storage

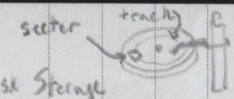
- Outside pln faster
- Unit of transfer is blk (i.e. sector)
- Overlapping tracks have greater density

Req. Time = $T_q + T_c + T_D$
 T_q : Device driver software queue
 T_c : Hardware controller
 T_D : Disk seek + rot. + transfer
 \hookrightarrow same track \Rightarrow no seek | $\frac{1}{2} \times \frac{1}{\text{rpm}}$
 Hardware checksum detects bad sectors
 Scheduling: FIFO has trouble seek times, ideally shortest seek and rot time (possibly leading to starvation)
 Ex: SCAN elevator algo to find closest request in direction of travel

SSD: Flash storage devices

- No moving parts; low power/light
- Operate on page-size chunks, can't address byte-level
- Sequential/random reads are same speed
- Writes limited in-order, must first erase numerous blks
 \hookrightarrow Controller manages free list
 \hookrightarrow Writes to free cells; erase entire writes
- Each blk has finite life of erases were logically managed by...
- Flash Translation Layer (FTL): maps os virt. blk \rightarrow phys pg #
 \hookrightarrow Can implement copy-on-write

Effective BW per op: $\frac{\text{transfer size}}{\text{response time}}$
 $x = \# \text{ ops}$
 $B = \text{time per op}$
 $S = \text{overhead}$
 $= \frac{x/S + x/B}{x/S + x/B}$
 $= B / (1 + SB/n)$



I/O Bus: Common set of wires for communication among devices, protocol for data transfer

- $O(n^2)$ wires for n devices, as $n \rightarrow \infty$ capacitance $\rightarrow \infty$ (slower)
- Communication monopolizes channel
- Ex: PCIe bus - Fast serial channels (parallel limited to slowest device)
 \hookrightarrow Account for vary access times/returns
- CPU \leftrightarrow Device controller communication
 \hookrightarrow Set of registers in device, possibly memory buffer
- Programmed IO: User cycles data size
1) Port-mapped IO: input/output
- 2) Memory-mapped IO: load/store insts
 \hookrightarrow Appear in physical addr space, can protect w/ translation
 \hookrightarrow Ex frame buffer / setting interrupt vector

Data Granularity

- Block: open/read/write/seek
 \hookrightarrow Disk drives, DVD ROM, raw I/O
- Character: get/put
 \hookrightarrow keyboard, mice, USB
- Network: socket/select
 \hookrightarrow Ethernet/WiFi, Bluetooth

Direct Memory Access [no CPU access audio]
Give controller access to mem. bus, let it transfer data blks directly
 \hookrightarrow Must signal/notify CPU when done
 \hookrightarrow Controller needs DRAM access
1) Interrupt: handles unexpected events well, but high overhead

- 2) Polling: OS periodically checks status register
lower overhead, but can waste cycles
 \hookrightarrow Use in kernel: inter. First packet, then poll rest

For doubly indirect:

$$\text{max-File} = \text{blk-size} \times \left[\text{indir.} + \frac{\text{blk-size}}{u} \times (u^2) \right]$$

For doubly indirect:
max-File = blk-size * [indir. + (blk-size/u) * (u^2)]