

CS162
Operating Systems and
Systems Programming
Lecture 18

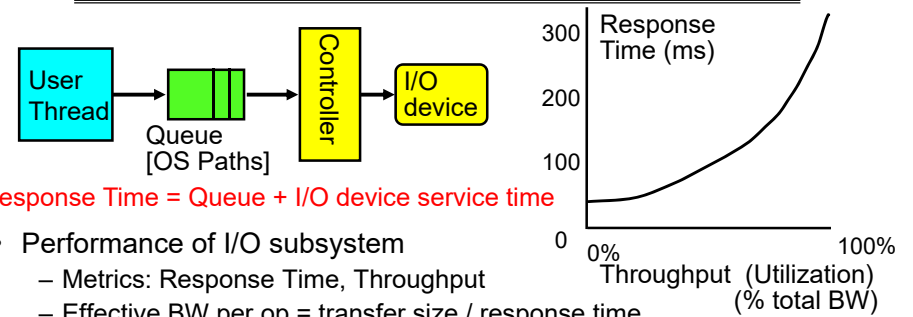
Queueing Theory,
Disk scheduling & File Systems

April 9th, 2019

Prof. John Kubiatowicz

<http://cs162.eecs.Berkeley.edu>

Recall: I/O Performance



Response Time = Queue + I/O device service time

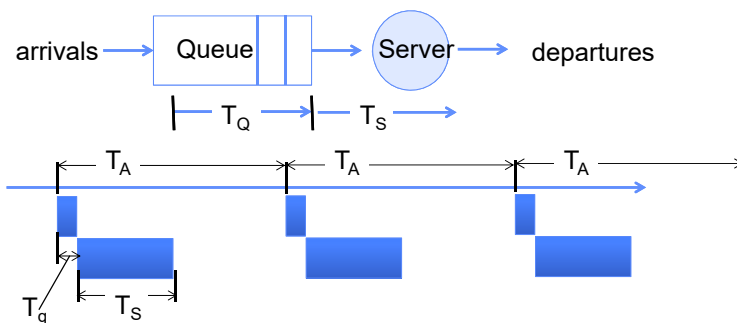
- Performance of I/O subsystem
 - Metrics: Response Time, Throughput
 - Effective BW per op = transfer size / response time
 - » $\text{EffBW}(n) = n / (S + n/B) = B / (1 + SB/n)$
 - Contributing factors to latency:
 - » Software paths (can be loosely modeled by a queue)
 - » Hardware controller
 - » I/O device service time
- Queuing behavior:
 - Can lead to big increases of latency as utilization increases
 - Solutions?

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.2

A Simple Deterministic World



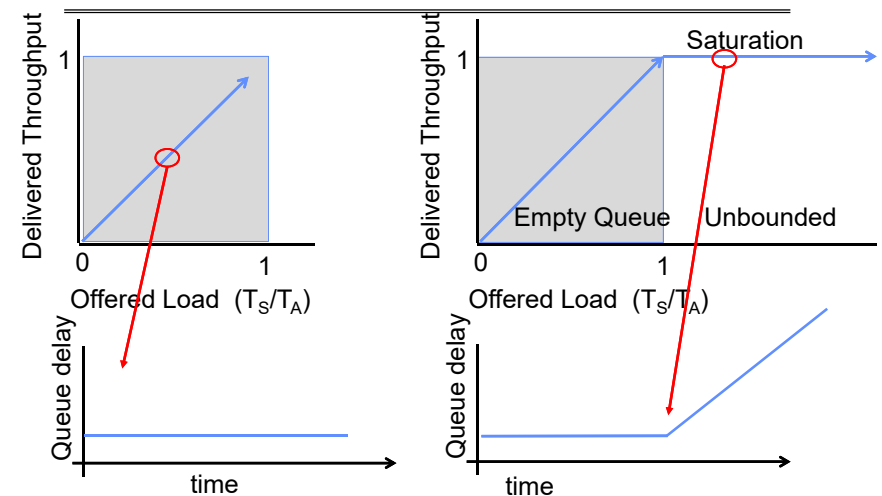
- Assume requests arrive at regular intervals, take a fixed time to process, with plenty of time between ...
- Service rate ($\mu = 1/T_S$) - operations per second
- Arrival rate: ($\lambda = 1/T_A$) - requests per second
- Utilization: $U = \lambda/\mu$, where $\lambda < \mu$
- Average rate is the complete story

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.3

A Ideal Linear World



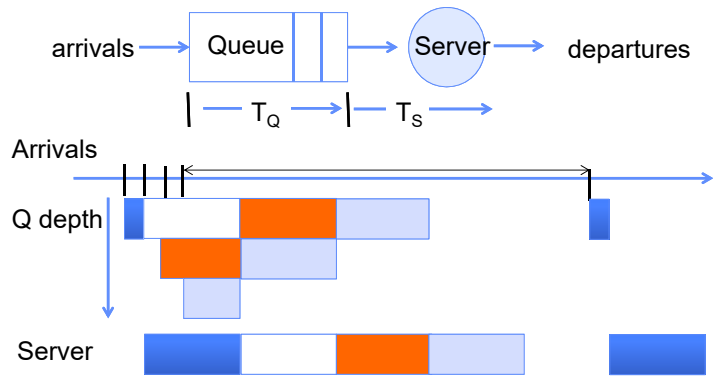
- What does the queue wait time look like during overload?
 - Grows unbounded at a rate $\sim (T_S/T_A)$ till request rate subsides

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.4

Reality: A Bursty World



- Requests arrive in a burst, must queue up till served
- Same average arrival time, but:
 - Almost all of the requests experience large queue delays
 - Even though average utilization is low!

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

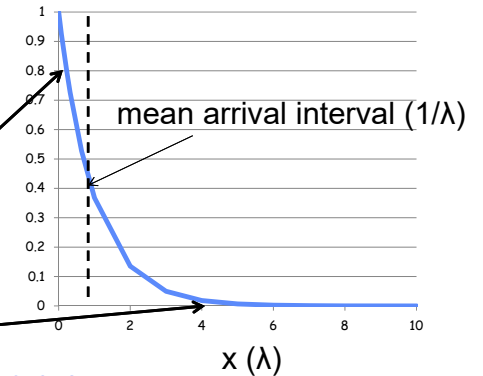
Lec 18.5

So how do we model the burstiness of arrival?

