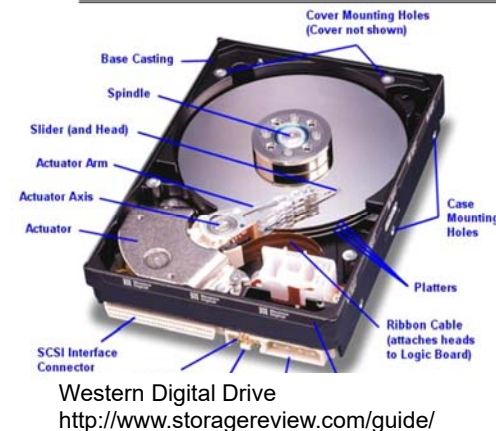## Storage Devices

- Magnetic disks
  - Storage that rarely becomes corrupted
  - Large capacity at low cost
  - Block level random access (except for SMR – later!)
  - Slow performance for random access
  - Better performance for sequential access

- Flash memory
  - Storage that rarely becomes corrupted
  - Capacity at intermediate cost (5-20x disk)
  - Block level random access
  - Good performance for reads; worse for random writes
  - Erasure requirement in large blocks
  - Wear patterns issue

## Hard Disk Drives (HDDs)



Western Digital Drive
http://www.storagereview.com/guide/

**Read/Write Head Side View**

**IBM/Hitachi Microdrive**

IBM Personal Computer/AT (1986)
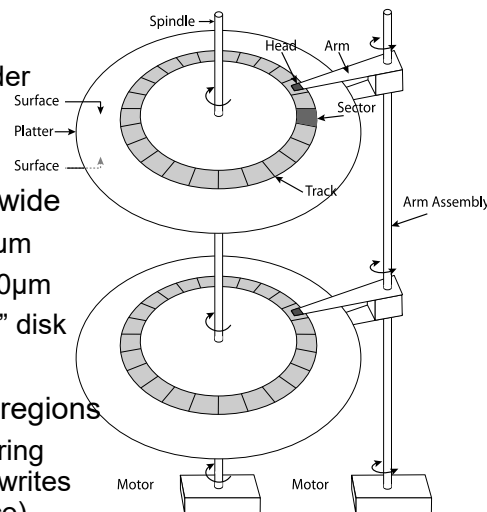  30 MB hard disk - $500
  30-40ms seek time
  0.7-1 MB/s (est.)

## The Amazing Magnetic Disk

- Unit of Transfer: Sector
  - Ring of sectors form a track
  - Stack of tracks form a cylinder
  - Heads position on cylinders

- Disk Tracks ~ 1µm (micron) wide
  - Wavelength of light is ~ 0.5µm
  - Resolution of human eye: 50µm
  - 100K tracks on a typical 2.5" disk

- Separated by unused guard regions
  - Reduces likelihood neighboring tracks are corrupted during writes (still a small non-zero chance)
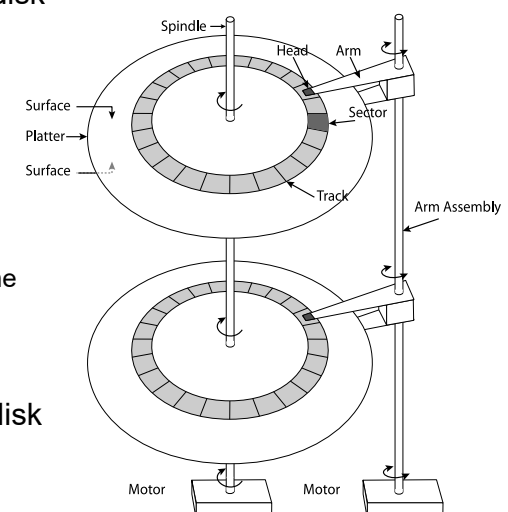
## The Amazing Magnetic Disk

- Track length varies across disk
  - Outside: More sectors per track, higher bandwidth
  - Disk is organized into regions of tracks with same # of sectors/track
  - Only outer half of radius is used
    » Most of the disk area in the outer regions of the disk

- Disks so big that some companies (like Google) reportedly only use part of disk for active data
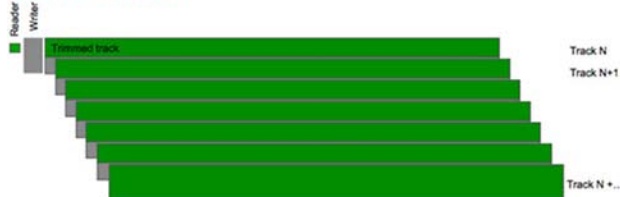  - Rest is archival data

