

- OS provides abstractions from hardware, makes it easier to reference
- Process
  - address space
  - its threads (id)
  - other (files, sockets, etc)
- 4 OS Concepts
  - Thread of Control
    - lots of registers, executing when present
    - suspended when not loaded
    - illusion of multicore
    - TCB - thread info, internal for now
  - Address Space
    - address state, not only files, maybe IO, etc
    - need some area of memory allocated to exec
    - ex. Base and Bound, translation, page memory
  - Process
    - execution environment w/ its threads
    - partitioned from each other and from OS
    - threads = concurrency
  - Privilege (Dual Mode Operation)
    - 2 modes: kernel, user
    - Kernel mode handle - syscall, interrupt, trap/exception
      - kernel specifies how to handle interrupt, ex vector table
    - PCB - state, registers, etc maintained by kernel
    - scheduler picks which to run
- Thread Abstraction
  - concurrency vs parallelism, OS handles multiple things at once
  - multiplexing = multithread
  - multiprogramming = multi-threading, concurrency by 2!
  - states: RUNNING, READY, BLOCKED
  - IO - sleep, to block thread until done
  - unique TCB per stack, but shared global in process
  - non-determinism - remembers my time
  - race condition - depends on interleaving
  - synchronization
    - mutual exclusion - only one thread doing thing
    - condition variable - only one thread can execute at once
    - lock - only one thread holds, others wait
  - Process
    - even shell is only process running after process
    - bootstrapping - kernel starts init, which starts other processes
    - syscalls - exit, fork, exec, wait, kill, signal
  - File / IO
    - Semaphore - generalized lock
      - unregistered int, about to p
      - PC == down()
      - VC() == up()
    - everything is a file
    - file & directory
      - each process has current working directory (CWD)
      - streams FILE\*, internal buffer, high level
    - file descriptor fd, low level
      - buffer in kernel
      - raw or block I/O
    - ioctl, dup, fprint
    - FILE\* internal buffer, no file buffer on disk
      - basically faster to minimize syscalls
    - open file descriptor maintained in kernel
      - file descriptor returned
      - forking copies the descriptor but both refer to same kernel file descriptor
      - have we can do IPC, pipes
      - never fork() in multithreaded, only 1 thread copied
      - fork() copies everything
      - exec() safe to replace, future address / open
      - don't fork() w/ FILE\* bc don't know if flushed
  - IPC, Pipes, Sockets
    - common form is File I/O
    - using actual file is slow, so pipes! (1-way)
      - EOF when last process holding fd is closed
    - also read protocol, usually part of API functionality
    - TCP Network Protocol, Sockets
      - 2 bytes
        - thread pool to serve clients
      - send socket: listen(), accept()
        - connection
    - Simplex
      - Source IP: source port - usually well known
      - Dest IP: client port usually random
      - Source Port
      - Dest Port
      - Protocol

- Synchronization
  - process life states - new, ready, running, waiting, terminated
  - multiple scheduler queues like for ready, disk, etc
  - scheduler / dispatcher gets control w/ preemption or voluntary
  - context switch - new stuff, move over
    - 10-100ns frequency
    - process switch takes 3-4 usec
    - thread switch takes 100 nsec
    - vs LLC in memory even cheaper
  - Thread Pool - lots of threads housekeeping
    - shared data corruption - locks bounded
      - also don't want lockstep, use buffers, ring
      - circular buffer - need
  - Lock Implementation
    - interrupt enable/disable - native hardware
    - atomic load/store
    - sleep doesn't need to unlock
  - Atomic Read Modify Write Instruction
    - test&set() - set bit and return old value
    - swap - swap, x86
    - compare&set - if true, set register
    - load-linked/compare-and-set
  - Busy waiting - does nothing
    - look for cache, CPU usage
    - priority inversion
    - better way to sleep()
  - Futex - fast user-space mutex
    - no syscalls, used in pthreads
  - Monitor - lock + condition variable
    - cv - list of threads waiting
    - writer - advance checkpoint, sleep, acquire
    - signal() - wake one
    - broadcast() - wake all
  - Hand semaphores - immediate transfer
  - More popular: reader-writer lock, writer on ready queue
    - spin - busy loop reacquiring (non-blocking), need
      - spin while loop to check spin
  - kernel thread we can use threads
  - Page Table for process
  - Device Driver
    - top half: collecting interrupt
    - bottom half: interrupt function
  - x86 Assembly (32-bit)
    - Registers: eax, esp, ebp, etc
    - AT&T Syntax
      - register %eax
      - immediate: \$X
      - pBuffer : (%eax)
      - pBuffer : %eax
    - int source destination
      - esp/16(%eax)
    - Caller
      - args on stack, even %ebp, 16 bytes total
      - push %ax
      - restore stack
    - C-Hat
      - b16(%esp), esp after esp
      - return value stored
      - restore esp
      - return
    - parallel add, sub, mov
      - push %eax, call, leave, ret, parallel, pop %eax
  - Extra
    - security ops in user space