- Elegant mathematical framework if you start with *exponential distribution*
 - Probability density function of a continuous random variable with a mean of $1/\lambda$
 - $f(x) = \lambda e^{-\lambda x}$
 - “Memoryless”

Likelihood of an event occurring is independent of how long we've been waiting

Lots of short arrival intervals (i.e., high instantaneous rate)
Few long gaps (i.e., low instantaneous rate)



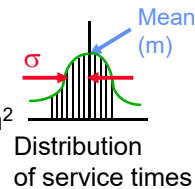
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

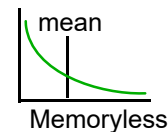
Lec 18.6

Background: General Use of Random Distributions

- Server spends variable time (T) with customers
 - Mean (Average) $m = \sum p(T) \times T$
 - Variance (stddev²) $\sigma^2 = \sum p(T) \times (T-m)^2 = \sum p(T) \times T^2 - m^2$
 - Squared coefficient of variance: $C = \sigma^2/m^2$



- Important values of C :
 - No variance or deterministic $\Rightarrow C=0$
 - “Memoryless” or exponential $\Rightarrow C=1$
 - » Past tells nothing about future
 - » Poisson process – *purely* or *completely* random process
 - » Many complex systems (or aggregates) are well described as memoryless
 - Disk response times $C \approx 1.5$ (majority seeks < average)

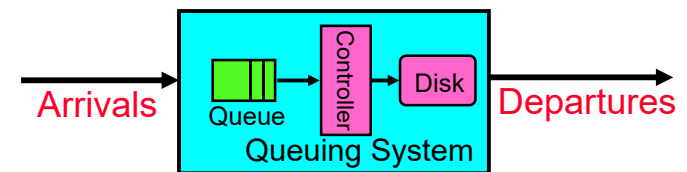


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.7

Introduction to Queuing Theory



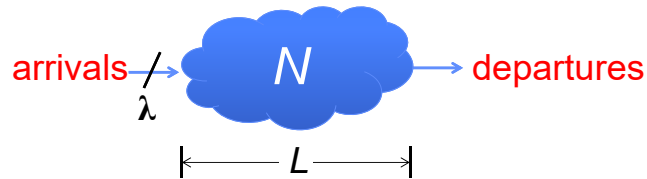
- What about queuing time??
 - Let's apply some queuing theory
 - Queuing Theory applies to long term, steady state behavior \Rightarrow Arrival rate = Departure rate
- Arrivals characterized by some probabilistic distribution
- Departures characterized by some probabilistic distribution

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.8

Little's Law



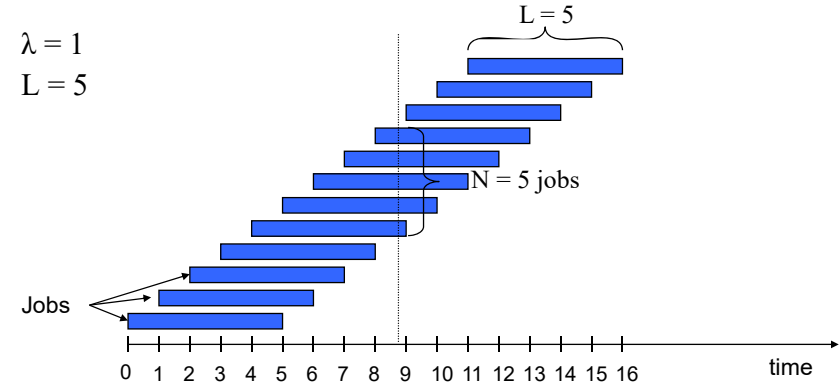
- In any *stable* system
 - Average arrival rate = Average departure rate
- The average number of jobs/tasks in the system (N) is equal to arrival time / throughput (λ) times the response time (L)
 - $N \text{ (jobs)} = \lambda \text{ (jobs/s)} \times L \text{ (s)}$
- Regardless of structure, bursts of requests, variation in service
 - Instantaneous variations, but it washes out in the average
 - Overall, requests match departures