## Shingled Magnetic Recording (SMR)

**Conventional Writes**



**SMR Writes**



- Overlapping tracks yields greater density, capacity
- Restrictions on writing, complex DSP for reading
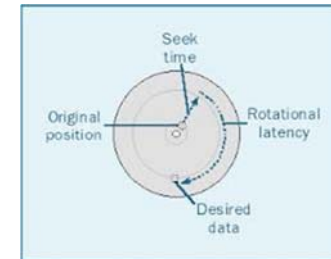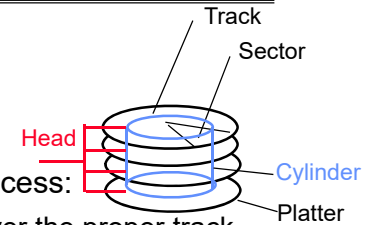- Examples: Seagate (8TB), Hitachi (10TB)

## Review: Magnetic Disks

- Cylinders: all the tracks under the head at a given point on all surface



- 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



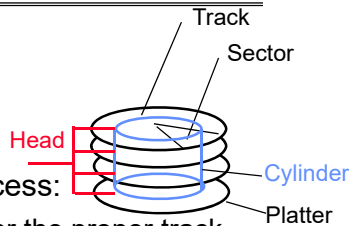**Seek time = 4-8ms**
**One rotation = 1-2ms**
**(3600-7200 RPM)**

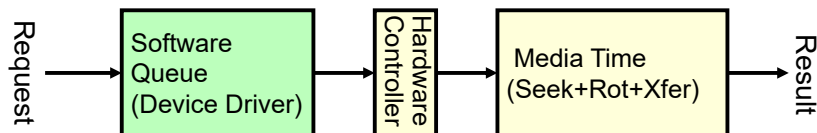## Review: Magnetic Disks

- Cylinders: all the tracks under the head at a given point on all surface



- 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

**Disk Latency = Queueing Time + Controller time + Seek Time + Rotation Time + Xfer Time**

## Typical Numbers for Magnetic Disk

| Parameter | Info / Range |
|---|---|
| Space/Density | Space: 14TB (Seagate), 8 platters, in 3½ inch form factor! Areal Density: ≥ 1Terabit/square inch! (PMR, Helium, …) |
| Average seek time | Typically 4-6 milliseconds. Depending on reference locality, actual cost may be 25-33% of this number. |
| Average rotational latency | Most laptop/desktop disks rotate at 3600-7200 RPM (16-8 ms/rotation). Server disks up to 15,000 RPM. Average latency is halfway around disk so 8-4 milliseconds |
| Controller time | Depends on controller hardware |
| Transfer time | Typically 50 to 250 MB/s. Depends on:<br>• Transfer size (usually a sector): 512B – 1KB per sector<br>• Rotation speed: 3600 RPM to 15000 RPM<br>• Recording density: bits per inch on a track<br>• Diameter: ranges from  1 in to 5.25 in |
| Cost | Used to drop by a factor of two every 1.5 years (or even faster); now slowing down |

## Disk Performance Example

- Assumptions:
  - Ignoring queuing and controller times for now
  - Avg seek time of 5ms,
  - 7200RPM $\Rightarrow$ Time for rotation: 60000 (ms/min) / 7200(rev/min) ~= 8ms
  - Transfer rate of 50MByte/s, block size of 4Kbyte $\Rightarrow$
    4096 bytes/$50\times10^6$ (bytes/s) = $81.92 \times 10^{-6}$ sec $\cong$ 0.082 ms for 1 sector
- Read block from random place on disk:
  - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.082ms) = 9.082ms
  - Approx 9ms to fetch/put data: 4096 bytes/$9.082\times10^{-3}$ s $\cong$ 451KB/s
- Read block from random place in same cylinder:
  - Rot. Delay (4ms) + Transfer (0.082ms) = 4.082ms
  - Approx 4ms to fetch/put data: 4096 bytes/$4.082\times10^{-3}$ s $\cong$ 1.03MB/s
- Read next block on same track:
  - Transfer (0.082ms): 4096 bytes/$0.082\times10^{-3}$ s $\cong$ 50MB/sec
- Key to using disk effectively (especially for file systems) is to minimize seek and rotational delays
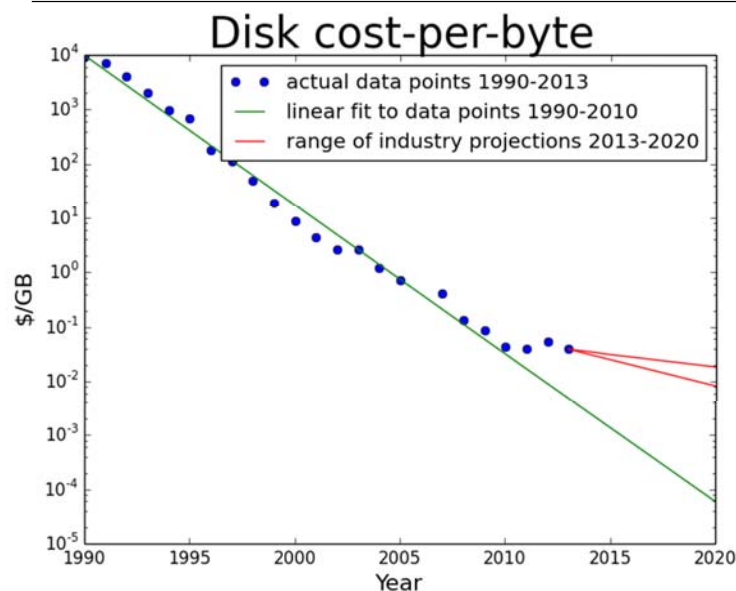
## (Lots of) Intelligence in the Controller