Matthew Tran  
CS162  
MTI

- Scheduling
    - A preemptive
    - 1 program/user
    - 1 thread/thread
    - independent prog
    - execution model
    - CPU bursts then IDL
  - Goals
    - minimize responsiveness
    - maximize throughput
    - Fairness - static vs. dynamic
  - FCFS - first come first serve, preemptible
  - Round robin
    - preemption = each process gets "time quota"
    - optimal & - big = FCFS, smaller becomes round robin
    - response time = time of finish - time queued
    - 10-15ms today, 0.1ms for cache switch
  - First Priority Scheduling
    - starvation - some don't run
    - priority inversion - highest priority need for lower priority
  - SJF/SCTF (Shortest Job First)
  - ERTF (Shortest Remaining Time First)
    - preemptive version of SJF
    - SJF/EFT both at minimizing average response time
    - cost/priority
      - cost/priority has long process takes / real world
      - adaptive - priority fluctuates in part
  - Lottery scheduling
    - each process gets 1/lot, randomly pick one
  - MLFQ (multi-level linked list scheduling)
    - multiple queue we diff priority w/ diff scheduler (RR, FCFS)
    - it's linearly expand, max down, if not move up
    - It's still single hps, CPU or batches, like ERTF
    - counter measure - us priorities to keep priority high
  - Multithreaded scheduling
    - very similar to rounds, but maybe per-core scheduler data
    - cache - want to reassign tasks to same core
    - cache coherence big issue
    - sporadic, busy wait on sometimes, more efficient when gang scheduling/pooling qps - want to schedule multi-thread together
    - want threads in linear perf scenario
- 
- Demand Driven - user's demand goes by half  $\rightarrow$  job backlog  $\rightarrow$  HWM
- Linear O(1) scheduler
  - 140 priorities, 100 "idle", 100 kernel/RT
  - lots of admin overhead, too complex
- Round Robin scheduling
  - Predictability 80% NOT PERFORMANCE 80%
  - Hard RT - meet all deadlines (if possible)
    - earliest deadline first (EDF)
    - Least-Likely-First (LLF)
    - Round-Robin Scheduling (RRS)
    - Round-Robin Scheduling (RRM)
  - Soft RT - attempt to reach deadlines w/ high priority
  - earliest deadline first (EDF)
    - tasks priorities w/ period P and computation C in each period
    - feasible if  $\frac{C}{P} \leq 1$
  - work-conserving scheduler - if work left, CPU never idle
  - Starvation - task intervals run but can't make right execution!
    - FCFS - nonpreemptive, no possible
    - RR - no, but thought fails
    - Priority Inversion - if stuck, ready, but always gets interrupted!
    - priority inversion!
  - SJF/MLFQ - constraints to start one first
- Proportional Share Scheduling
  - Idea: share CPU, but don't strict, "proportionally"
  - no starvation, but remember
- Shale Scheduling
  - various "fairness" numbers
  - "shale" =  $\frac{1}{n}$ , max load 1/n = interleaved
  - each "pair" counter, priorities, runs, add shales per core
    - can be many
- Linux Completely Fair Scheduler (CFS)
  - currently used, goal to perfectly divide
  - O(log n) heap, randomized w/ start time
  - good for interactive
  - target latency =  $\text{latency} = \text{time} / n + \sqrt{\text{variance}} / n$
  - throughput - want overhead, min time slice
- Deadline
  - circular dependency
  - no fair - fine resource acquire order
  - Requests
    1. initialization
    2. Hit and wait
    3. no preemption (voluntary release)
    4. circular wait