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.9

Example



A: $N = \lambda \times L$

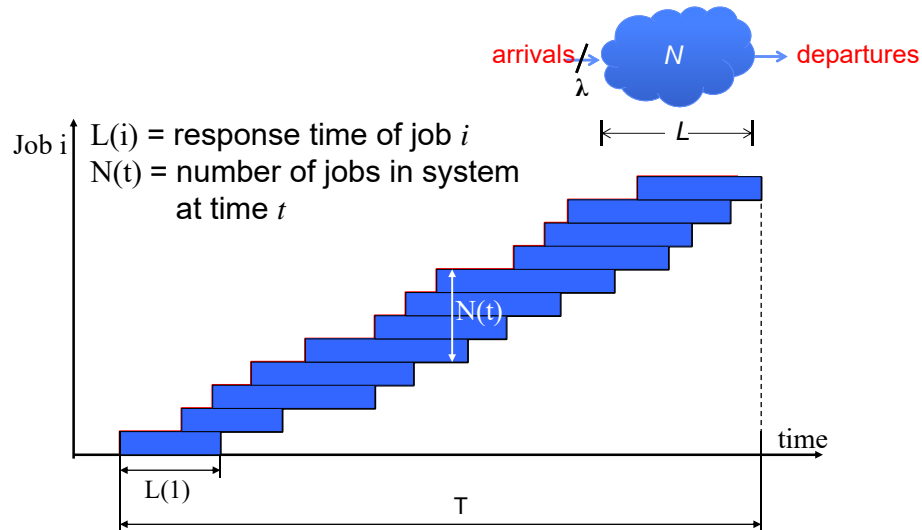
- E.g., $N = \lambda \times L = 5$

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.10

Little's Theorem: Proof Sketch

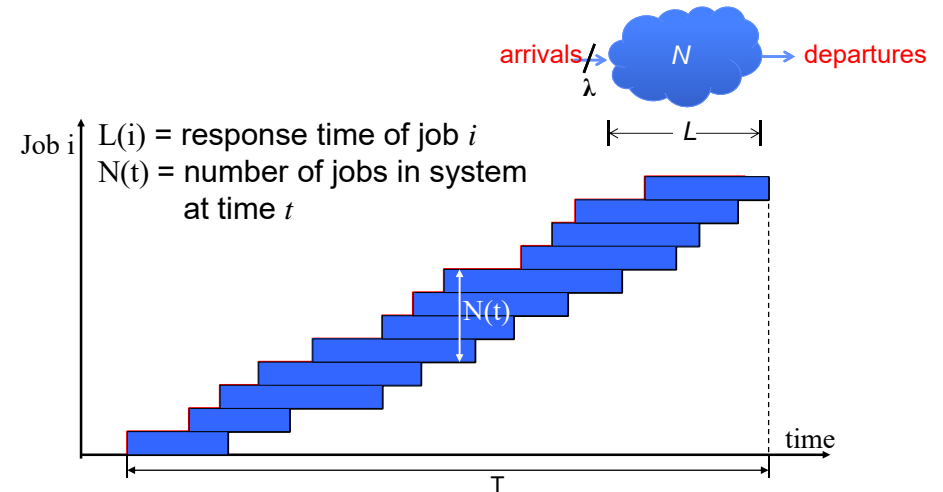


4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.11

Little's Theorem: Proof Sketch



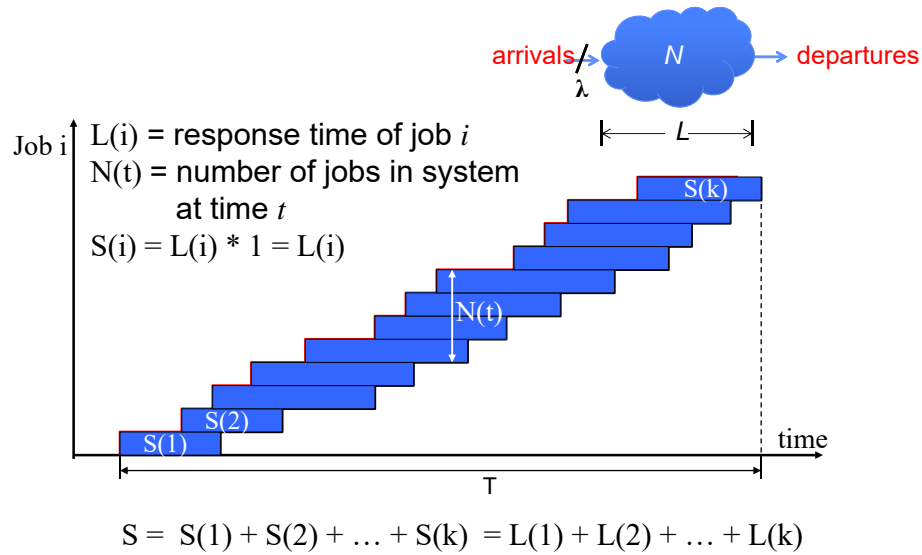
What is the system occupancy, i.e., average number of jobs in the system?

