

CS162
Operating Systems and
Systems Programming
Lecture 21

Filesystems 1: Performance,
Queueing Theory, Filesystem Design

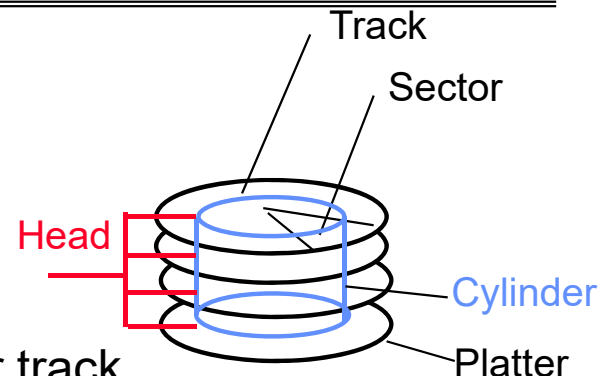
April 9th, 2024

Prof. John Kubiatowicz

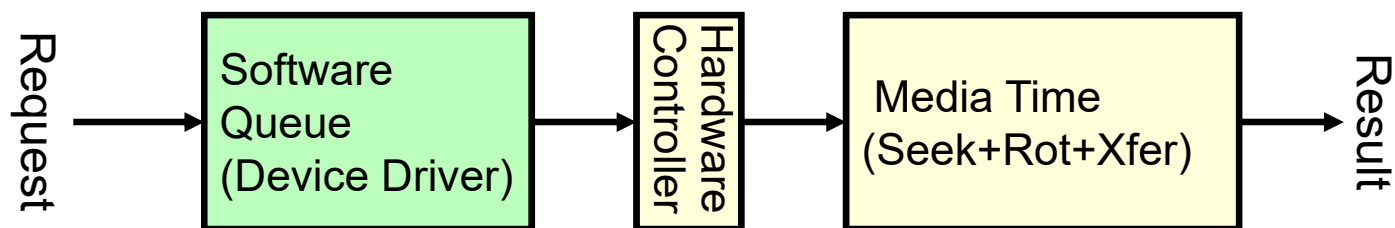
<http://cs162.eecs.Berkeley.edu>

Recall: Magnetic Disks

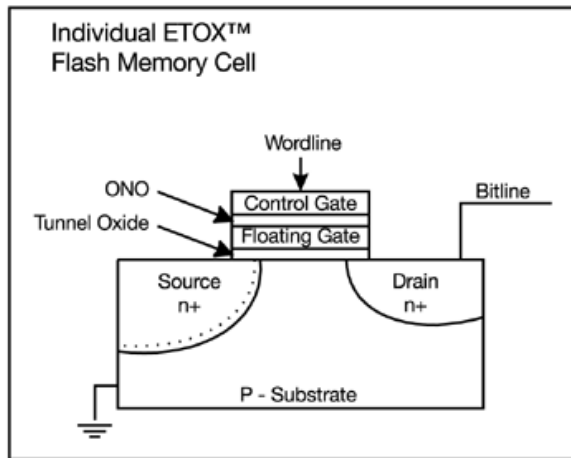
- **Cylinders:** all the tracks under the head at a given point on all surfaces
- Read/write data is a three-stage process:
 - **Seek time:** position the head/arm over the proper track
 - **Rotational latency:** wait for desired sector to rotate under r/w head
 - **Transfer time:** transfer a block of bits (sector) under r/w head



$$\text{Disk Latency} = \text{Queueing Time} + \text{Controller time} + \text{Seek Time} + \text{Rotation Time} + \text{Xfer Time}$$



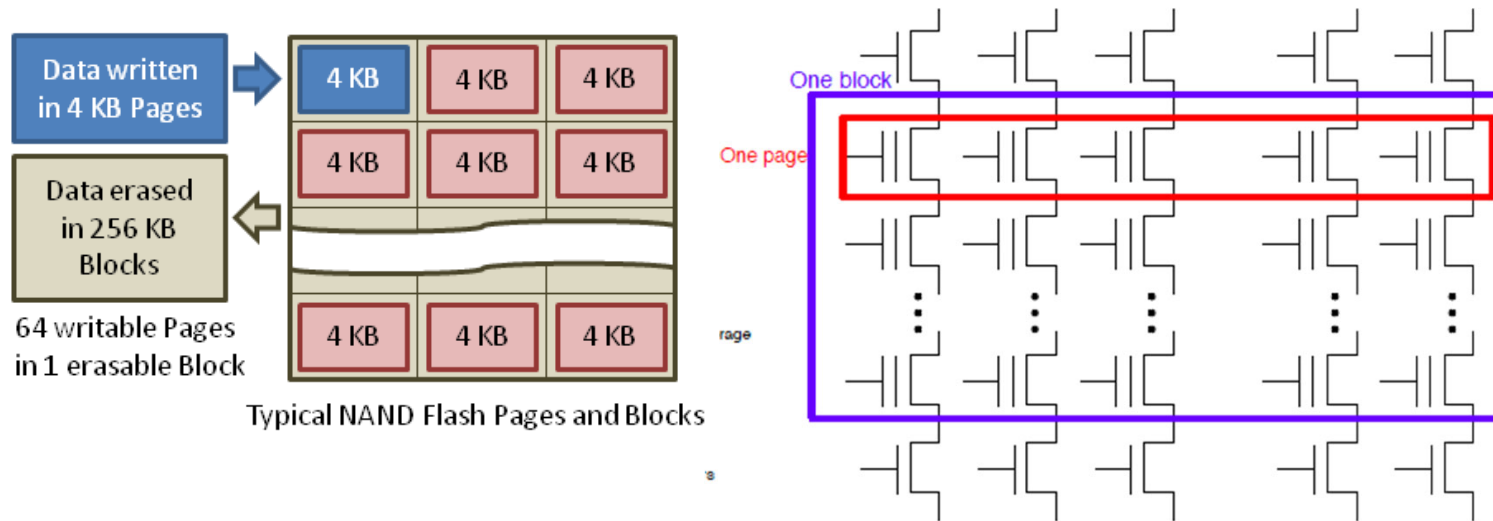
Recall: FLASH Memory



Samsung 2015: 512GB, NAND Flash

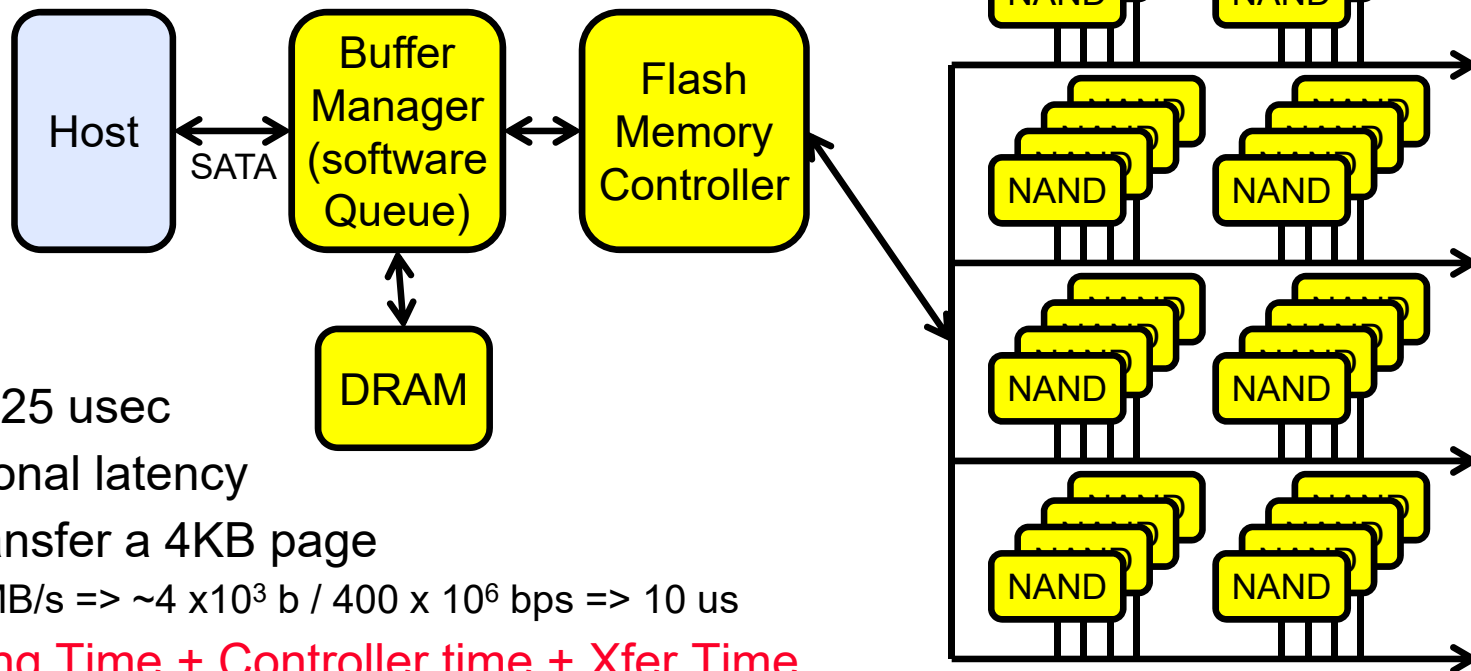
- Like a normal transistor but:
 - Has a floating gate that can hold charge
 - To write: raise or lower wordline high enough to cause charges to tunnel
 - To read: turn on wordline as if normal transistor
 - » presence of charge changes threshold and thus measured current
- Two varieties:
 - NAND: denser, must be read and written in blocks
 - NOR: much less dense, fast to read and write
- V-NAND: 3D stacking (Samsung claims 1TB possible in 1 chip)

Flash Memory (Con't)



- Data read and written in page-sized chunks (e.g. 4K)
 - Cannot be addressed at byte level
 - Random access at block level for reads (no locality advantage)
 - Writing of new blocks handled in order (kinda like a log)
- Before writing, must be *erased* (256K block at a time)
 - Requires free-list management
 - CANNOT write over existing block (Copy-on-Write is normal case)

SSD Architecture – Reads



Read 4 KB Page: ~25 usec

– No seek or rotational latency

– Transfer time: transfer a 4KB page

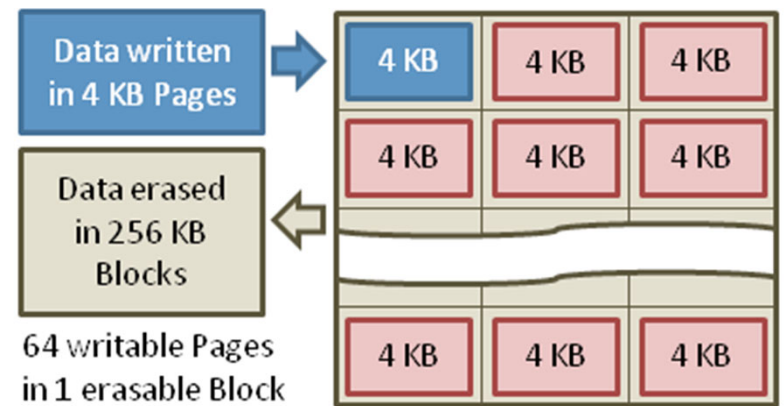
» SATA: $300\text{-}600\text{MB/s} \Rightarrow \sim 4 \times 10^3 \text{ b} / 400 \times 10^6 \text{ bps} \Rightarrow 10 \text{ us}$

– **Latency = Queuing Time + Controller time + Xfer Time**

– **Highest Bandwidth:** Sequential OR Random reads

SSD Architecture – Writes

- Writing data to NAND Flash is complex!
 - Can only write empty pages in a block (~ 200µs)
 - Erasing a block takes ~1.5ms
 - Controller maintains pool of empty blocks by coalescing used pages (read, erase, write), also reserves some % of capacity
 - Rule of thumb: writes 10x reads, erasure 10x writes
- SSDs provide same interface as HDDs: read and write chunk (4KB) at a time
- Why not just erase and rewrite new version of entire 256KB block?
 - Erasure is very slow (milliseconds)
 - Each block has a finite lifetime, can only be erased and rewritten about 10K times
 - Heavily used blocks likely to wear out quickly

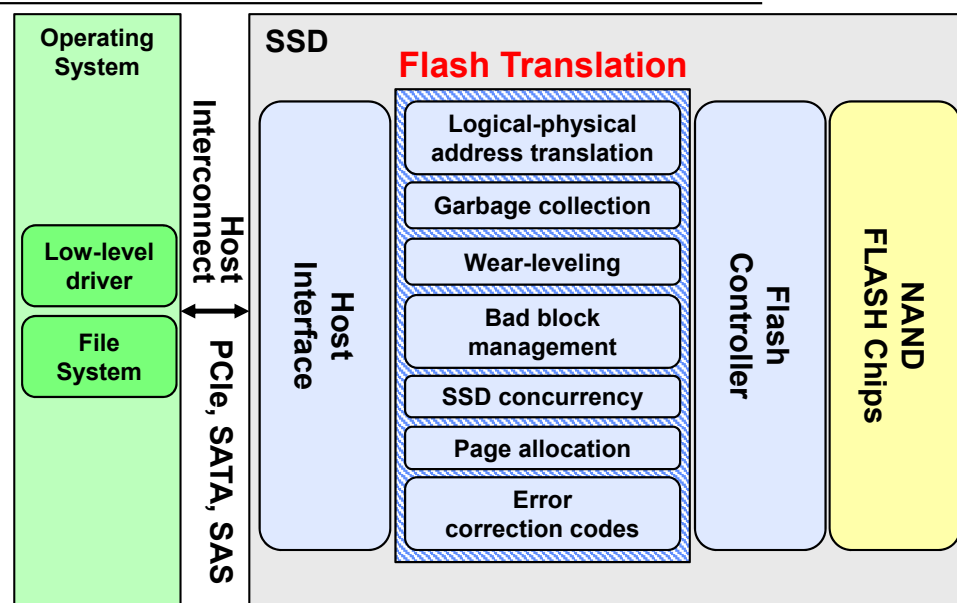


Typical NAND Flash Pages and Blocks

https://en.wikipedia.org/wiki/Solid-state_drive

Managing Writes: Flash Translation Layer

- Maintain *Flash Translation Layer (FTL)* in SSD
 - Layer of Indirection between OS and FLASH
 - Map virtual block numbers (which OS uses) to physical page numbers (which flash mem. controller uses)
 - **Can now freely relocate data w/o OS knowing**
- FTL advantages/mechanism:
 - Copy on Write: No need to immediately erase entire 256K block when modifying 4K page
 - » Don't overwrite page when OS updates data
 - » Instead, write new version in a free page
 - » Update FTL mapping to point to new location
 - Wear Levelling: Try to wear out NAND evenly
 - » SSD controller can assign mappings to spread workload across pages
 - What to do with old versions of pages?
 - » *Garbage Collection* in background
 - » Erase blocks with old pages, add to free list



Some “Current” (large) 3.5in SSDs

- Seagate Exos SSD: 15.36TB (2017)
 - Dual 12Gb/s interface
 - Seq reads 860MB/s
 - Seq writes 920MB/s
 - Random Reads (IOPS): 102K
 - Random Writes (IOPS): 15K
 - Price (Amazon): \$5495 (\$0.36/GB)
- Nimbus SSD: 100TB (2019)
 - Dual port: 12Gb/s interface
 - Seq reads/writes: 500MB/s
 - Random Read Ops (IOPS): 100K
 - *Unlimited writes for 5 years!*
 - Price: ~ \$40K? (\$0.4/GB)
 - » However, 50TB drive costs \$12500 (\$0.25/GB)



Amusing calculation: Is a full Kindle heavier than an empty one?