- Virtual Memory
  - memory multiplexing
  - protection
  - translation
  - centralized memory
  - Bus & Direct
    - add bus, check board
  - Memory Translation
    - fragmentation, space leak, sharing issues
  - Segmentation
    - diff regions have separate regions, anywhere physically
    - segmentation in processor  $\rightarrow$  page limit
    - entry per chunk
    - ex: 96 no entries for this, but 16 bytes will give 16 pages
    - many 16 bytes = variable chunks, no wins, fragmentation
  - Paging - fixed?
    - PTE - ppn + bbb (mult, msk, etc.)
    - thrashing? - thrash? PTE switch to same page, read diff?
    - kernel PT and PTE, contexts overspace
    - multi-level PT
      - simple talk for long - communicate DT
      - ex: 16-12-12 (4K/2K/1K) (16MB)
      - higher levels can talk on to disk too
    - PTE - problem with PT or page + permission
    - > 8K PTE - PPN (2^14)
    - P: present/valid
    - W: writeable
    - U: user accessible
    - PWT: page-walk transparent
    - PCD: page-change detect/flush
    - As accessed
    - D: dirty (PT entry), if physical PT
    - PS: exclusive (directory entry)
    - Multilevel Translation: Segments / Pages / Address
      - segment - memory w/ width 16 (block limit)
      - page - equal size blocks - same memory area!
      - msg  $\rightarrow$  PT + PPN
      - context switch - save top level segment / PT
      - 8K segment descriptors
      - 6 usage register = SE, DS, CS, ES, SS, GS
        - pointer to segment descriptor
        - G/L bit with GOT/DOT (global/local descriptor table)
        - RPL - regular privilege level
        - GDT/RDT/ldt pointers to them below
        - G1 granularity of segment (16B, 4KB)
        - DB: shadowed size (1GB, 32B)
      - As anything
      - P: Accessed/valid
      - DPL: descriptor privilege level
      - S: virtual segment (r/w, o/d, width)
      - Type: code, data, segment
      - 1 GOT for all, diff of LDT for components
        - usually pointers for code segments - legacy
        - modern - flat address, 1 segment 4GB
        - 16 bits in GOT + 16 bits in Thread Local Scope (TLS)
        - return to stack - save base, save TLS
      - Invited Page Table
        - host table, guest table, guest for each bit
        - pre-allocate, copy to host table
      - MMU - does all this translation
      - TLB: cache for translation: Vpn + PPN
        - PT might be in normal cache
        - even multilevel
        - dirty entries, we don't map in first
        - multiply TLB?
        - context switch? invalid TLB?
    - Deadlock detection - Resource Allocation graph
      - a Directed, m-Partite graph
        - two instances of each
      - request(s), use(s), release(s)
      - request edge: T  $\rightarrow$  R
      - assignment edge: R  $\rightarrow$  T
      - alg (cycle (deadlock) detection
        - add all nodes to uninitial
        - executable resource array
        - loop thru initial
          - if node requires executable
          - remove
          - add requires to available
          - no deadlock
        - logically needs all nodes can initially require
    - Deadlock deal with
      1. partition
      2. recovery
      3. avoidance
      4. denial
      - O.S. prevents own deadlock handles!
      - Context Algorithm
        - only good if no deadlock (deadlocks)
        - max - alias conflict of request