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.12

Little's Theorem: Proof Sketch

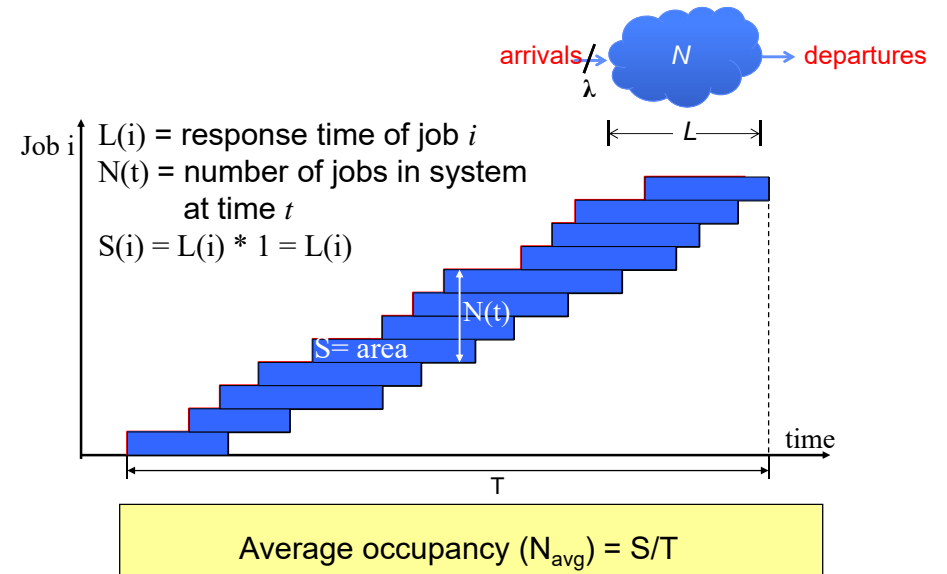


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.13

Little's Theorem: Proof Sketch

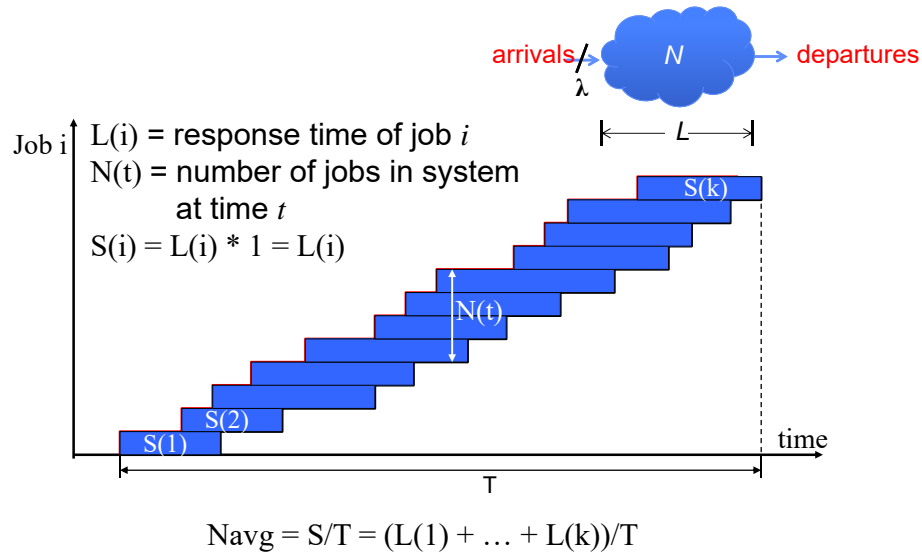


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.14

Little's Theorem: Proof Sketch

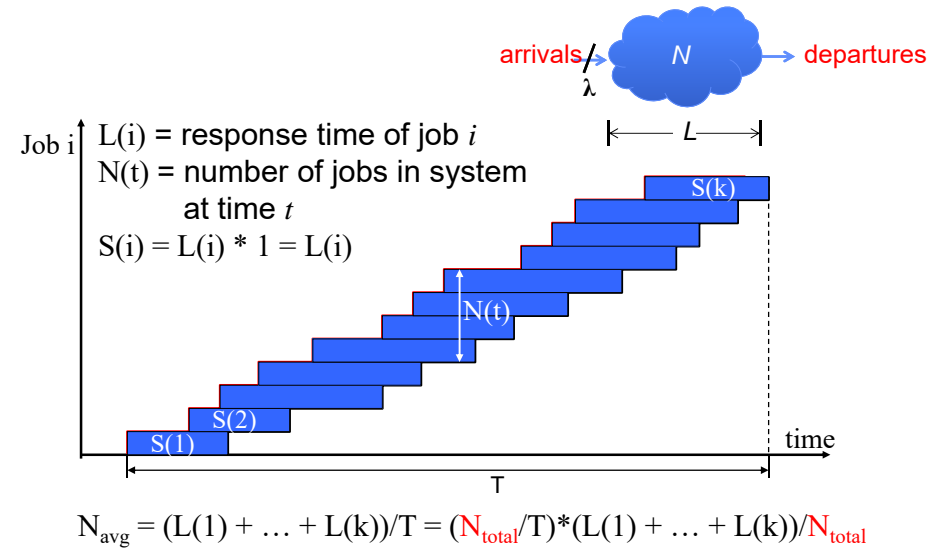


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.15

Little's Theorem: Proof Sketch

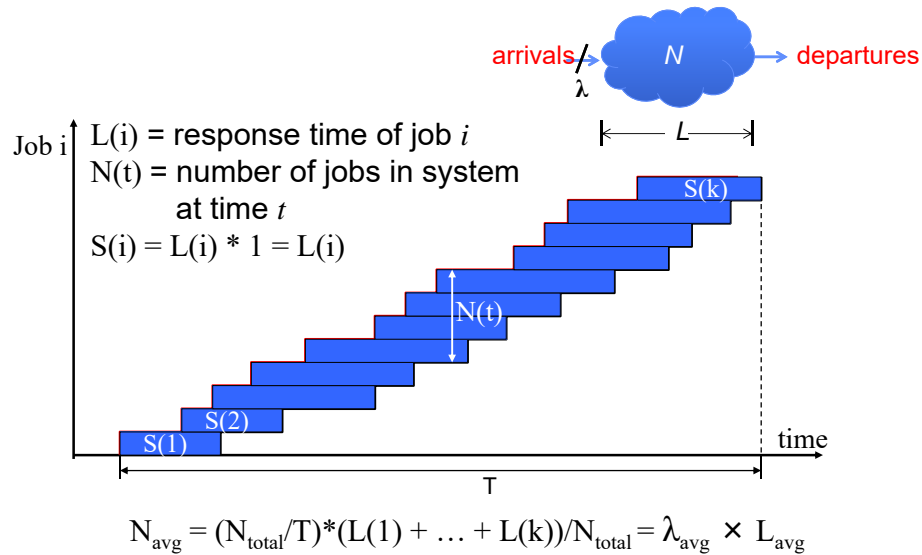


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.16

Little's Theorem: Proof Sketch

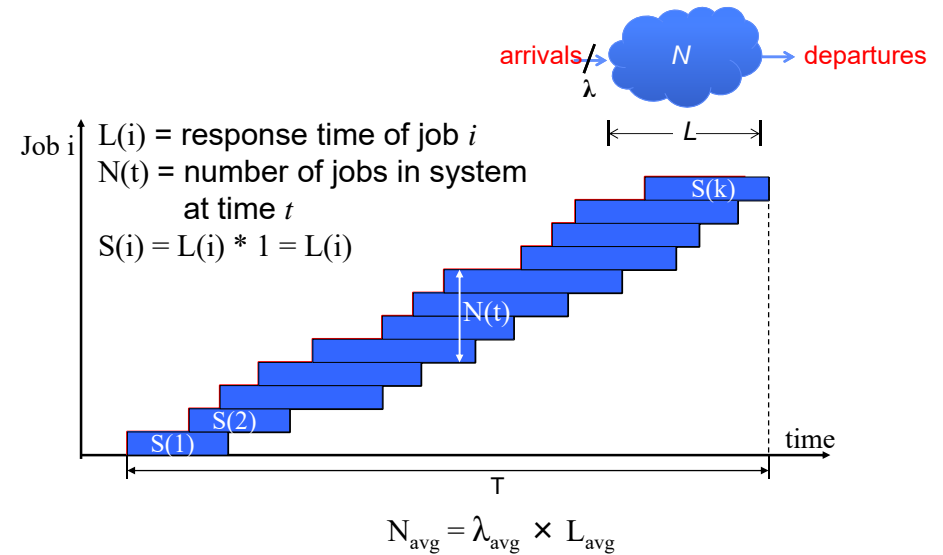


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.17

Little's Theorem: Proof Sketch

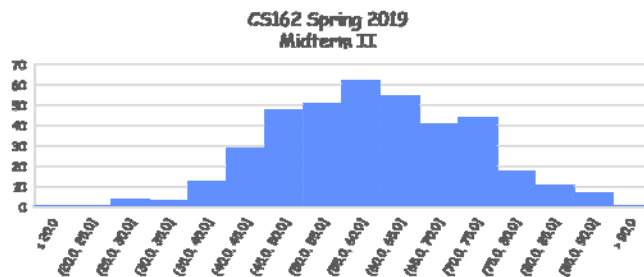


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.18

Administrivia



- Midterm II: Too long! Sorry!
 - Yup, we misjudged that one
- Midterm II Statistics:
 - Mean: 58.9, STD: 12.7, Max: 94.0
- Solutions are up
 - Regrade requests close on Friday 4/12, 11:59pm

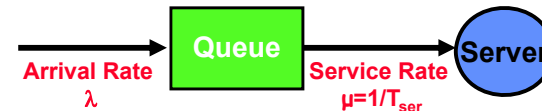
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.19

A Little Queuing Theory: Some Results

- Assumptions:
 - System in equilibrium; No limit to the queue
 - Time between successive arrivals is random and memoryless



- Parameters that describe our system:
 - λ : mean number of arriving customers/second
 - T_{ser} : mean time to service a customer ("m1")
 - C : squared coefficient of variance = σ^2/m^2
 - μ : service rate = $1/T_{\text{ser}}$
 - u : server utilization ($0 \leq u \leq 1$): $u = \lambda/\mu = \lambda \times T_{\text{ser}}$
- Parameters we wish to compute:
 - T_q : Time spent in queue
 - L_q : Length of queue = $\lambda \times T_q$ (by Little's law)
- Results:
 - Memoryless service distribution ($C = 1$): (an "M/M/1 queue"):
 - » $T_q = T_{\text{ser}} \times u / (1 - u)$
 - General service distribution (no restrictions), 1 server (an "M/G/1 queue"):
 - » $T_q = T_{\text{ser}} \times \frac{1}{2}(1+C) \times u / (1 - u)$

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

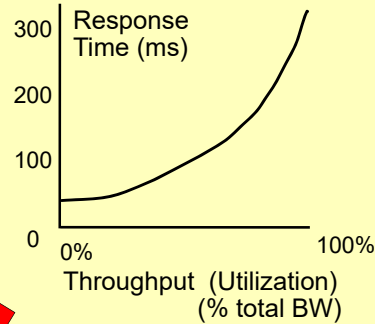
Lec 18.20

A Little Queuing Theory: Some Results

- Assumptions:
 - System in equilibrium; No limit to the number of customers in the system
 - Time between successive arrivals is independent of the number of customers in the system



Why does response/queueing delay grow unboundedly even though the utilization is < 1 ?



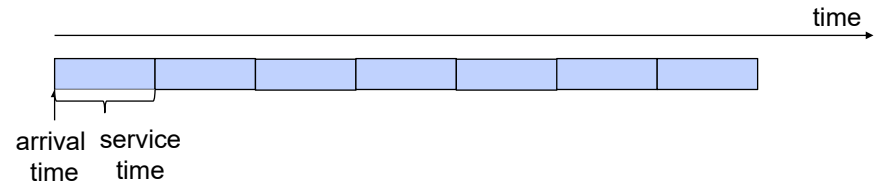
- Parameters that describe our system:
 - λ : mean number of arriving customers per second
 - T_{ser} : mean time to service a customer
 - C : squared coefficient of variation of service times
 - μ : service rate = $1/T_{ser}$
 - u : server utilization ($0 \leq u \leq 1$): $u = \lambda T_{ser}$

- Parameters we wish to compute:
 - T_q : Time spent in queue
 - L_q : Length of queue = $\lambda \times T_q$

- Results:
 - Memoryless service distribution ($C=1$): (an "M/M/1 queue"):
 - $T_q = T_{ser} \times u / (1 - u)$
 - General service distribution (no restrictions), 1 server (an "M/G/1 queue"):
 - $T_q = T_{ser} \times \frac{1}{2}(1 + C) \times u / (1 - u)$

Why unbounded response time?