- Actually, “Yes”, but not by much
- Flash works by trapping electrons:
 - So, erased state lower energy than written state
- Assuming that:
 - Kindle has 4GB flash
 - $\frac{1}{2}$ of all bits in full Kindle are in high-energy state
 - High-energy state about 10^{-15} joules higher
 - Then: Full Kindle is 1 attogram (10^{-18} gram) heavier (Using $E = mc^2$)
- Of course, this is less than most sensitive scale can measure (it can measure 10^{-9} grams)
- Of course, this weight difference overwhelmed by battery discharge, weight from getting warm,
- **Source: John Kubiawicz (New York Times, Oct 24, 2011)**

SSD Summary

- Pros (vs. hard disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:
 - » Very light weight, low power, silent, very shock insensitive
 - Read at memory speeds (limited by controller and I/O bus)
- Cons
 - Small storage (0.1-0.5x disk), expensive (3-20x disk)
 - » Hybrid alternative: combine small SSD with large HDD

SSD Summary

- Pros (vs. hard disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:
 - » Very light weight, low power, silent, very shock insensitive
 - Read at memory speeds (limited by controller and I/O bus)
- Cons
 - ~~Small storage (0.1-0.5x disk), expensive (5-20x disk)~~
 - » Hybrid alternative: combine small SSD with large HDD
 - Asymmetric block write performance: read pg/erase/write pg
 - » Controller garbage collection (GC) algorithms have major effect on performance
 - Limited drive lifetime
 - » 1-10K writes/page for MLC NAND
 - » Avg failure rate is 6 years, life expectancy is 9–11 years
- These are changing rapidly!

No
longer
true!

Administrivia (4/9/2024)

- Midterm 3: Thursday April 25th
 - All topics up to and including lecture on the 23rd
 - 3 sheets of notes, double-sided
- Extra (fun!) lecture on Tuesday April 30th
 - Topics TBA
- Class attendance: No credit for people who use the same photo!
- Data4All@Berkeley: This Friday!
 - Friday 4/12, 12:00-1:00 in Soda 510
 - Undergraduate or Masters students interested in Systems broadly defined (DB, Arch, Sec, Networking, Systems, etc.) who identify as an URM in Computer Science
 - Come by for free lunch to meet fellow interested students
 - Talk to relevant faculty, discuss possible classes, research opportunities in systems, as well as the best pizza topping!



<https://tinyurl.com/3r3cj3ya>

Ways of Measuring Performance: Times (s) and Rates (op/s)

- **Latency** – time to complete a task
 - Measured in units of time (s, ms, us, ..., hours, years)
- **Response Time** - time to initiate and operation and get its response
 - Able to issue one that *depends* on the result
 - Know that it is done (anti-dependence, resource usage)
- **Throughput** or **Bandwidth** – rate at which tasks are performed
 - Measured in units of things per unit time (ops/s, GFLOP/s)
- **Start up or “Overhead”** – time to initiate an operation
- Most I/O operations are roughly linear in b bytes
 - $\text{Latency}(b) = \text{Overhead} + b/\text{TransferCapacity}$
- Performance???
 - Operation time (4 mins to run a mile...)
 - Rate (mph, mpg, ...)

Example: Overhead in Fast Network

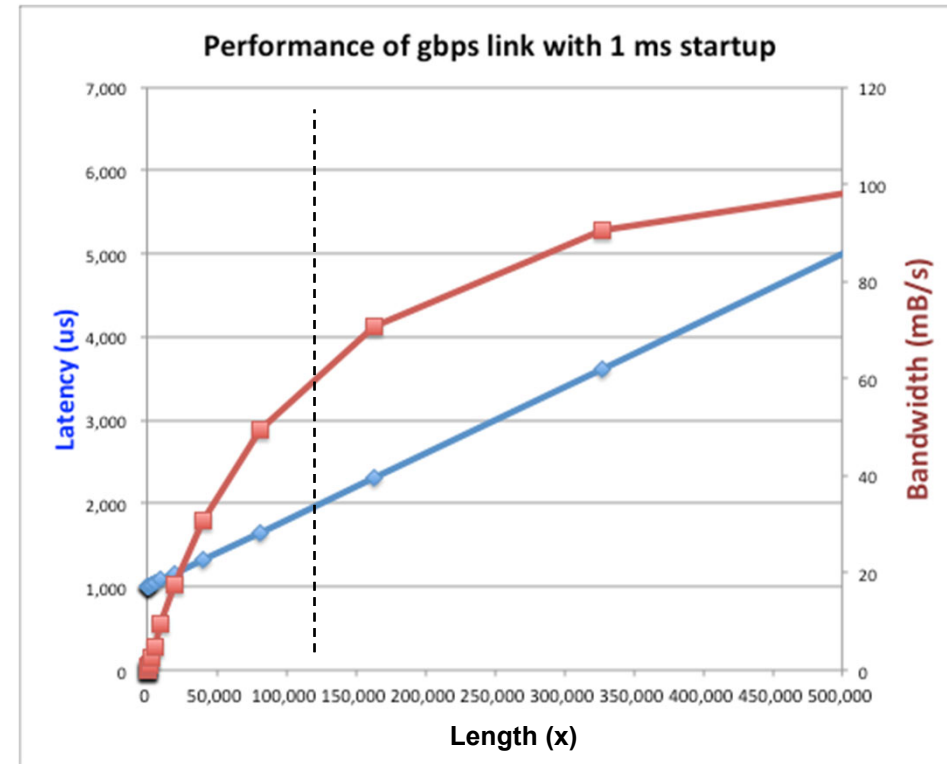
- Consider a 1 Gb/s link ($B_W = 125 \text{ MB/s}$) with startup cost $S = 1 \text{ ms}$

- Latency: $L(x) = S + \frac{x}{B_W}$

- Effective Bandwidth:

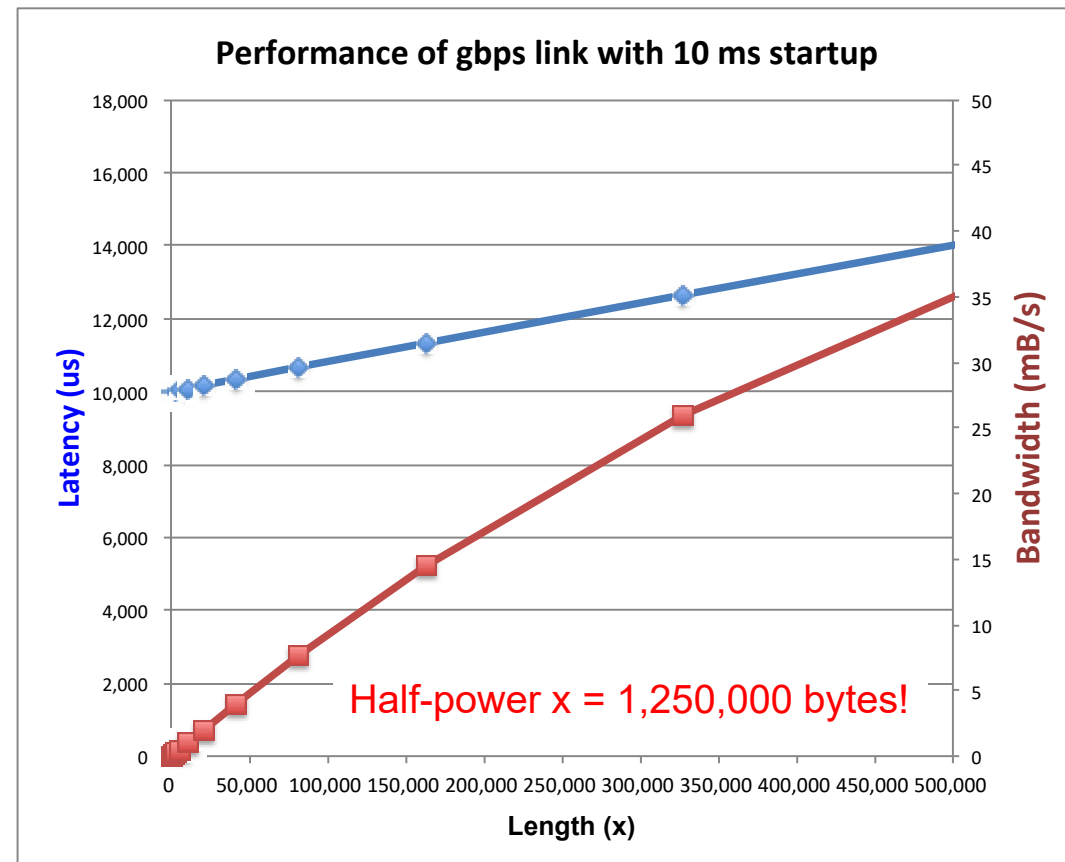
$$E(x) = \frac{x}{S + \frac{x}{B_W}} = \frac{B_W \cdot x}{B_W \cdot S + x} = \frac{B_W}{\frac{B_W \cdot S}{x} + 1}$$

- Half-power Bandwidth: $E(x) = \frac{B_W}{2}$
- For this example, half-power bandwidth occurs at $x = 125 \text{ KB}$



Example: 10 ms Startup Cost (e.g., Disk)

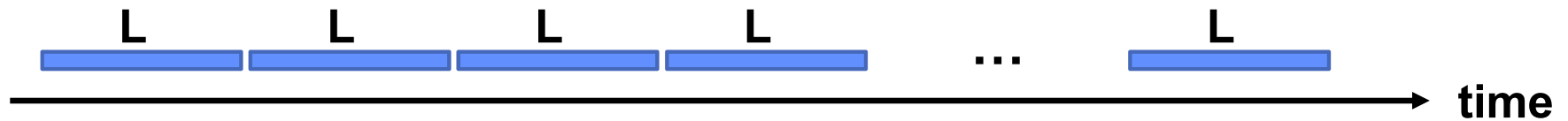
- Half-power bandwidth at $x = 1.25$ MB
- Large startup cost can degrade effective bandwidth
- Amortize it by performing I/O in larger blocks



What Determines Peak BW for I/O?

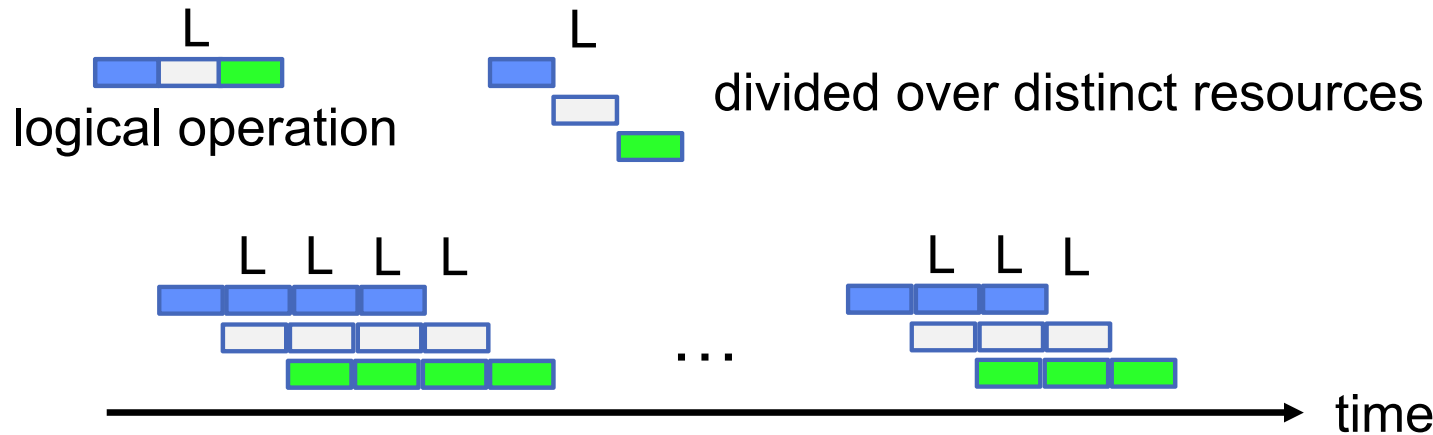
- Bus Speed
 - PCI-X: 1064 MB/s = 133 MHz x 64 bit (per lane)
 - ULTRA WIDE SCSI: 40 MB/s
 - Serial Attached SCSI & Serial ATA & IEEE 1394 (firewire): 1.6 Gb/s full duplex (200 MB/s)
 - USB 3.0 – 5 Gb/s
 - Thunderbolt 3 – 40 Gb/s
- Device Transfer Bandwidth
 - Rotational speed of disk
 - Write / Read rate of NAND flash
 - Signaling rate of network link
- Whatever is the bottleneck in the path...