- Matthew Tran  
C1162  
MT2
- Cache
    - L1, L2, L3
    - compulsory - most
    - capacity - how much
    - conflict - 2 may access
    - coherence - other update
  - Address
    - blocks - offset
    - index - which block
    - tag - ID
  - Types
    - Direct Mapped - 1 address per block
    - Set Association - lots of blocks
    - Fully Associative - looks, read tag
  - Replacement
    - Random
    - LRU
    - Worst
    - write through - d. both
    - write back - cache writes, but dirty
  - Physically-Addressed vs. Virtual-Address Cache
    - VA  $\rightarrow$  VA + TLB  $\rightarrow$  Cache
    - VA + Cache (bad flesh on switch?)
    - overlap TLB and cache
      - cache by tag = VA, offset
      - virtual cache faster?
  - Demand Paging
    - page fault - no translation fault/trap
    - main memory is "cache" for disks
    - block size = 1 page
    - Fully associative
    - replacement policy?
    - write back to disk if not modified
    - cache real copy pages
    - also mmap?
    - O.S. maintains blocks in disk
    - working set model
      - pages in certain sections at a time
  - Replacement Policy
    - compulsion, capacity, conflict, aging, min
    - FIFO - fair, static
    - RANDOM
    - MRU - probably optimal, replace page w/ most used in longest time
    - LRU - least recently used, typical of MSW
    - Belady Anomaly - from Euclidean
      - affects FIFO
    - Clock Algorithm
      - typical of LRU - replace old, not oldest
      - all pages in table, have use bit
      - new gets old bit, youngest
      - on fault - check bit, if 1 - increment
      - on replace
        - in next iteration
    - NRU Cache algorithm
      - clock bit is cleaner
      - if counter = 0, replace; no swap
      - give dirty page extra chance?
    - Clock Alg Variation
      - no "modified" bit
        - use "used only" bit
        - use "old" bit
      - if 0, use, otherwise write, no use modify
      - on "use" to 1
        - old available
        - mark all page invalid first
    - Second Chance Alg (VAX/VMS)
      - actual list, SC list
      - PES
        - LRU
        - SL (shortest page first)
        - move to back of queue
        - move active to SC
        - return PC, LRU and free, go in
      - Free List
        - similar to SC
        - filled by diff alg like clock
    - Reverse Dom Mapping ("fixing")
      - which PTE to invalidate?
      - known object based reverse mapping
    - Page Frame Allocation
      - each process real minimum
      - global/local placement?
      - equal/proportional/priority allocation
      - Consistency? being invalid/legit?
        - track working set?
    - Multibyte - explicit pipeline to next kernel mem
      - patches - different bits for different kernel
      - flush TLB twice on sync call, then