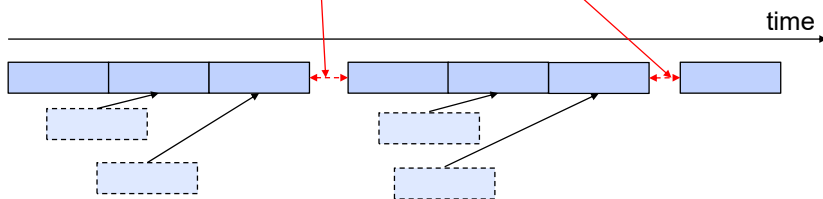
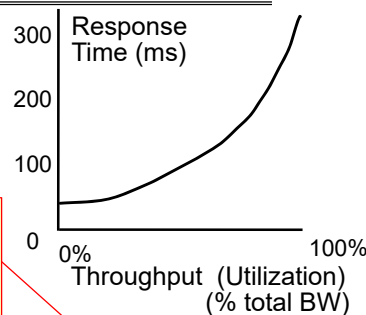
- Assume deterministic arrival process and service time
 - Possible to sustain utilization = 1 with bounded response time!



Why unbounded response time?

- Assume stochastic arrival process (and service time)
 - No longer possible to achieve utilization = 1

This wasted time can never be reclaimed! So cannot achieve $u = 1$!



A Little Queuing Theory: An Example

- Example Usage Statistics:
 - User requests 10 x 8KB disk I/Os per second
 - Requests & service exponentially distributed ($C=1.0$)
 - Avg. service = 20 ms (From controller+seek+rot+trans)
- Questions:
 - How utilized is the disk?
 - Ans: server utilization, $u = \lambda T_{ser}$
 - What is the average time spent in the queue?
 - Ans: T_q
 - What is the number of requests in the queue?
 - Ans: L_q
 - What is the avg response time for disk request?
 - Ans: $T_{sys} = T_q + T_{ser}$

- Computation:
 - λ (avg # arriving customers/s) = 10/s
 - T_{ser} (avg time to service customer) = 20 ms (0.02s)
 - u (server utilization) = $\lambda \times T_{ser} = 10/s \times .02s = 0.2$
 - T_q (avg time/customer in queue) = $T_{ser} \times u / (1 - u) = 20 \times 0.2 / (1 - 0.2) = 20 \times 0.25 = 5 \text{ ms} (0.005s)$
 - L_q (avg length of queue) = $\lambda \times T_q = 10/s \times .005s = 0.05$
 - T_{sys} (avg time/customer in system) = $T_q + T_{ser} = 25 \text{ ms}$

Queuing Theory Resources

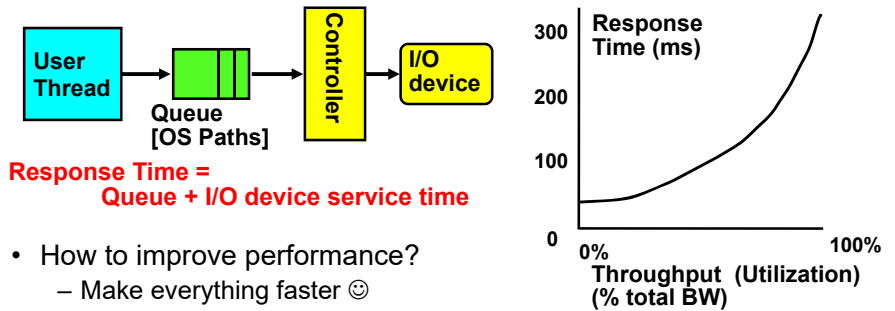
- Resources page contains Queuing Theory Resources (under Readings):
 - Scanned pages from Patterson and Hennessy book that gives further discussion and simple proof for general equation: https://cs162.eecs.berkeley.edu/static/readings/patterson_queue.pdf
 - A complete website full of resources: <http://web2.uwindsor.ca/math/hlynka/qonline.html>
- Some previous midterms with queuing theory questions
- Assume that Queuing Theory is fair game for Midterm III!

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.25

Optimize I/O Performance



- How to improve performance?
 - Make everything faster ☺
 - More Decoupled (Parallelism) systems
 - multiple independent buses or controllers
 - Optimize the bottleneck to increase service rate
 - Use the queue to optimize the service
 - Do other useful work while waiting
- Queues absorb bursts and smooth the flow
- Admissions control (finite queues)
 - Limits delays, but may introduce unfairness and livelock

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.26

When is Disk Performance Highest?

- When there are big sequential reads, or
- When there is so much work to do that they can be piggy backed (reordering queues—one moment)
- OK to be inefficient when things are mostly idle
- Bursts are both a threat and an opportunity
- <your idea for optimization goes here>
 - Waste space for speed?
- Other techniques:
 - Reduce overhead through user level drivers
 - Reduce the impact of I/O delays by doing other useful work in the meantime

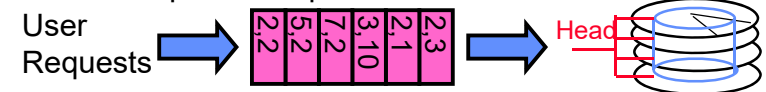
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

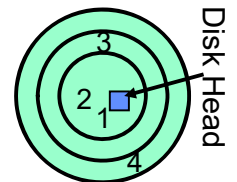
Lec 18.27

Disk Scheduling (1/2)

- Disk can do only one request at a time; What order do you choose to do queued requests?



- FIFO Order
 - Fair among requesters, but order of arrival may be to random spots on the disk \Rightarrow Very long seeks
- SSTF: Shortest seek time first
 - Pick the request that's closest on the disk
 - Although called SSTF, today must include rotational delay in calculation, since rotation can be as long as seek
 - Con: SSTF good at reducing seeks, but may lead to starvation



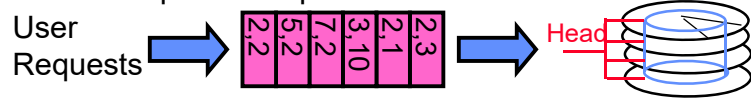
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

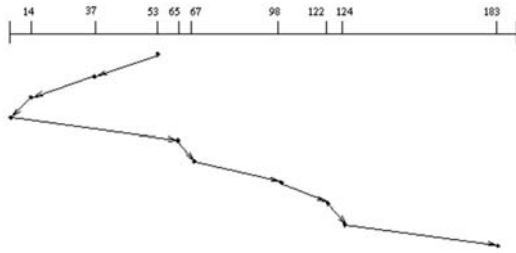
Lec 18.28

Disk Scheduling (2/2)

- Disk can do only one request at a time; What order do you choose to do queued requests?