Sequential Server Performance



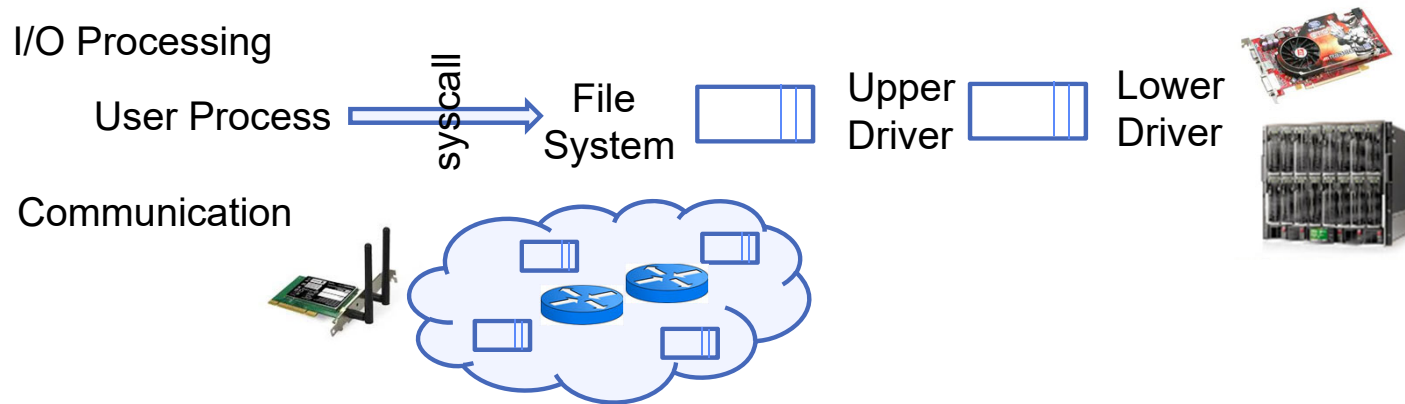
- Single sequential “server” that can deliver a task in time L operates at rate $\leq \frac{1}{L}$ (on average, in steady state, ...)
 - $L = 10 \text{ ms} \rightarrow B = 100 \text{ op/s}$
 - $L = 2 \text{ yr} \rightarrow B = 0.5 \text{ op/yr}$
- Applies to a processor, a disk drive, a person, a TA, ...

Single Pipelined Server



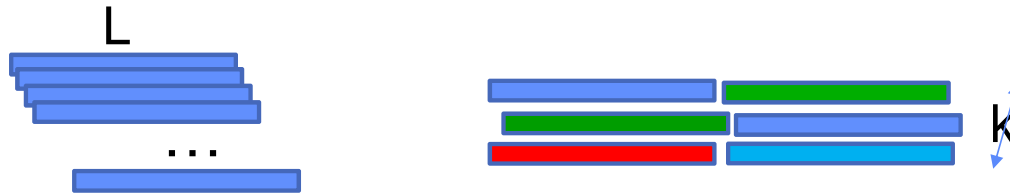
- Single pipelined server of k stages for tasks of length L (i.e., time L/k per stage) delivers at rate $\leq k/L$.
 - $L = 10 \text{ ms}, k = 4 \rightarrow B = 400 \text{ op/s}$
 - $L = 2 \text{ yr}, k = 2 \rightarrow B = 1 \text{ op/yr}$

Example Systems “Pipelines”



- Anything with queues between operational process behaves roughly “pipeline like”
- Important difference is that “initiations” are decoupled from processing
 - May have to queue up a burst of operations
 - Not synchronous and deterministic like in 61C

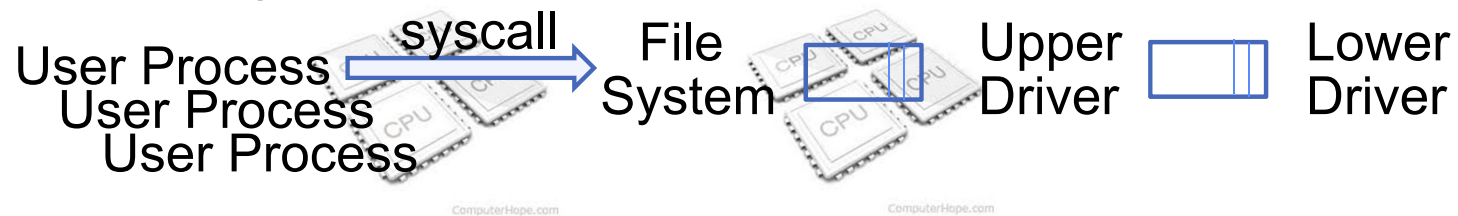
Multiple Servers



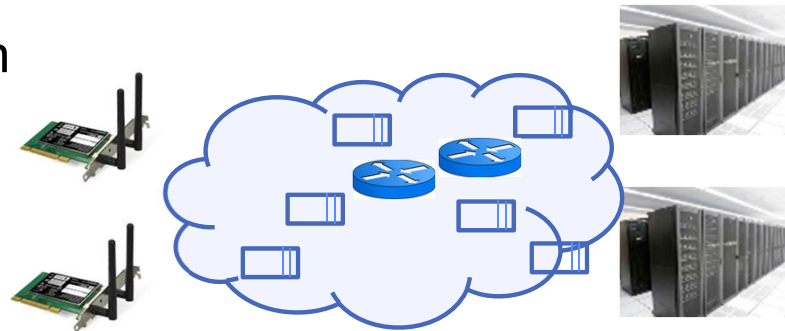
- k servers handling tasks of length L delivers at rate $\leq k/L$.
 - $L = 10 \text{ ms}$, $k = 4 \rightarrow B = 400 \text{ op/s}$
 - $L = 2 \text{ yr}$, $k = 2 \rightarrow B = 1 \text{ op/yr}$
- In 61C you saw multiple processors (cores)
 - Systems present lots of multiple parallel servers
 - Often with lots of queues

Example Systems “Parallelism”

I/O Processing

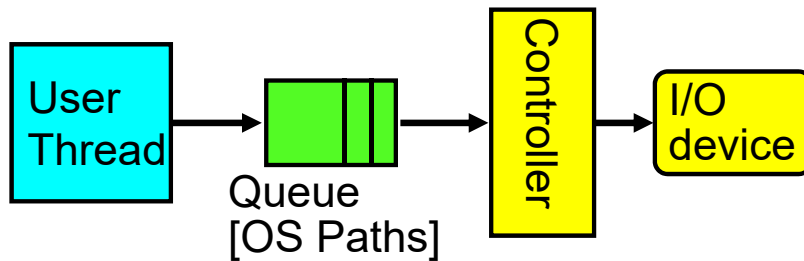


Communication

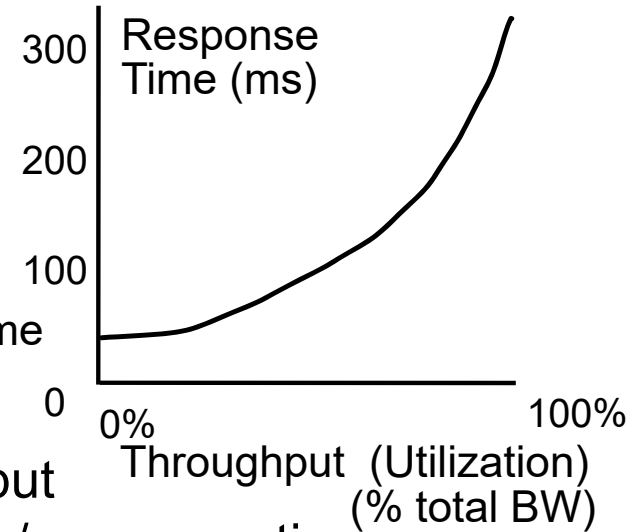


Parallel Computation, Databases, ...

I/O Performance



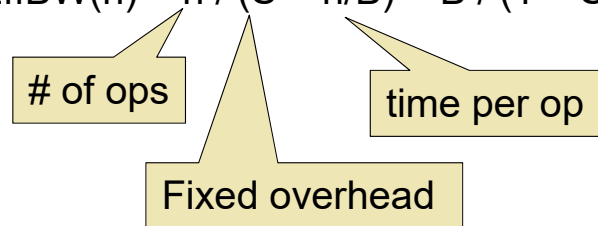
Response Time = Queue + I/O device service time



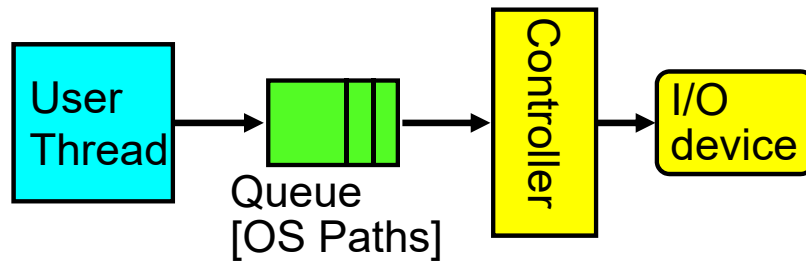
- Performance of I/O subsystem

- Metrics: Response Time, Throughput
- Effective BW per op = transfer size / response time

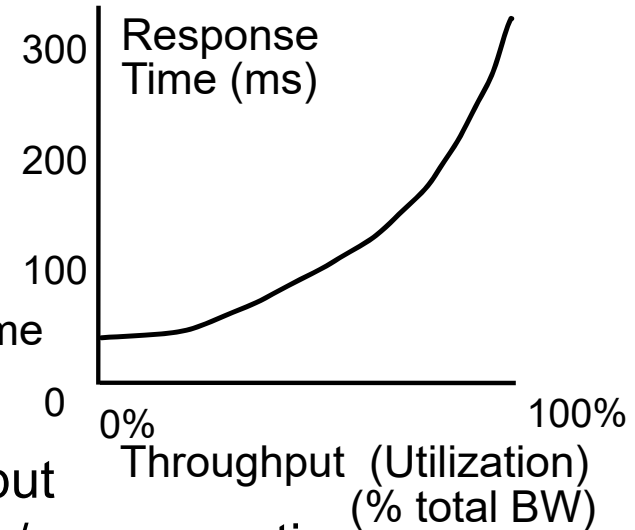
» $\text{EffBW}(n) = n / (S + n/B) = B / (1 + SB/n)$



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

- » $\text{EffBW}(n) = n / (S + n/B) = B / (1 + SB/n)$

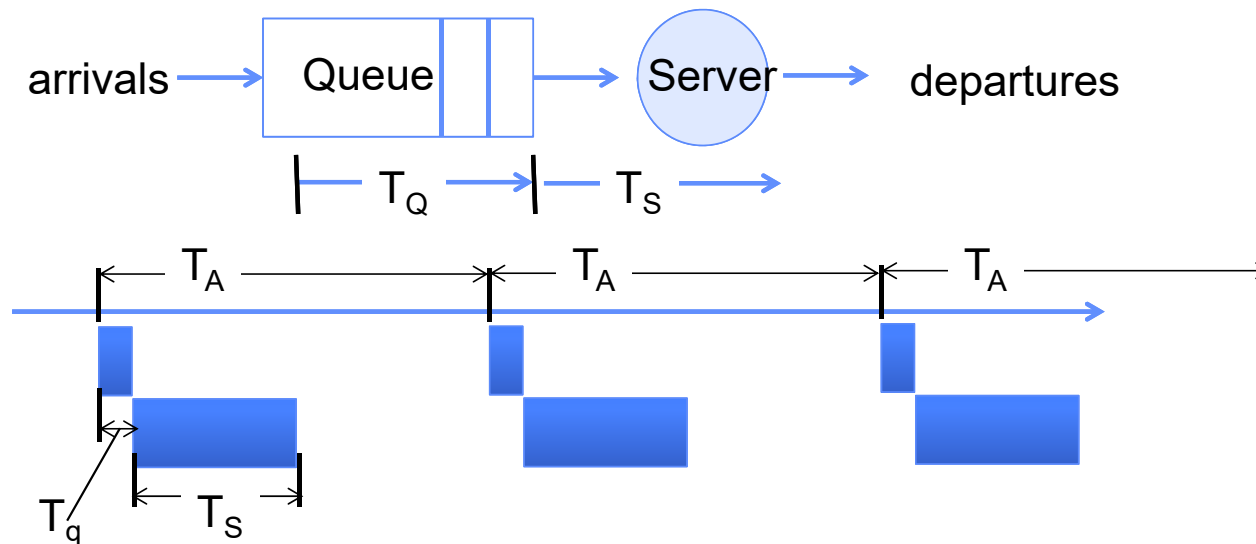
- Contributing factors to latency:

- » Software paths (can be loosely modeled by a queue)
- » Hardware controller
- » I/O device service time

- Queuing behavior:

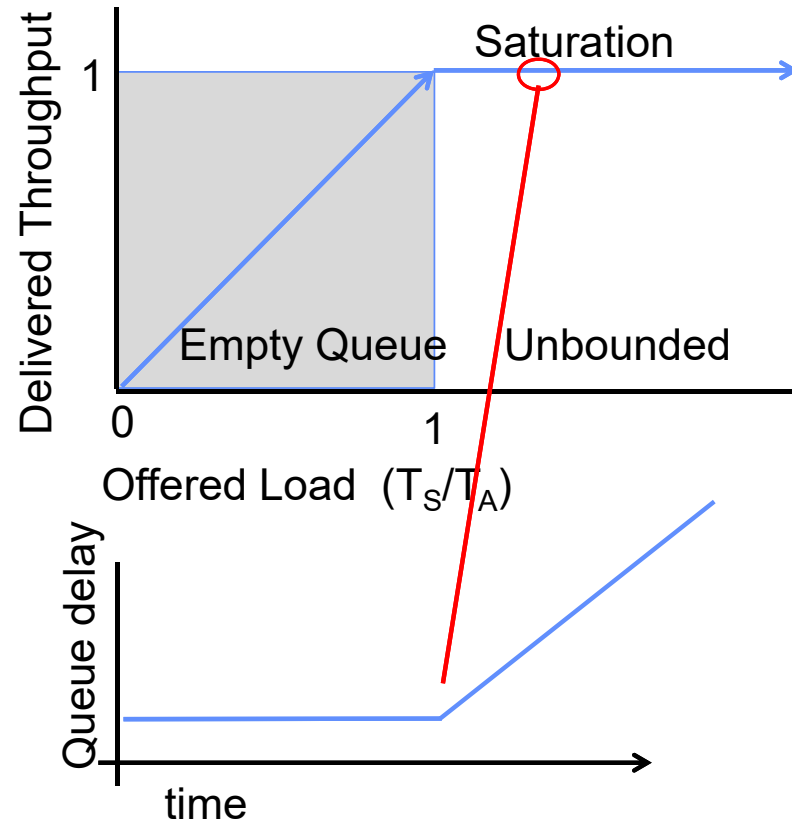
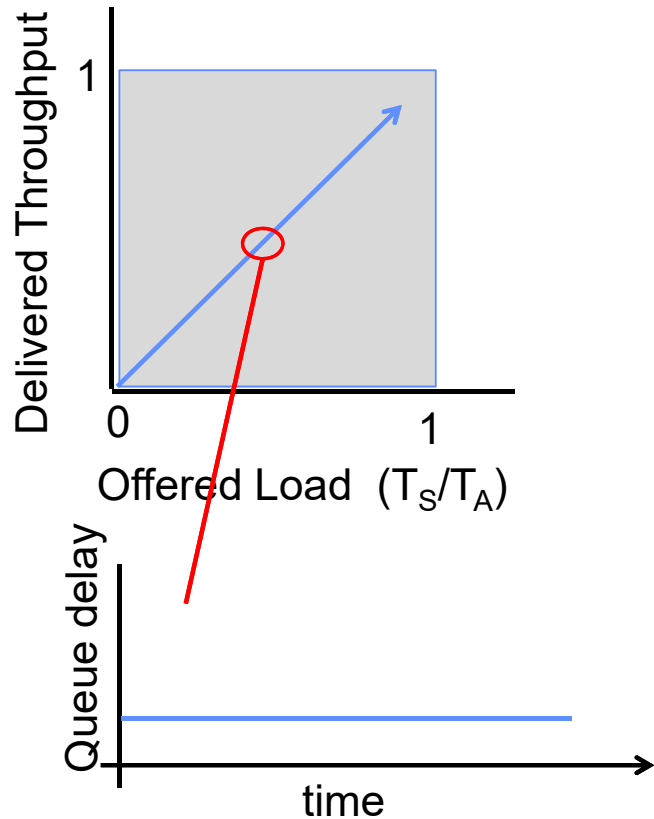
- Can lead to big increases of latency as utilization increases
- Solutions?

A Simple Deterministic World



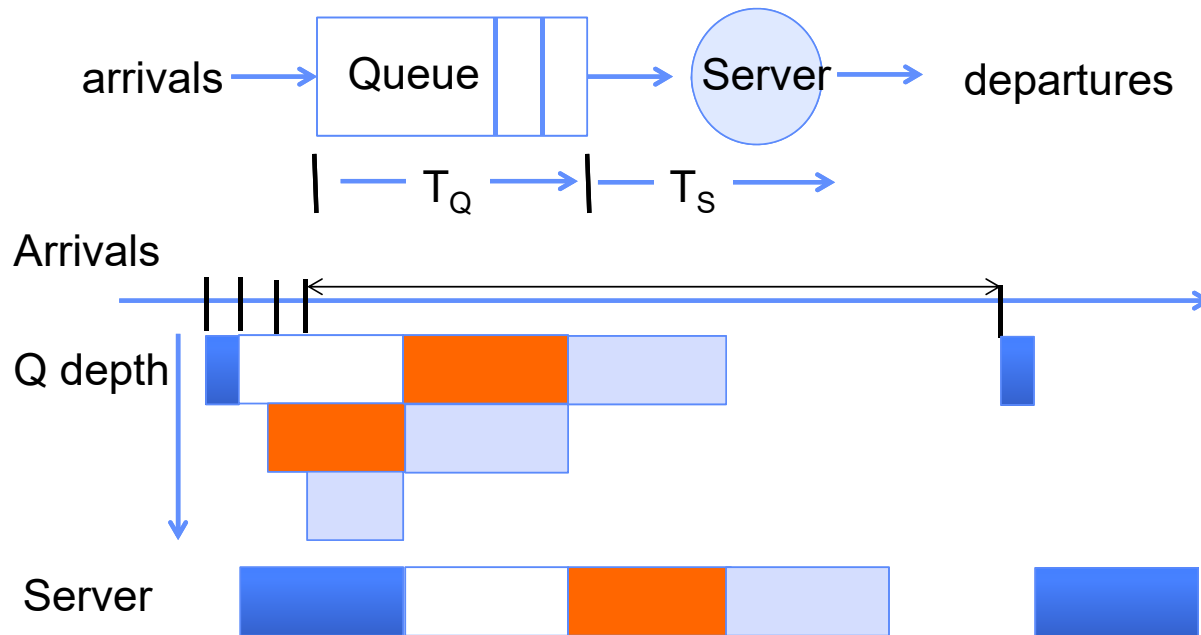
- Assume requests arrive at regular intervals, take a fixed time to process, with plenty of time between ...
- Service rate ($\mu = 1/T_S$) - operations per second
- Arrival rate: ($\lambda = 1/T_A$) - requests per second
- Utilization: $U = \lambda/\mu$, where $\lambda < \mu$
- Average rate is the complete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

A Ideal Linear World



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

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

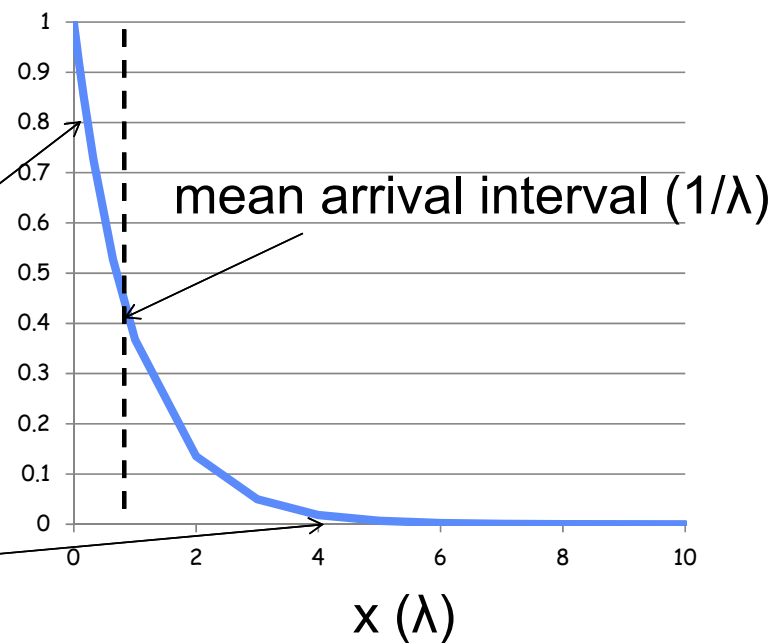
So how do we model the burstiness of arrival?

- 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

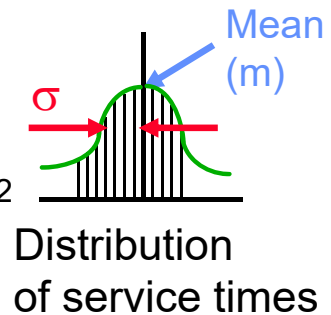
Lots of short arrival intervals (i.e., high instantaneous rate)

Few long gaps (i.e., low instantaneous rate)

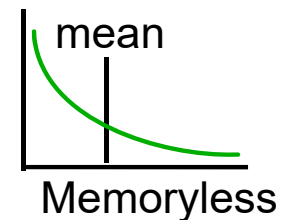


Background: General Use of Random Distributions

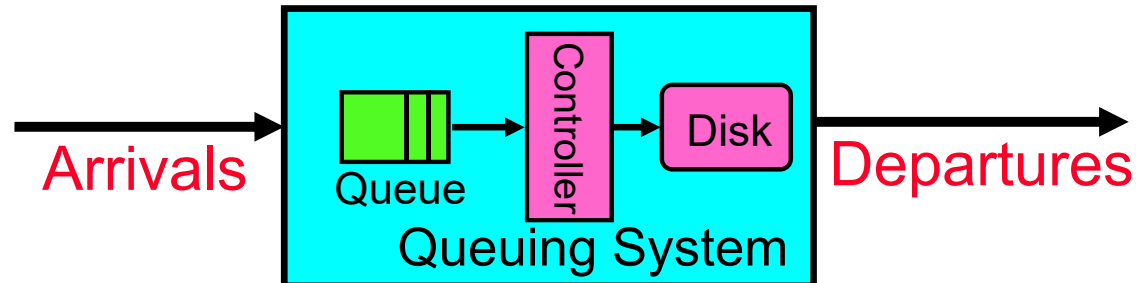
- Server spends variable time (T) with customers
 - Mean (Average) $m = \sum p(T) \times T$
 - Variance (stddev²) $\sigma^2 = \sum p(T) \times (T-m)^2 = \sum p(T) \times T^2 - m^2$
 - Squared coefficient of variance: $C = \sigma^2/m^2$Aggregate description of the distribution



- 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)

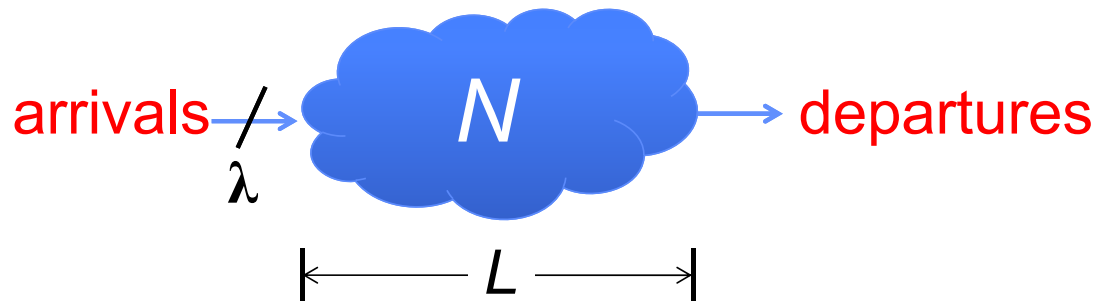


Introduction to Queuing Theory



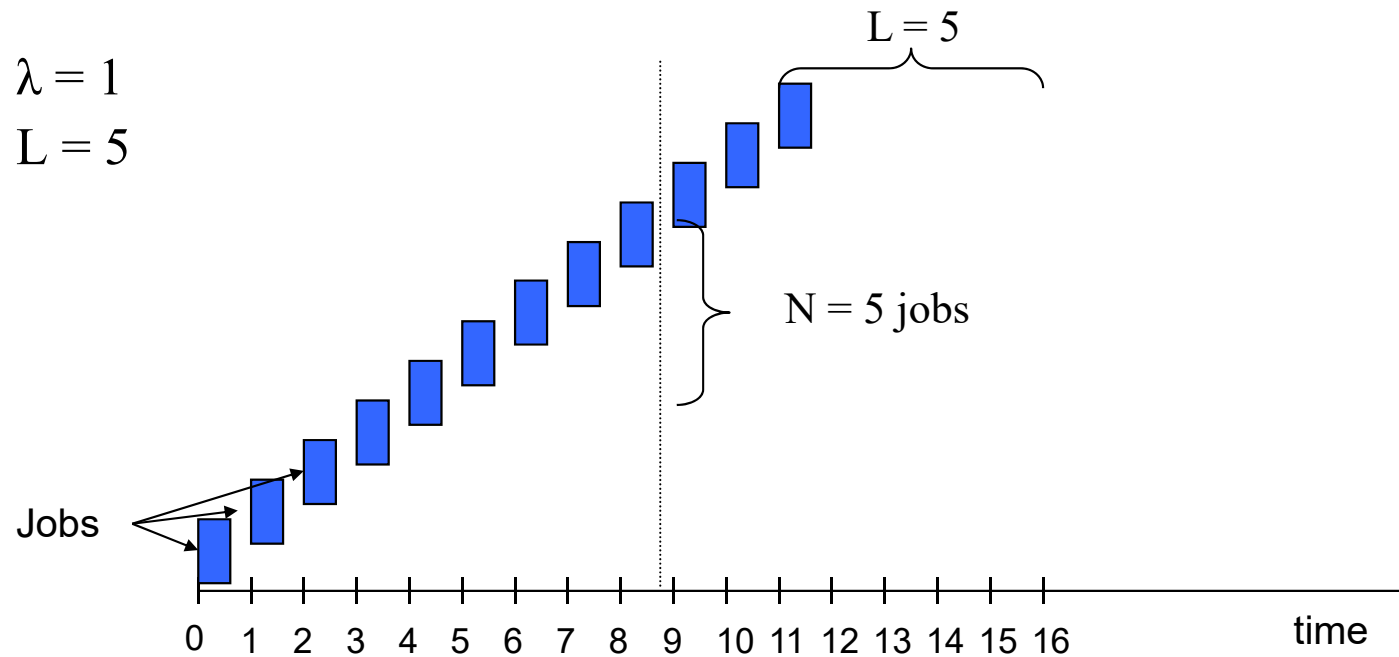
- What about queuing time??
 - Let's apply some queuing theory
 - Queuing Theory applies to long term, steady state behavior \Rightarrow Arrival rate = Departure rate
- Arrivals characterized by some probabilistic distribution
- Departures characterized by some probabilistic distribution

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 (jobs) = λ (jobs/s) \times L (s)
- Regardless of structure, bursts of requests, variation in service
 - Instantaneous variations, but it washes out in the average
 - Overall, requests match departures