- I/O
  - Block, character, network devices
    - set/pkt swap
  - blocking, nonblocking, synchronous interface
  - HDD
    - surface ring (track) separated by head regions
    - outside reading tracks have more sectors
    - SIRL - overlapping tracks = T density/capacity
      - bit width header, DIF for reading
    - cylinder = all tracks under head at point (multiplexed)
    - latency = queue time + controller + seek + rotation + transfer
      - seek = move head over track  $\sim 4-6ms$
      - rotational = wait for sector  $\sim 4-10ms$  (Latency)
      - transfer = transfer bits  $\sim 0.0001-0.0002ms$
      - total  $\sim 10ms$
    - SSD
      - latency = queue time + controller + transfer
        - write only  $\sim 25-60ms$  at a time, slow, limited!
      - FTL - Flash Translation Layer, like VMM
        - don't overallocate page, write to new and indirect
        - wear leveling, garbage collection
      - MLC/MLC code module, DRAM speed is even
      - Performance measure metrics
        - latency - time to complete
        - response time - time to start & respond
        - throughput/track width - rate
        - storage overhead - time to initialize
    - File system
      - LRU(Least Recently Used) system in use
      - $N_{\text{jobs}} = \lambda \cdot T_{\text{idle}} + L_{\text{latency}}$
      - utilization  $\rho = \frac{\lambda}{T_{\text{idle}}} \cdot \min(\text{max job length}, \text{idle time})$ 
        - $\Delta T_{\text{idle}} = \text{queue time} + \text{idle time}$
        - $L_{\text{latency}} = \text{idle time}$
      - queuing theory
        - $C = \frac{\sum_{i=1}^n \rho_i}{m^2} \cdot \text{avg. determination}$
        - memory  $T_{\text{idle}} = \frac{\rho}{1-\rho} \cdot T_s$
        - general  $T_{\text{idle}} = \frac{1-\rho}{2} \cdot \frac{\rho}{1-\rho} \cdot T_s$
        - utilization  $L_{\text{idle}} = \lambda T_{\text{idle}} = \frac{\rho}{T_s} \cdot T_{\text{idle}} = \frac{\rho^2}{1-\rho} \cdot T_s$
      - Disk Scheduling
        - FIFO - using next head
        - SSTF - possible starvation
        - SCAN - elevator, closest request to direction, C track
        - C-SCAN - elevator, skips in same track, no middle track
      - block  $\rightarrow$  file/directory
        - also look free blocks, idle blocks, directory
        - expand file size
        - open file with file pointer (head/tail/pointer)
        - fd - directory, index, header, always write from beginning
        - read/write raw bytes w/ directory
        - cwd - other relative path
        - write - file info on disk, don't copy in memory
      - FAT (File Allocation Table)
        - indexed link of blocks
          - directory has file
          - file number = index in table
          - metadata?
          - root / or block 2
      - Berkeley FFS (Unix)
        - file number = index in disk array
        - block/file holds metadata?
        - direct, indirect pointers
        - core - index all 1 place, having file  $\rightarrow$  of FAT
        - locality - block groups, long blocks close (HDD)
        - core additional layer
          - free block pointers or read ahead
          - core - lots of cache for large files
      - Ext2/3 - Linux
        - block groups - locality, do lots free space
        - GFS has journaling
      - Hard links - rename file name
        - multiple references count to keep writers deleted
      - Symbolic - link to file names
      - B-trees for large directories
      - Windows NTFS
        - variable length entries
        - master file table - like FAT on earlier
          - metadata & controls if small
          - 1KB entries
        - memory file
          - lots + sharing?
        - log file cache - for files, in software
          - replaced policy? LRU?
          - read ahead?
          - evict heuristic - remove file
    - File system relationship/recovery
      - availability, durability, reliability
      - RAID
        - RAID 1 - full copy
        - RAID 5 - striping, ECC
        - RAID 6 - 2 sub disks
        - RAID branching?
          - careful ordering - FAT, FFS
            - operating in memory
            - copying in memory
            - copy as write (COW)
            - FFS, write
            - make new file, move over
          - atom transaction
            - commit changes
        - Journaling
          - update in log, apply later
          - NFS, AMFS, XFS, ZFS, ext4
          - after commit, complete log apply
            - recover different portion
            - exposure - usually only metadata update
        - Log FS (LFS)
          - log is the storage
            - all sequential, wrap around
            - buffer cache makes fast
            - segment summary to read if needed
            - free block - what will release, read & record (replay)
        - FFS - flat file system
          - for SSDs - reduce sequential reads
            - random writes bad (capital R mark)
            - page addressable (read & write) instead of FTL
        - Distributed Systems
          - IP - and network layer protocols
        - Distributed Data - Moving
          - General's Problem
            - contiguous & continuous
            - Two-phase Commit (2PC)
              - pessimistic log (optimistic)
              - pre-prepare, commit phase
              - all or nothing
              - no half guarantee
              - two blocking
              - after - 3PC, Paxos, RAFT
            - Byzantine General Problem
              - malicious nodes
              - Paxos, P2P
              - Sliders
                - application
                - Transport - TCP/UDP
                - Network - IP
                - Database
                - physical
              - TCP
                - 4 quad = 2 per
                - window size = remaining space
                - never read more than can receive
                - ACK no. & length  $\uparrow$ 
                  - after ack.
                - 3 way handshake
                - 4 way termination
              - RPC
                - Virtual FS (VFS)
                  - expandable, portable, directory, file
                  - shadow - all info need to be intact
              - NFS
                - UNIX, VMS, Netware
                - write through
                - shadow, intercept
              - AFS
                - Network
                  - write through cache
                - key value store  $\uparrow$
                - content-based hashing
                  - deal peer-to-peer machine
                  - ring, shadow node w/ random ID
                  - local directory

$$\begin{aligned}\lambda &= \frac{L_{\text{idle}}}{T_{\text{idle}}} = \frac{0.75}{0.001} \\ T_{\text{idle}} &= \frac{L_{\text{idle}}}{\lambda} = \frac{0.75}{0.75} \\ L_{\text{idle}} &= \lambda \cdot T_{\text{idle}} = 1.25 \cdot \frac{0.75}{1-0.75} \cdot T_s \\ L_{\text{idle}} &= \lambda (1.25 T_s - \frac{\lambda}{1-\lambda})\end{aligned}$$

$$\begin{aligned}L_{\text{idle}} &= L_{\text{idle}}^2 = 1.25 T_s \lambda^2 \\ L_{\text{idle}}^2 - L_{\text{idle}} \lambda &= 1.25 T_s \lambda^2 \\ 1.25 T_s \lambda^2 + L_{\text{idle}} \lambda - L_{\text{idle}}^2 &= 0\end{aligned}$$