- SCAN: Implements an Elevator Algorithm: take the closest request in the direction of travel
 - No starvation, but retains flavor of SSTF



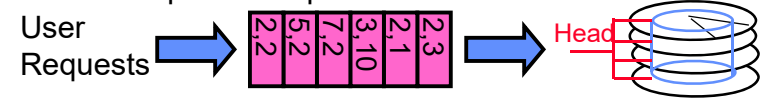
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

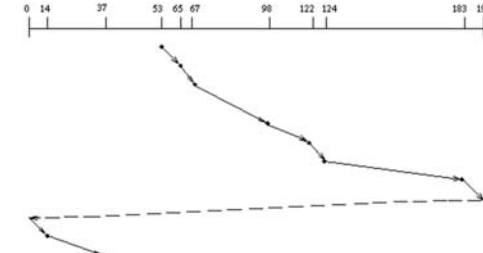
Lec 18.29

Disk Scheduling (2/2)

- Disk can do only one request at a time; What order do you choose to do queued requests?



- C-SCAN: Circular-Scan: only goes in one direction
 - Skips any requests on the way back
 - Fairer than SCAN, not biased towards pages in middle



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.30

Recall: How do we Hide I/O Latency?

- Blocking Interface:** “Wait”
 - When request data (e.g., read() system call), put process to sleep until data is ready
 - When write data (e.g., write() system call), put process to sleep until device is ready for data
- Non-blocking Interface:** “Don’t Wait”
 - Returns quickly from read or write request with count of bytes successfully transferred to kernel
 - Read may return nothing, write may write nothing
- Asynchronous Interface:** “Tell Me Later”
 - When requesting data, take pointer to user’s buffer, return immediately; later kernel fills buffer and notifies user
 - When sending data, take pointer to user’s buffer, return immediately; later kernel takes data and notifies user

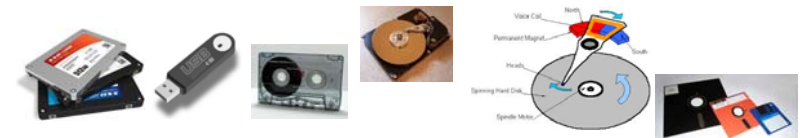
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.31

I/O & Storage Layers

Operations, Entities and Interface



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.32

Recall: C Low level I/O

- Operations on File Descriptors – as OS object representing the state of a file
 - User has a “handle” on the descriptor

```
#include <fcntl.h>
#include <unistd.h>
#include <sys/types.h>

int open (const char *filename, int flags [, mode_t mode])
int create (const char *filename, mode_t mode)
int close (int filedes)
```

Bit vector of:

- Access modes (Rd, Wr, ...)
- Open Flags (Create, ...)
- Operating modes (Appends, ...)

Bit vector of Permission Bits:

- User|Group|Other X R|W|X

http://www.gnu.org/software/libc/manual/html_node/Opening-and-Closing-Files.html

Recall: C Low Level Operations

```
ssize_t read (int filedes, void *buffer, size_t maxsize)
- returns bytes read, 0 => EOF, -1 => error
ssize_t write (int filedes, const void *buffer, size_t size)
- returns bytes written
off_t lseek (int filedes, off_t offset, int whence)
- set the file offset
* if whence == SEEK_SET: set file offset to “offset”
* if whence == SEEK_CUR: set file offset to crt location + “offset”
* if whence == SEEK_END: set file offset to file size + “offset”
int fsync (int filedes)
- wait for i/o of filedes to finish and commit to disk
void sync (void) - wait for ALL to finish and commit to disk
```

- When write returns, data is on its way to disk and can be read, but it may not actually be permanent!

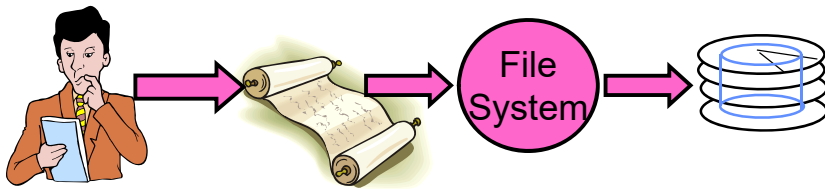
Building a File System

- **File System:** Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.
- File System Components
 - Naming: Interface to find files by name, not by blocks
 - Disk Management: collecting disk blocks into files
 - Protection: Layers to keep data secure
 - Reliability/Durability: Keeping of files durable despite crashes, media failures, attacks, etc.

Recall: User vs. System View of a File

- User’s view:
 - Durable Data Structures
- System’s view (system call interface):
 - Collection of Bytes (UNIX)
 - Doesn’t matter to system what kind of data structures you want to store on disk!
- System’s view (inside OS):
 - Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
 - Block size \geq sector size; in UNIX, block size is 4KB

Recall: Translating from User to System View



- What happens if user says: give me bytes 2—12?
 - Fetch block corresponding to those bytes
 - Return just the correct portion of the block
- What about: write bytes 2—12?
 - Fetch block
 - Modify portion
 - Write out Block
- Everything inside File System is in whole size blocks
 - For example, `getc()`, `putc()` \Rightarrow buffers something like 4096 bytes, even if interface is one byte at a time
- From now on, file is a collection of blocks

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.37

Disk Management Policies (1/2)

- Basic entities on a disk:
 - **File**: user-visible group of blocks arranged sequentially in logical space
 - **Directory**: user-visible index mapping names to files
- Access disk as linear array of sectors. Two Options:
 - Identify sectors as vectors [cylinder, surface, sector], sort in cylinder-major order, not used anymore
 - **Logical Block Addressing (LBA)**: Every sector has integer address from zero up to max number of sectors
 - » First case: OS/BIOS must deal with bad sectors
 - » Second case: hardware shields OS from structure of disk

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.38

Disk Management Policies (2/2)

- Need way to track free disk blocks
 - Link free blocks together \Rightarrow too slow today
 - Use bitmap to represent free space on disk
- Need way to structure files: **File Header**
 - Track which blocks belong at which offsets within the logical file structure
 - **Optimize placement of files' disk blocks to match access and usage patterns**

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.39

Designing a File System ...