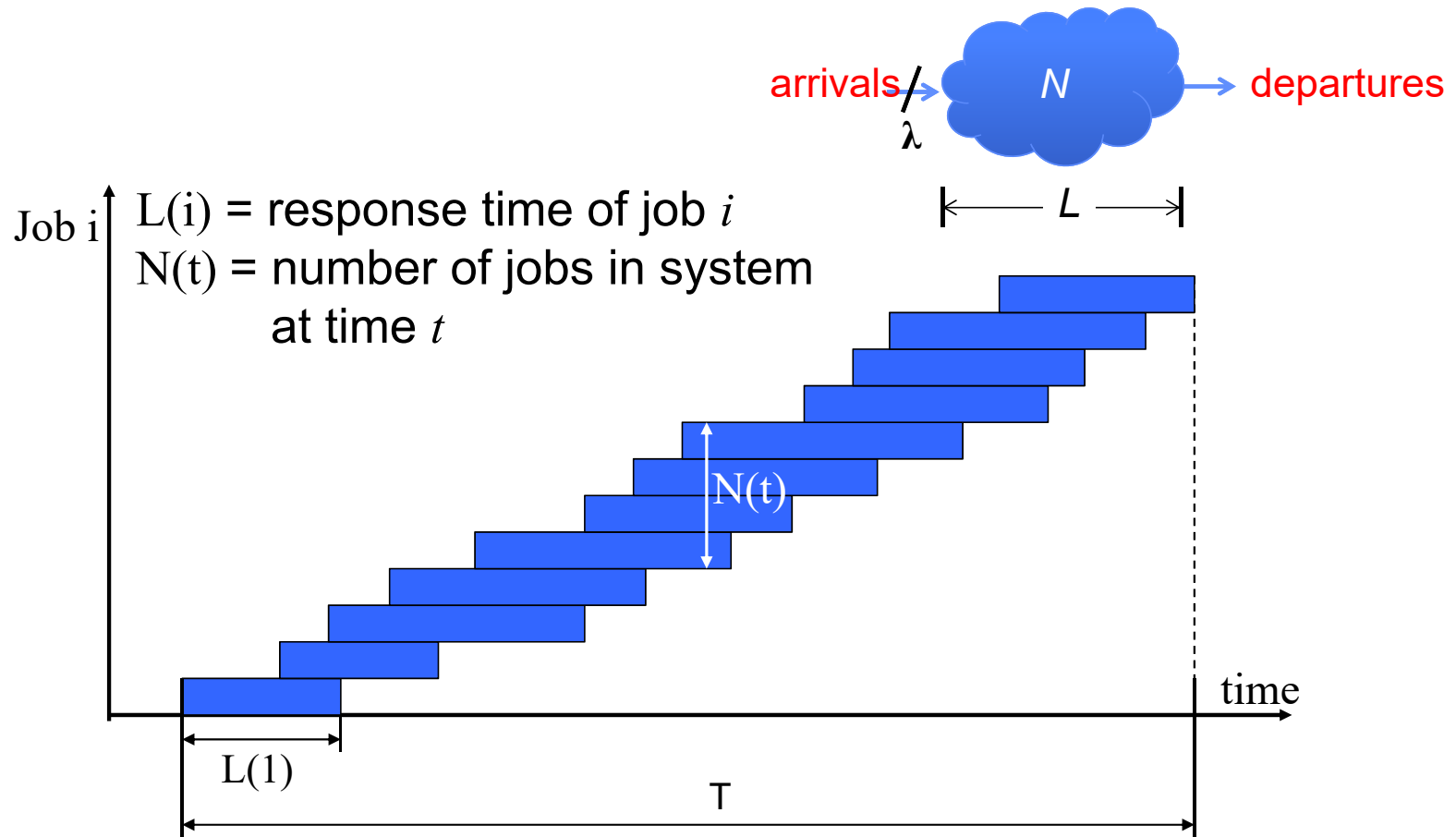
Example



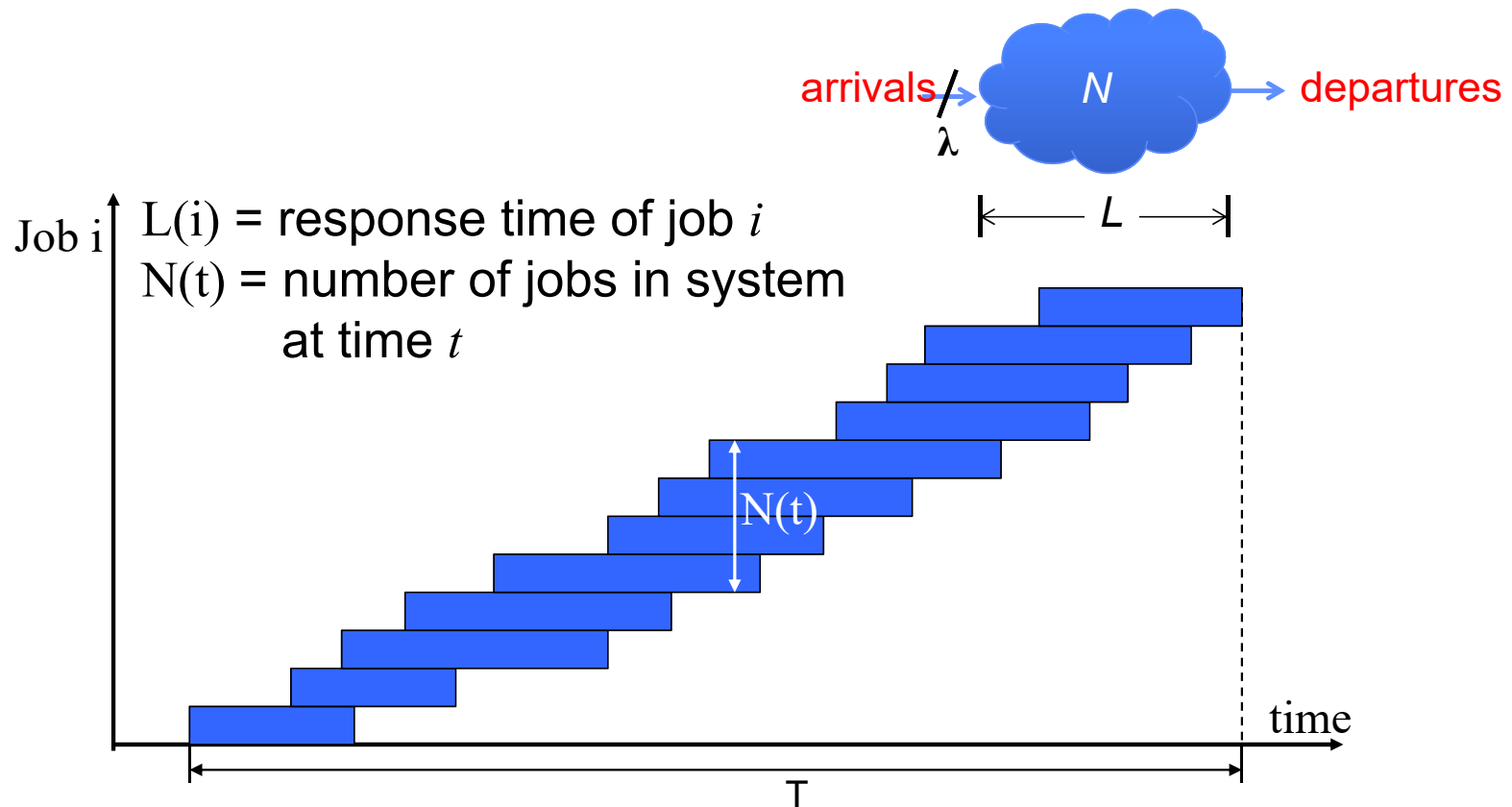
A: $N = \lambda \times L$

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

Little's Theorem: Proof Sketch

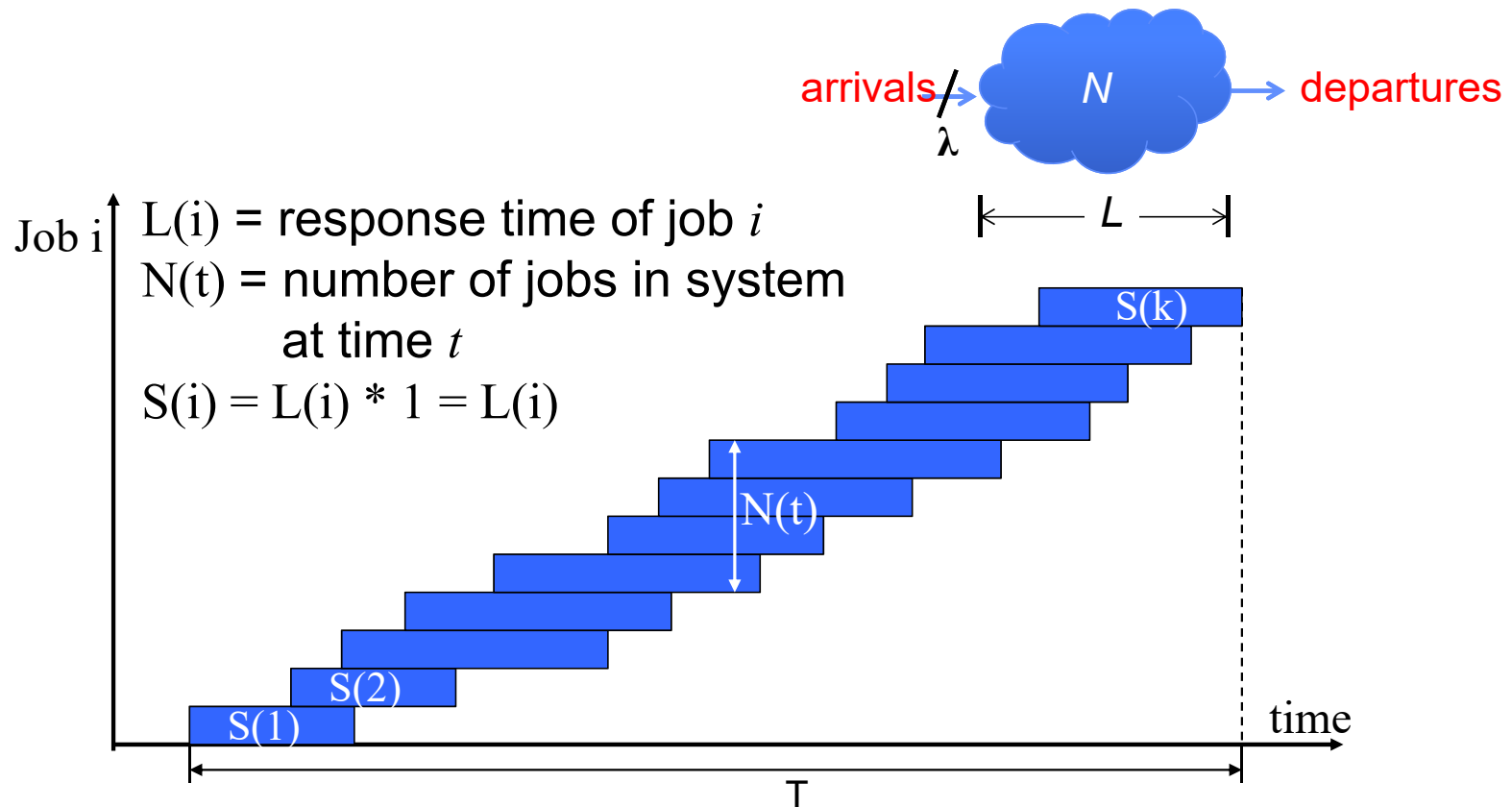


Little's Theorem: Proof Sketch



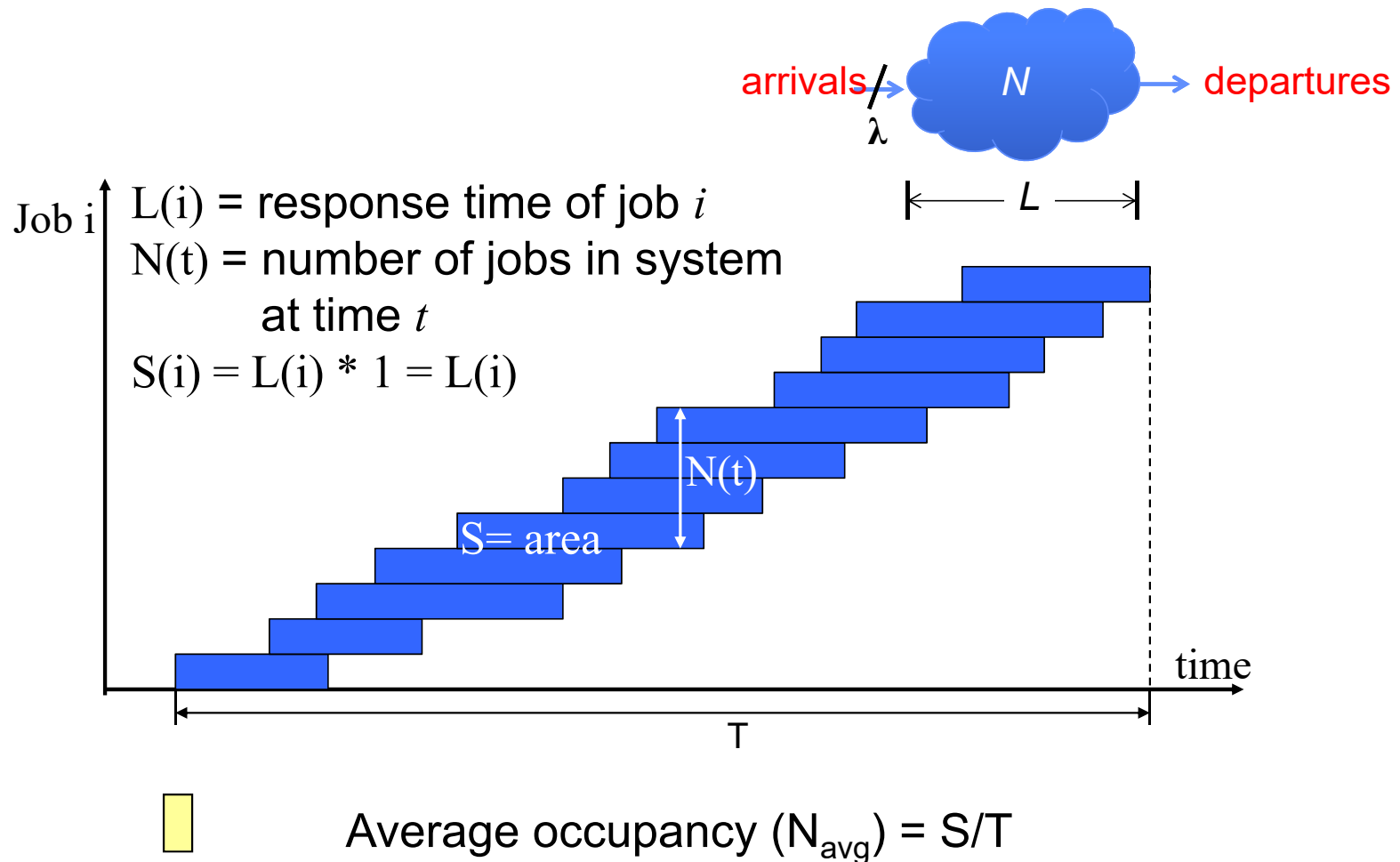
What is the system occupancy, i.e., average number of jobs in the system?