- Sectors contain sophisticated error correcting codes
  - Disk head magnet has a field wider than track
  - Hide corruptions due to neighboring track writes

- Sector sparing
  - Remap bad sectors transparently to spare sectors on the same surface

- Slip sparing
  - Remap all sectors (when there is a bad sector) to preserve sequential behavior

- Track skewing
  - Sector numbers offset from one track to the next, to allow for disk head movement for sequential ops

- …

## Hard Drive Prices over Time

## Example of Current HDDs

- Seagate Exos X14 (2018)
  - 14 TB hard disk
    » 8 platters, 16 heads
    » Helium filled: reduce friction and power
  - 4.16ms average seek time
  - 4096 byte physical sectors
  - 7200 RPMs
  - 6 Gbps SATA /12Gbps SAS interface
    » 261MB/s MAX transfer rate
    » Cache size: 256MB
  - Price: $615 (< $0.05/GB)

- IBM Personal Computer/AT (1986)
  - 30 MB hard disk
  - 30-40ms seek time
  - 0.7-1 MB/s (est.)
  - Price: $500 ($17K/GB, 340,000x more expensive !!)
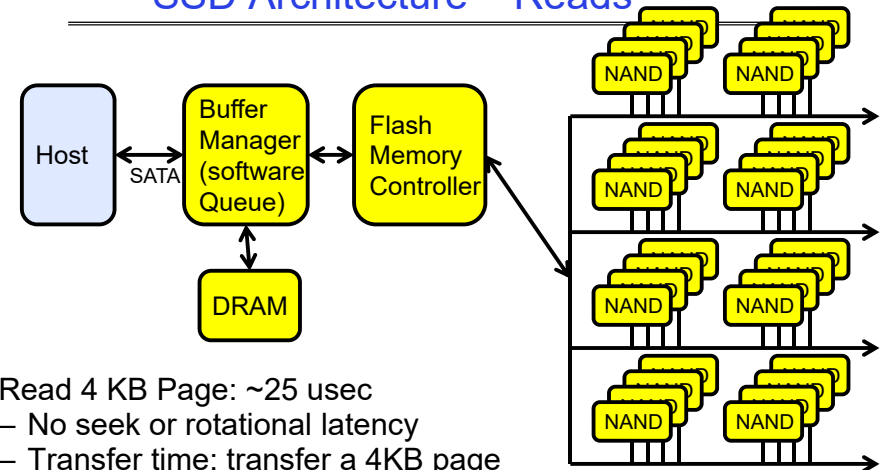
# Solid State Disks (SSDs)



- 1995 – Replace rotating magnetic media with non-volatile memory (battery backed DRAM)
- 2009 – Use NAND Multi-Level Cell (2 or 3-bit/cell) flash memory
  - Sector (4 KB page) addressable, but stores 4-64 "pages" per memory block
  - Trapped electrons distinguish between 1 and 0
- No moving parts (no rotate/seek motors)
  - Eliminates seek and rotational delay (0.1-0.2ms access time)
  - Very low power and lightweight
  - Limited "write cycles"
- Rapid advances in capacity and cost ever since!