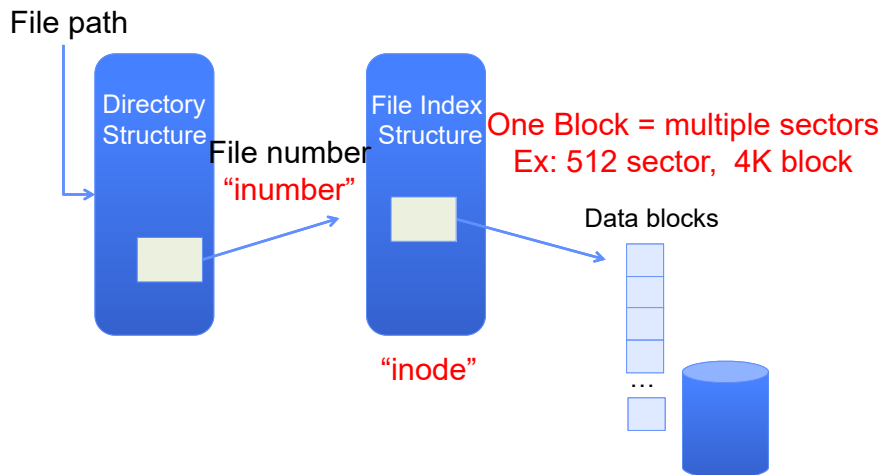
- What factors are critical to the design choices?
- Durable data store \Rightarrow it's all on disk
- (Hard) Disks Performance !!!
 - Maximize sequential access, minimize seeks
- Open before Read/Write
 - Can perform protection checks and look up where the actual file resource are, in advance
- Size is determined as they are used !!!
 - Can write (or read zeros) to expand the file
 - Start small and grow, need to make room
- Organized into directories
 - What data structure (on disk) for that?
- Need to allocate / free blocks
 - Such that access remains efficient

4/9/19

Kubiawicz CS162 ©UCB Spring 2019

Lec 18.40

Components of a File System



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.41

Components of a file system

file name \longrightarrow file number \longrightarrow Storage block
 offset directory offset index structure

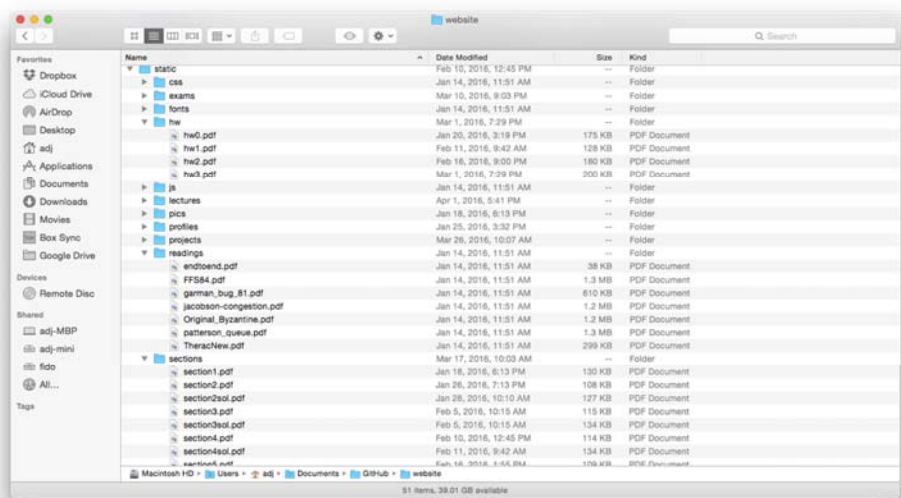
- Open performs **Name Resolution**
 - Translates pathname into a “file number”
 - » Used as an “index” to locate the blocks
 - Creates a file descriptor in PCB within kernel
 - Returns a “handle” (another integer) to user process
- Read, Write, Seek, and Sync operate on handle
 - Mapped to file descriptor and to blocks

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.42

Directories



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.43

Directory

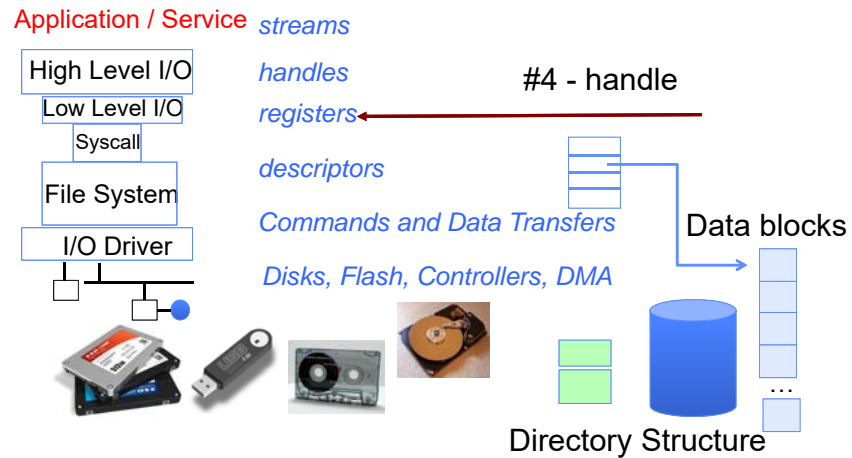
- Basically a hierarchical structure
- Each directory entry is a collection of
 - Files
 - Directories
 - » A link to another entries
- Each has a name and attributes
 - Files have data
- Links (hard links) make it a DAG, not just a tree
 - Softlinks (aliases) are another name for an entry

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.44

I/O & Storage Layers



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

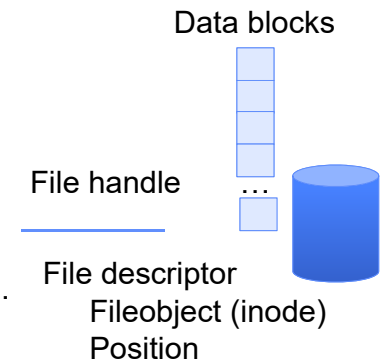
Lec 18.45

File

- Named permanent storage

- Contains

- Data
 - » Blocks on disk somewhere
- Metadata (Attributes)
 - » Owner, size, last opened, ...
 - » Access rights
 - R, W, X
 - Owner, Group, Other (in Unix systems)
 - Access control list in Windows system

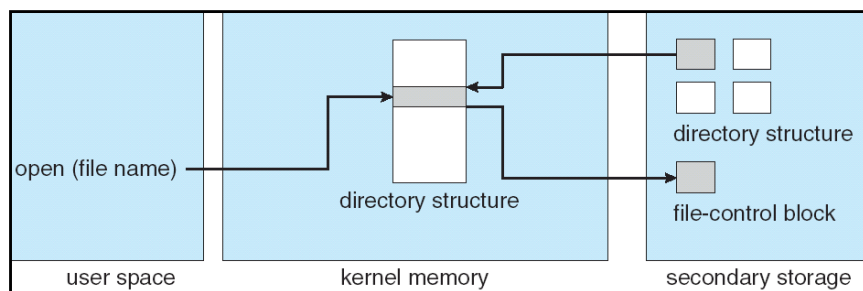


4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.46

In-Memory File System Structures



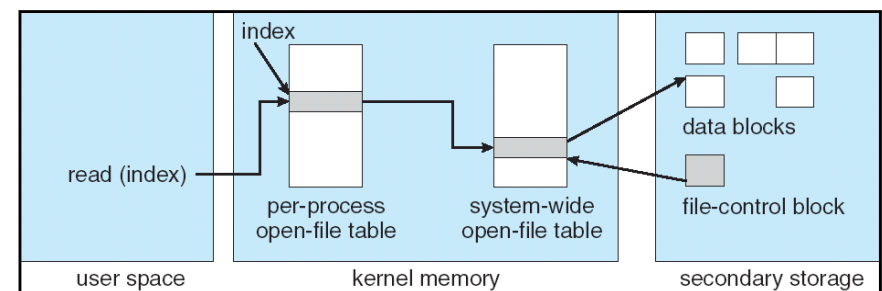
- Open system call:
 - Resolves file name, finds file control block (inode)
 - Makes entries in per-process and system-wide tables
 - Returns index (called “file handle”) in open-file table