Little's Theorem: Proof Sketch

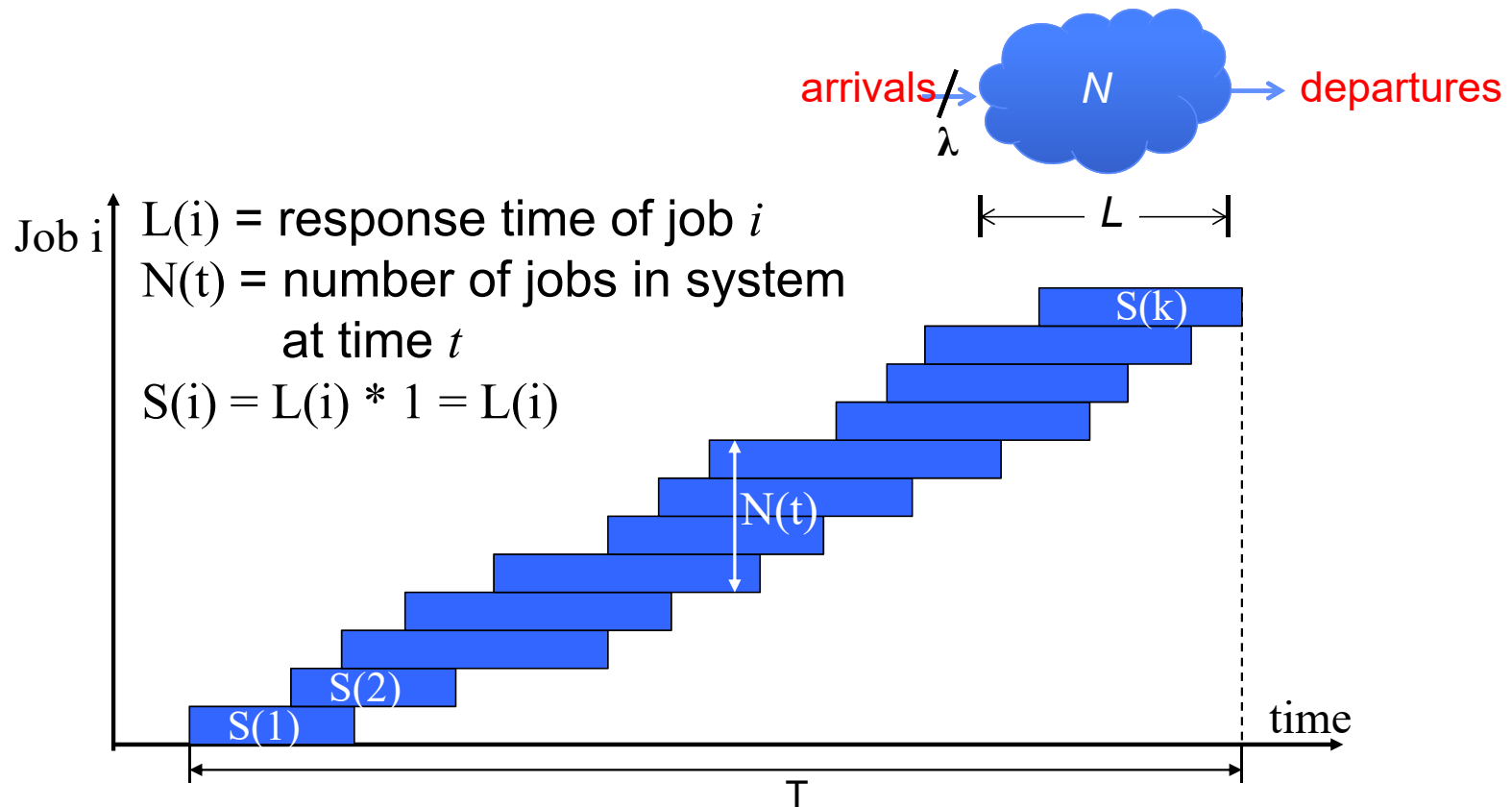


$$S = S(1) + S(2) + \dots + S(k) = L(1) + L(2) + \dots + L(k)$$

Little's Theorem: Proof Sketch

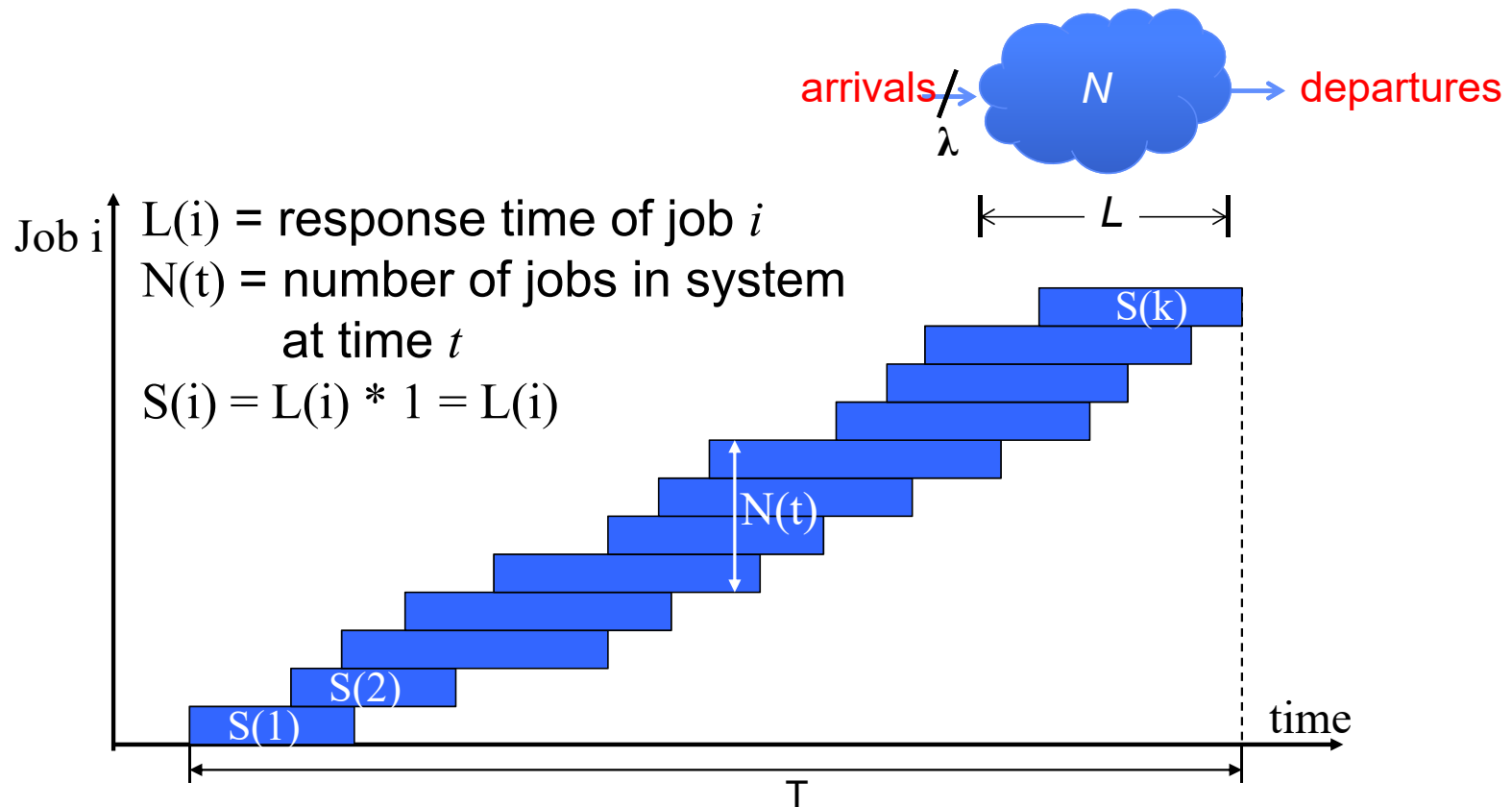


Little's Theorem: Proof Sketch



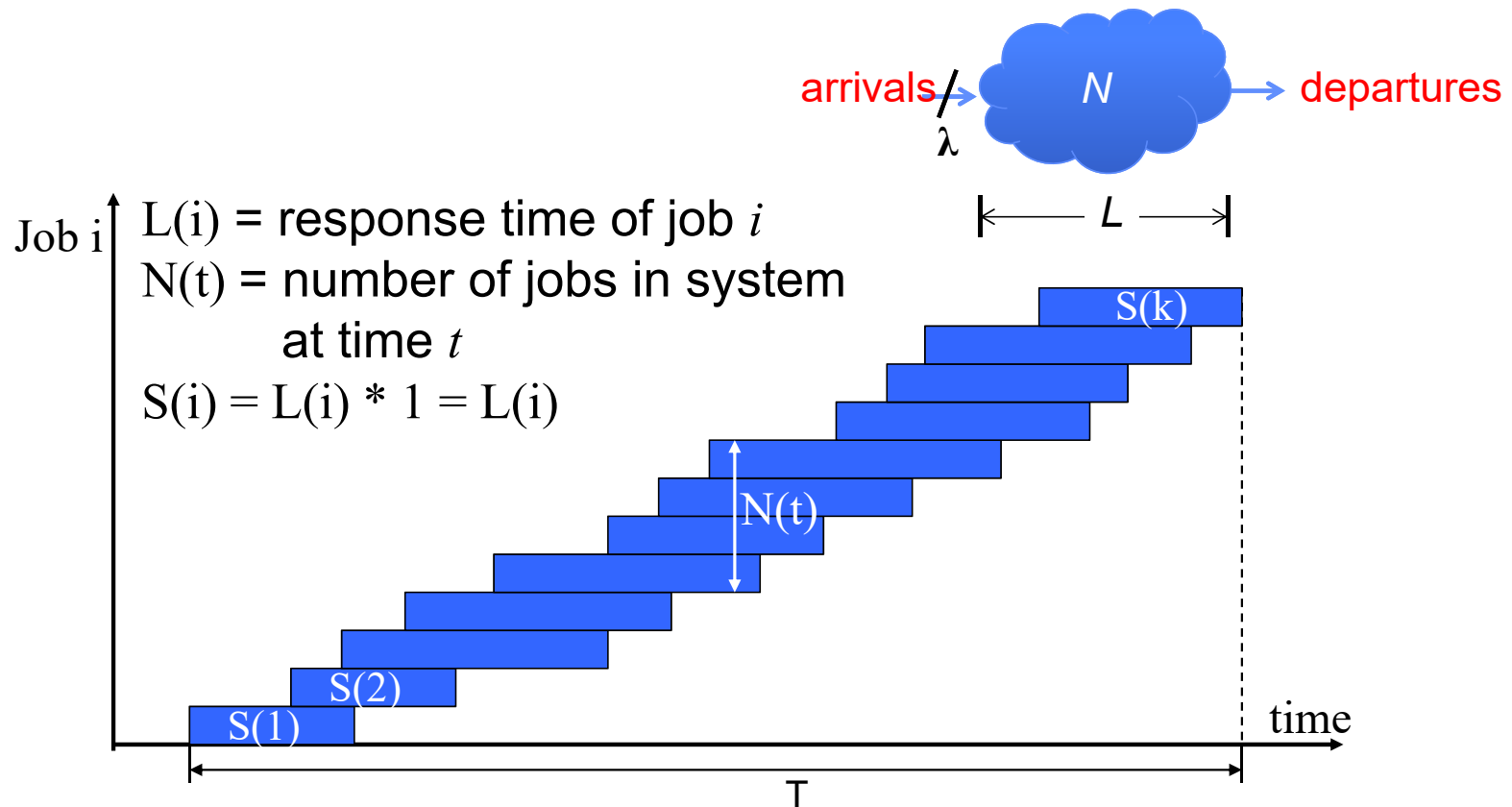
$$N_{\text{avg}} = S/T = (L(1) + \dots + L(k))/T$$

Little's Theorem: Proof Sketch



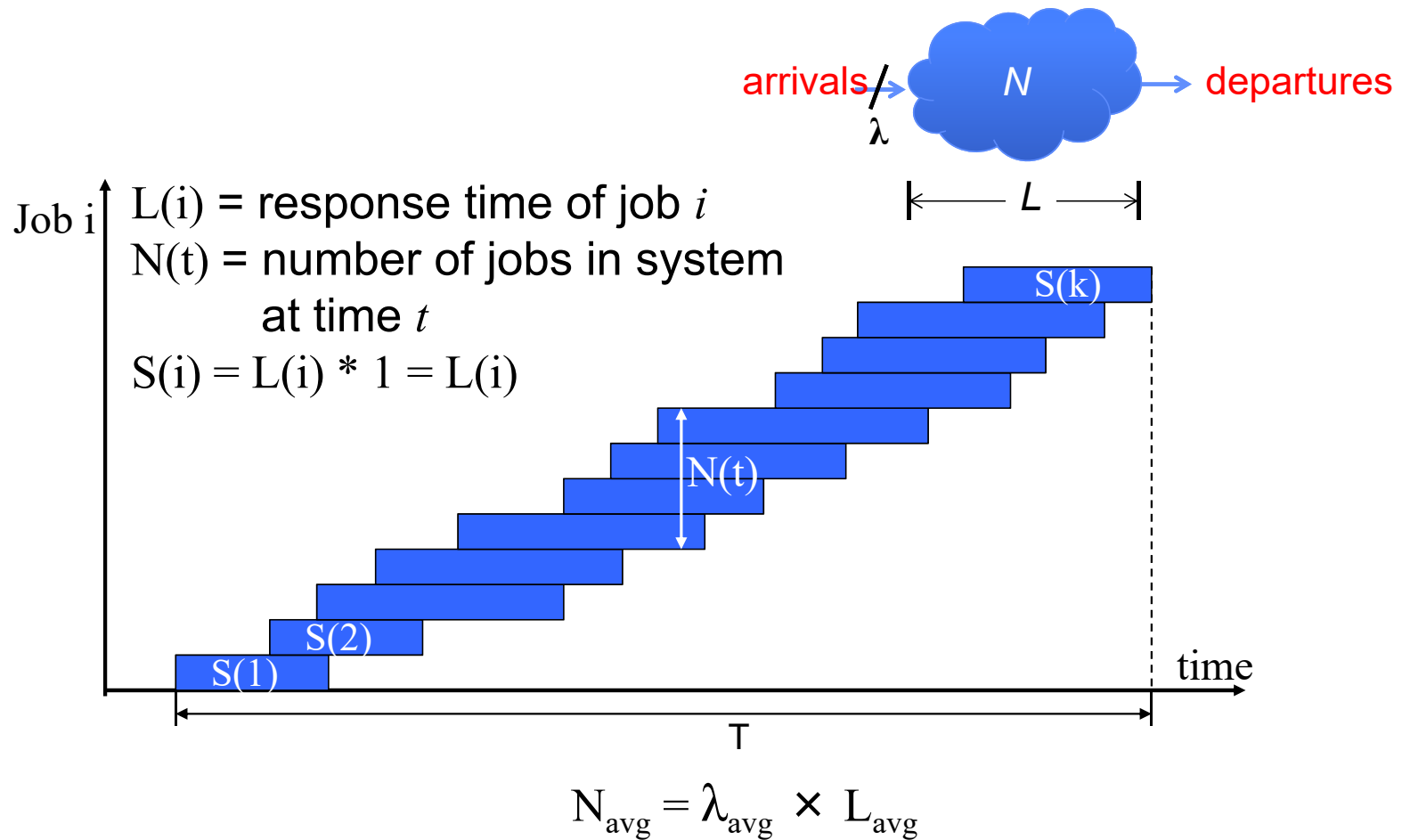
$$N_{\text{avg}} = (L(1) + \dots + L(k))/T = (N_{\text{total}}/T) * (L(1) + \dots + L(k))/N_{\text{total}}$$

Little's Theorem: Proof Sketch



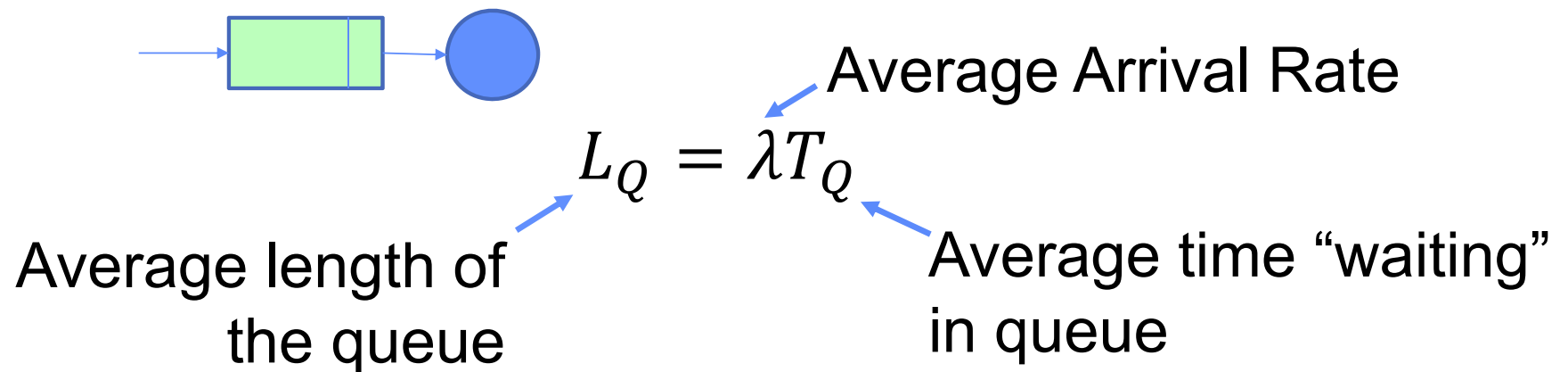
$$N_{\text{avg}} = (N_{\text{total}}/T) * (L(1) + \dots + L(k)) / N_{\text{total}} = \lambda_{\text{avg}} \times L_{\text{avg}}$$

Little's Theorem: Proof Sketch



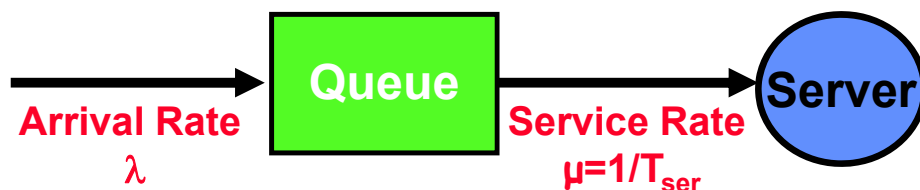
Little's Law Applied to a Queue

- When Little's Law applied to a queue, we get:

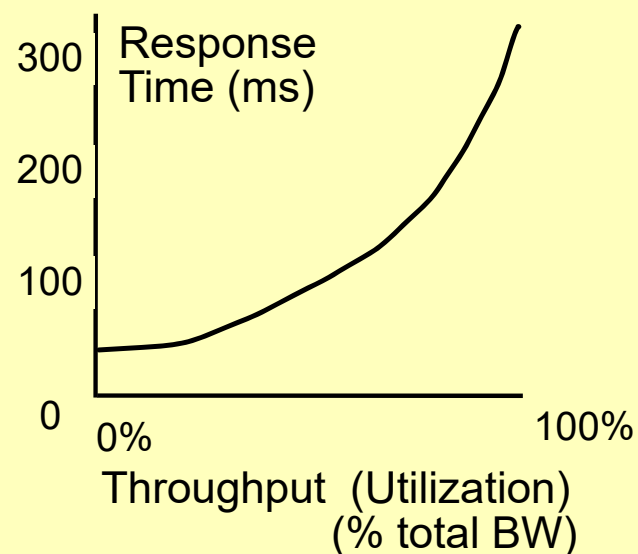


A Little Queuing Theory: Computing T_q

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



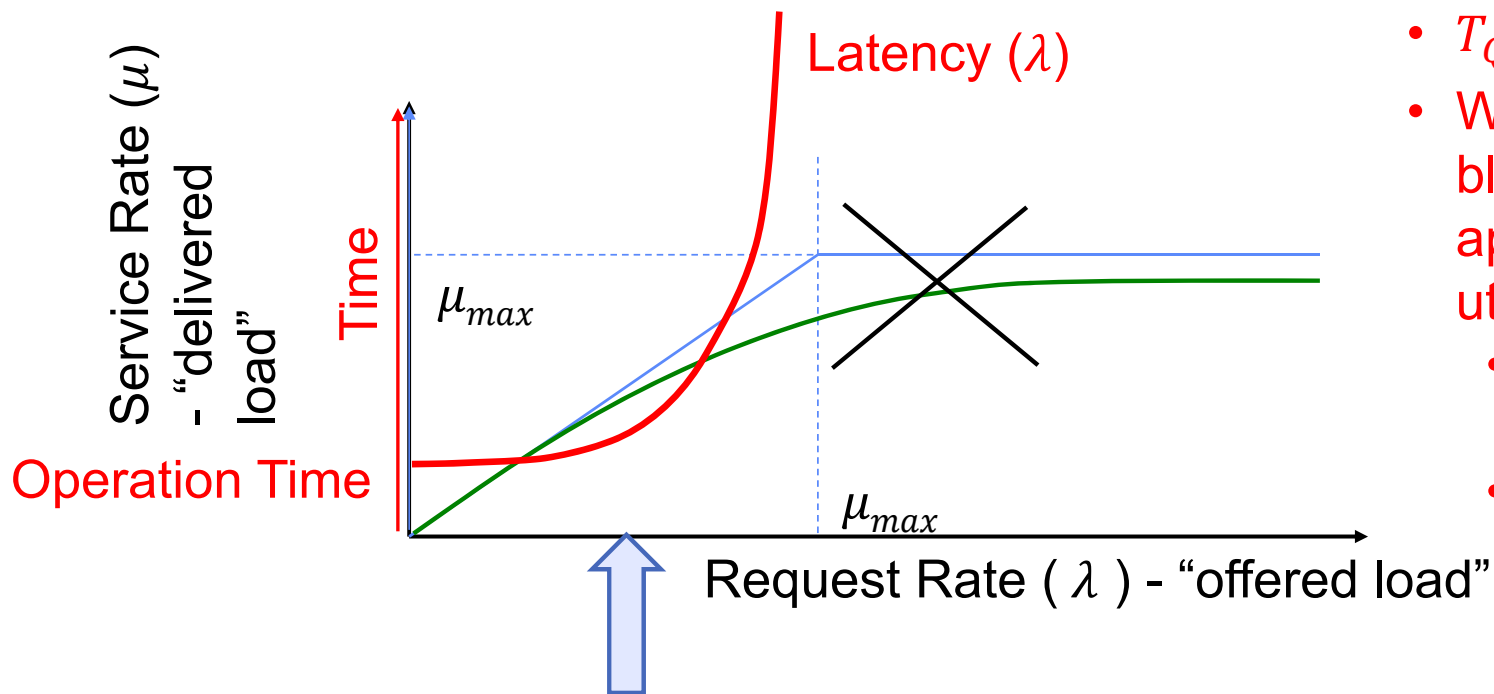
Why does response/queueing delay grow unboundedly even though the utilization is < 1 ?



- Parameters that describe our system:
 - λ : mean number of arriving customers/second
 - T_{ser} : mean time to service a customer ("m")
 - C : squared coefficient of variance σ^2/m^2
 - μ : service rate = $1/T_{ser}$
 - u : server utilization ($0 \leq u < 1$). $u = \lambda/\mu = \lambda \times T_s$

- Results:
 - Memoryless service distribution ($C = 1$): (an "M/M/1 queue"):
 - $T_q = T_{ser} \times \frac{u}{(1-u)}$
 - General service distribution server (an "M/G/1 queue"):
 - $T_q = T_{ser} \times \frac{1}{2}(1+C) \times \frac{u}{(1-u)}$

System Performance In presence of a Queue



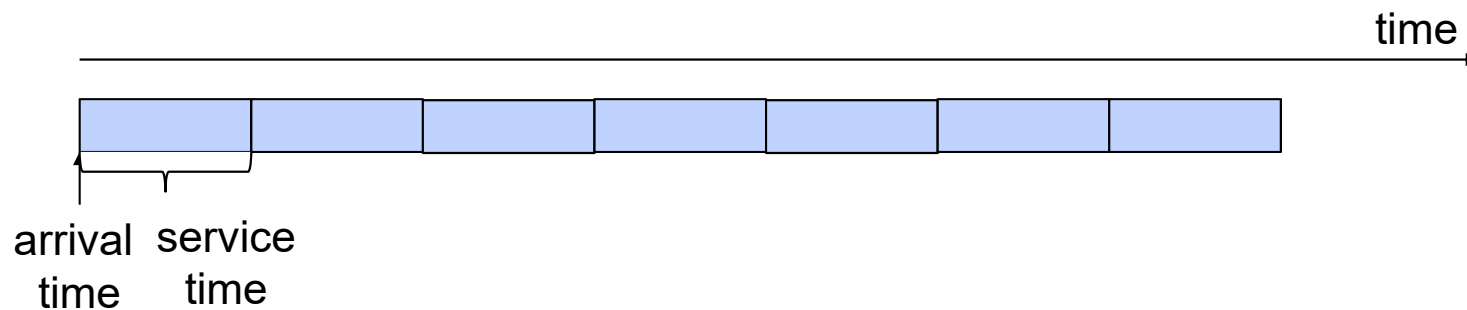
- $T_Q \sim \frac{u}{1-u}$, $u = \lambda / \mu_{max}$
- Why does latency blow up as we approach 100% utilization?
 - Queue builds up on each burst
 - But very rarely (or never) gets a chance to drain

"Half-Power Point" : load at which system delivers half of peak performance

- Design and provision systems to operate roughly in this regime
- Latency low and predictable, utilization good: ~50%

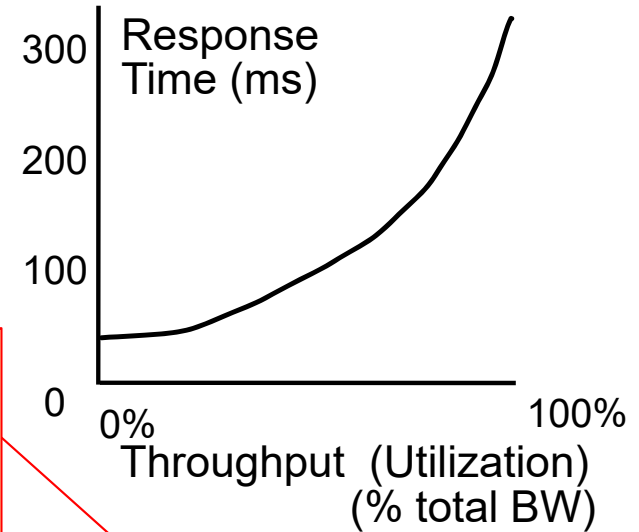
Why unbounded response time?

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

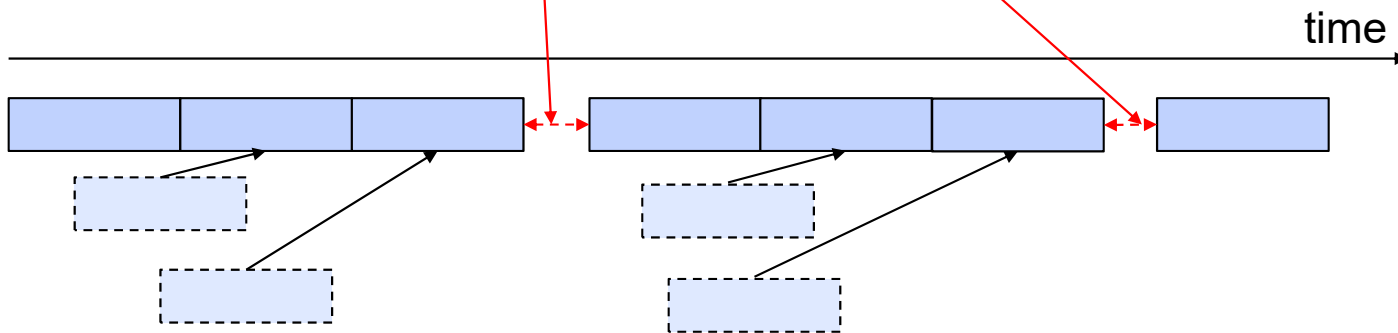


Why unbounded response time?