4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.47

In-Memory File System Structures



- Read/write system calls:
 - Use file handle to locate inode
 - Perform appropriate reads or writes

4/9/19

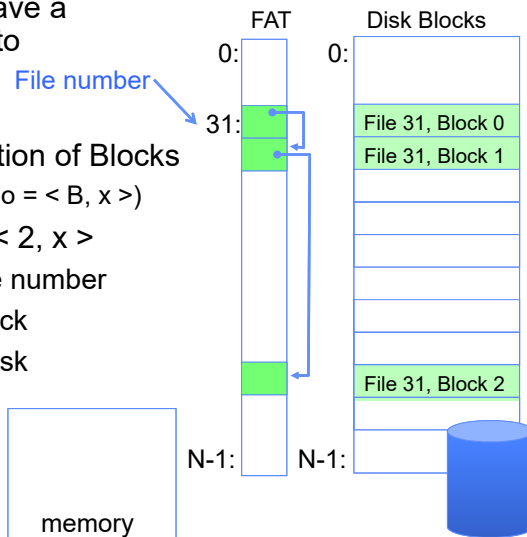
Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.48

Our first filesystem: FAT (File Allocation Table)

- The most commonly used filesystem in the world!

- Assume (for now) we have a way to translate a path to a "file number"
 - i.e., a directory structure
- Disk Storage is a collection of Blocks
 - Just hold file data (offset $o = \langle B, x \rangle$)
- Example: file_read 31, $\langle 2, x \rangle$
 - Index into FAT with file number
 - Follow linked list to block
 - Read the block from disk into memory



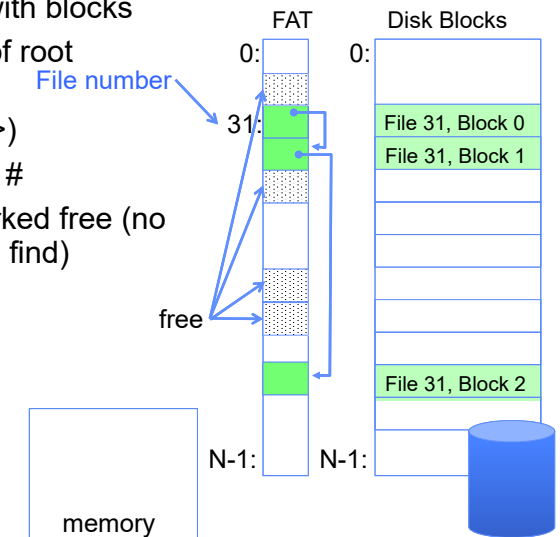
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.49

FAT Properties

- File is collection of disk blocks
- FAT is linked list 1-1 with blocks
- File Number is index of root of block list for the file
- File offset ($o = \langle B, x \rangle$)
- Follow list to get block #
- Unused blocks \Leftrightarrow Marked free (no ordering, must scan to find)



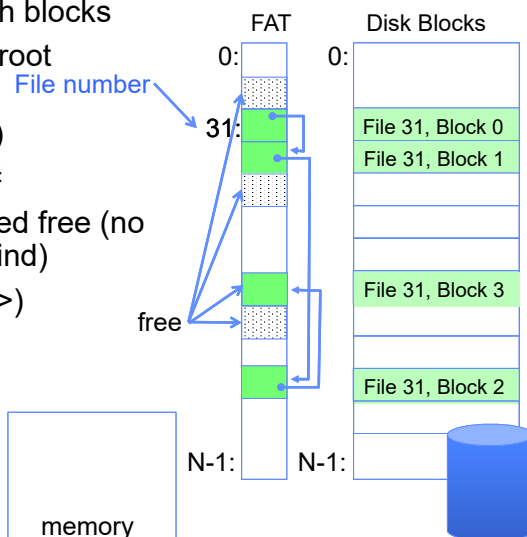
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.50

FAT Properties

- File is collection of disk blocks
- FAT is linked list 1-1 with blocks
- File Number is index of root of block list for the file
- File offset ($o = \langle B, x \rangle$)
- Follow list to get block #
- Unused blocks \Leftrightarrow Marked free (no ordering, must scan to find)
- Ex: file_write(31, $\langle 3, y \rangle$)
 - Grab free block
 - Linking them into file



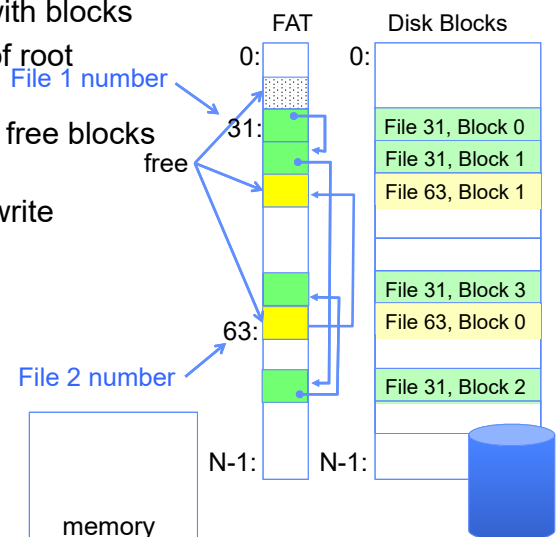
4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.51

FAT Properties

- File is collection of disk blocks
- FAT is linked list 1-1 with blocks
- File Number is index of root of block list for the file
- Grow file by allocating free blocks and linking them in
- Ex: Create file, write, write



4/9/19

Kubiatowicz CS162 ©UCB Spring 2019

Lec 18.52

Many Huge FAT Security Holes!

- FAT has no access rights
- FAT has no header in the file blocks
- Just gives an index into the FAT
 - (file number = block number)

Summary

- Bursts & High Utilization introduce queuing delays
- Queuing Latency:
 - M/M/1 and M/G/1 queues: simplest to analyze
 - As utilization approaches 100%, latency $\rightarrow \infty$
 $T_q = T_{ser} \times \frac{1}{2}(1+C) \times u/(1-u)$
- File System:
 - Transforms blocks into Files and Directories
 - Optimize for access and usage patterns
 - Maximize sequential access, allow efficient random access
- File (and directory) defined by header, called “inode”
- File Allocation Table (FAT) Scheme
 - Linked-list approach
 - Very widely used: Cameras, USB drives, SD cards
 - Simple to implement, but poor performance and no security
- Look at actual file access patterns – many small files, but large files take up all the space!