- 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: $T_{sys} = T_q + T_{ser}$

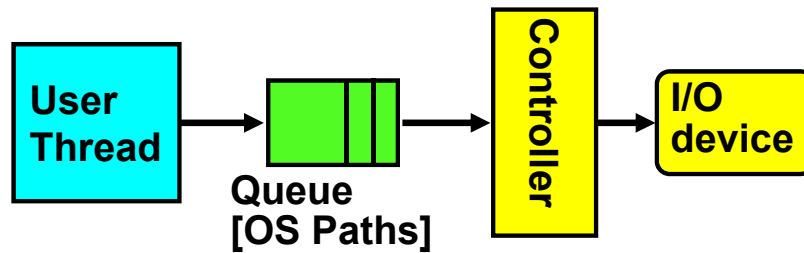
- Computation:

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

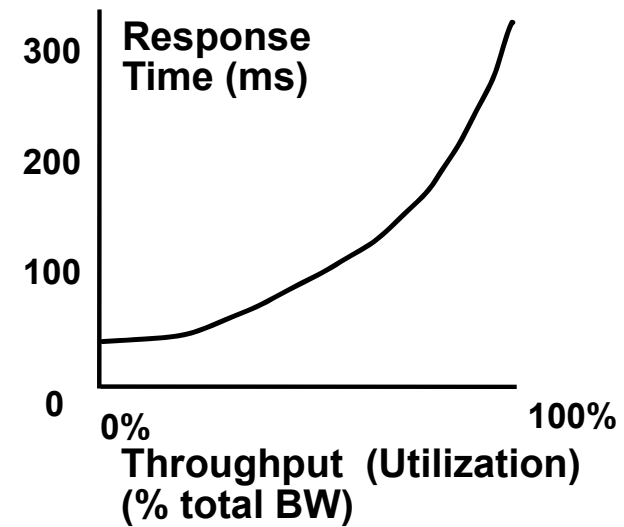
Queuing Theory Resources

- Resources page contains Queueing 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 queueing theory questions
- Assume that Queueing Theory is fair game for Midterm III!

Optimize I/O Performance



**Response Time =
Queue + I/O device service time**



- 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

Recall: I/O and Storage Layers

Application / Service

High Level I/O

Streams

Low Level I/O

File Descriptors

Syscall

*open(), read(), write(), close(), ...
Open File Descriptions*

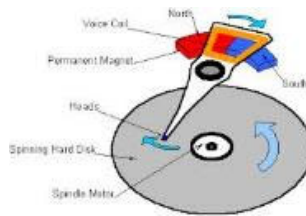
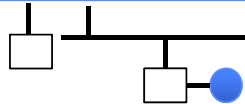
File System

Files/Directories/Indexes

I/O Driver

Commands and Data Transfers

Disks, Flash, Controllers, DMA

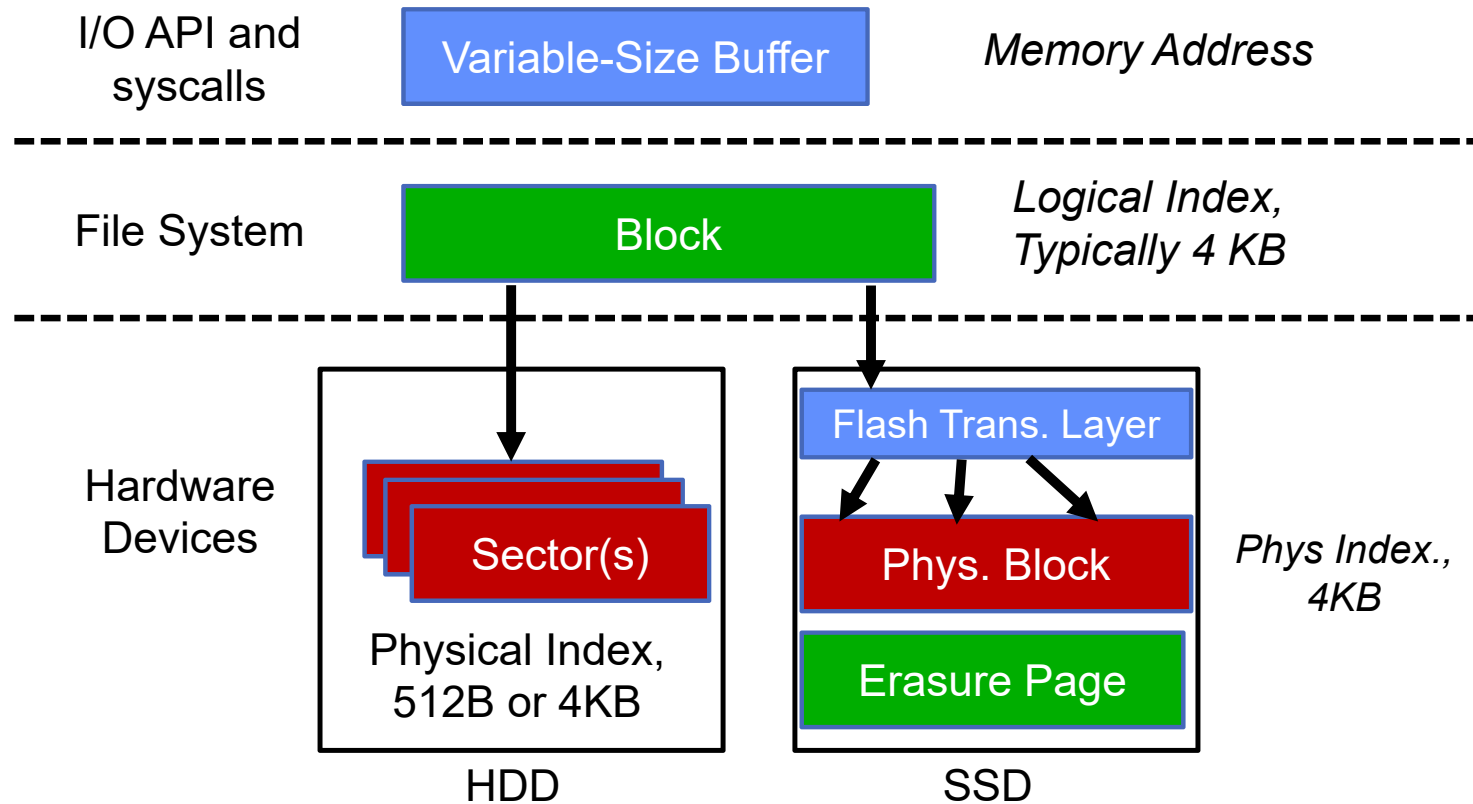


What we covered in Lecture 4

What we will cover next...

What we just covered...

From Storage to File Systems



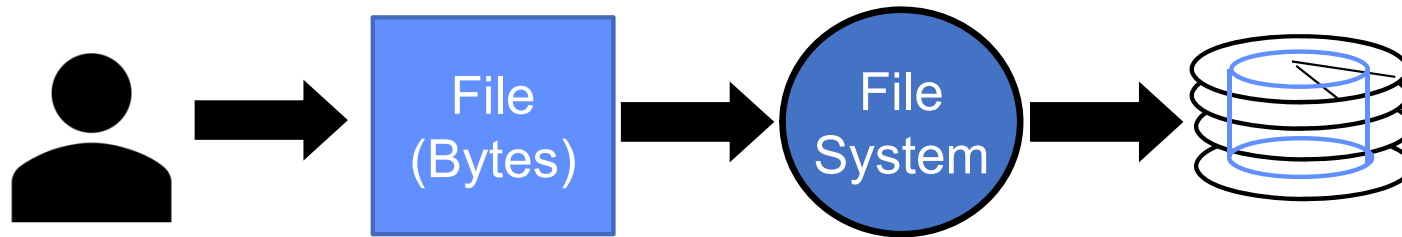
Building a File System

- **File System:** Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.
- Classic OS situation: Take limited hardware interface (array of blocks) and provide a more convenient/useful interface with:
 - Naming: Find file by name, not block numbers
 - Organize file names with directories
 - Organization: Map files to blocks
 - Protection: Enforce access restrictions
 - Reliability: Keep files intact despite crashes, hardware failures, 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

Translation 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 writing bytes 2 – 12?
 - Fetch block, modify relevant portion, write out block
- Everything inside file system is in terms of whole-size blocks
 - Actual disk I/O happens in blocks
 - read/write smaller than block size needs to translate and buffer

Disk Management

- Basic entities on a disk:
 - **File**: user-visible group of blocks arranged sequentially in logical space
 - **Directory**: user-visible index mapping names to files
- The disk is accessed as linear array of sectors
- How to identify a sector?
 - Physical position
 - » Sectors is a vector [cylinder, surface, sector]
 - » Not used anymore
 - » OS/BIOS must deal with bad sectors
 - **Logical Block Addressing (LBA)**
 - » Every sector has integer address
 - » Controller translates from address \Rightarrow physical position
 - » Shields OS from structure of disk

What Does the File System Need?

- Track free disk blocks
 - Need to know where to put newly written data
- Track which blocks contain data for which files
 - Need to know where to read a file from
- Track files in a directory
 - Find list of file's blocks given its name
- Where do we maintain all of this?
 - Somewhere on disk

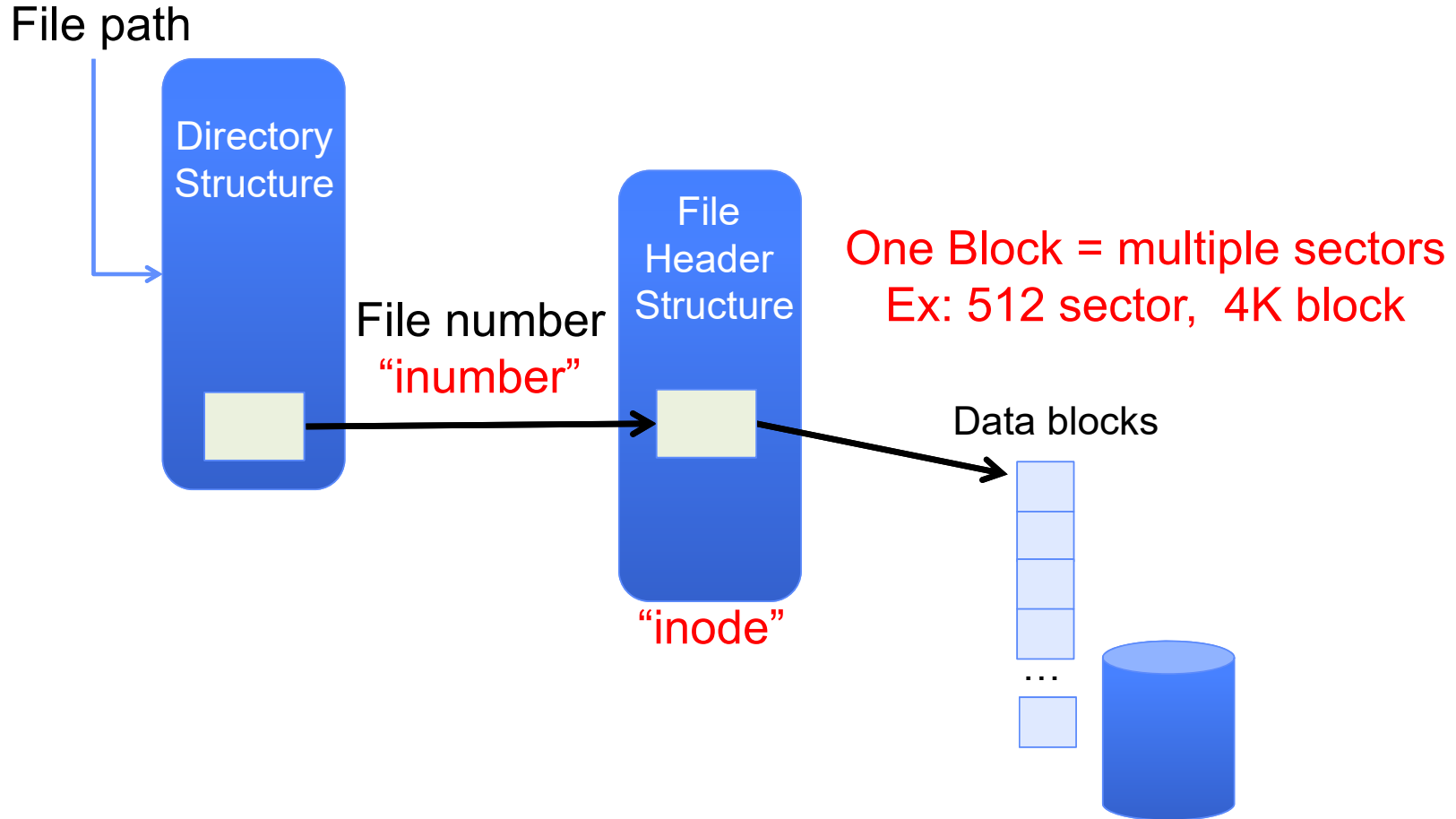
Data Structures on Disk

- Somewhat different from data structures in memory
- Access a block at a time
 - Can't efficiently read/write a single word
 - Have to read/write full block containing it
 - Ideally want sequential access patterns
- Durability
 - Ideally, file system is in meaningful state upon shutdown
 - This obviously isn't always the case...

Critical Factors in File System Design

- (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 carefully allocate / free blocks
 - Such that access remains efficient

Components of a File System



Conclusion

- Devices have complex interaction and performance characteristics
 - Response time (Latency) = Queue + Overhead + Transfer
 - » Effective BW = $BW * T/(S+T)$
 - HDD: Queuing time + controller + seek + rotation + transfer
 - SSD: Queuing time + controller + transfer (erasure & wear)
- 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