# SSD Architecture – Reads



Read 4 KB Page: ~25 usec
- No seek or rotational latency
- Transfer time: transfer a 4KB page
  - » SATA: 300-600MB/s => ~4 x$10^3$ b / 400 x $10^6$ bps => 10 us
- Latency = Queuing Time + Controller time + Xfer Time
- Highest Bandwidth: Sequential OR Random reads

# SSD Architecture – Writes

- Writing data is complex! (~200µs – 1.7ms )
  - Can only write empty pages in a block
  - 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



Data written in 4 KB Pages

4 KB | 4 KB | 4 KB
4 KB | 4 KB | 4 KB

Data erased in 256 KB Blocks

64 writable Pages in 1 erasable Block

4 KB | 4 KB | 4 KB

Typical NAND Flash Pages and Blocks

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

# Some "Current" 3.5in SSDs

- Seagate Nytro SSD: 15TB (2017)
  - Dual 12Gb/s interface
  - Seq reads 860MB/s
  - Seq writes 920MB/s
  - Random Reads (IOPS): 102K
  - Random Writes (IOPS): 15K
  - Price (Amazon): $6325 ($0.41/GB)



- Nimbus SSD: 100TB (2019)
  - Dual port: 12Gb/s interface
  - Seq reads/writes: 500MB/s
  - Random Read Ops (IOPS): 100K
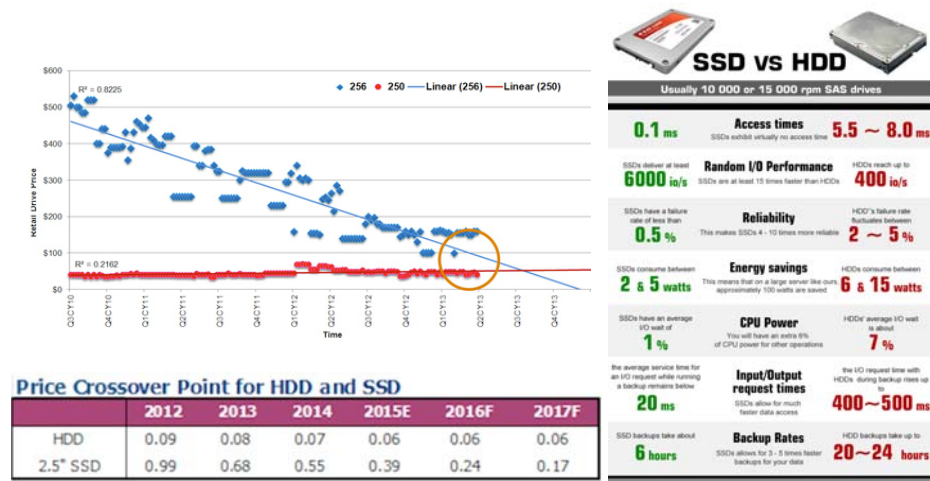  - *Unlimited writes for 5 years!*
  - Price: ~ $50K? ($0.50/GB)

## HDD vs SSD Comparison



**Price Crossover Point for HDD and SSD**

| | 2012 | 2013 | 2014 | 2015E | 2016F | 2017F |
|---|---|---|---|---|---|---|
| HDD | 0.09 | 0.08 | 0.07 | 0.06 | 0.06 | 0.06 |
| 2.5" SSD | 0.99 | 0.68 | 0.55 | 0.39 | 0.24 | 0.17 |

### SSD prices drop much faster than HDD

---

## 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
  - ½ 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 Kubiatowicz (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

**No longer true!**

- Cons
  - Small storage (0.1-0.5x disk), expensive (3-20x 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!
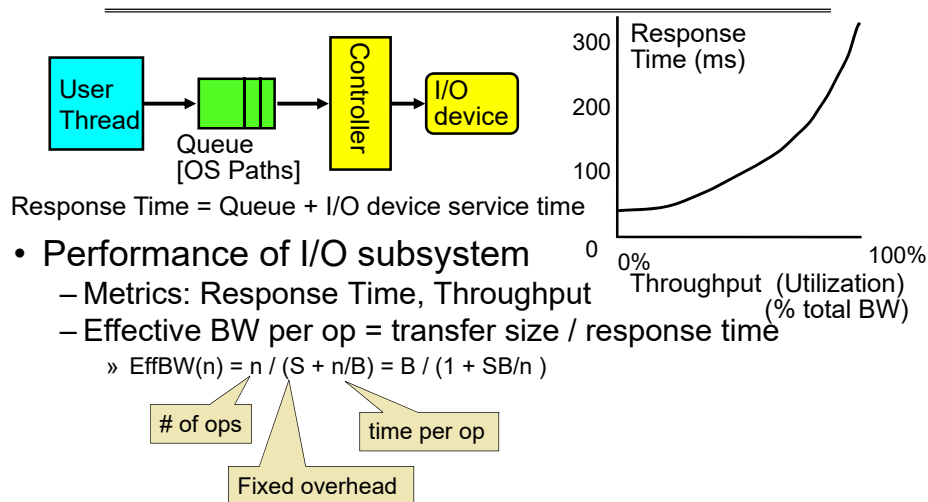
## I/O Performance



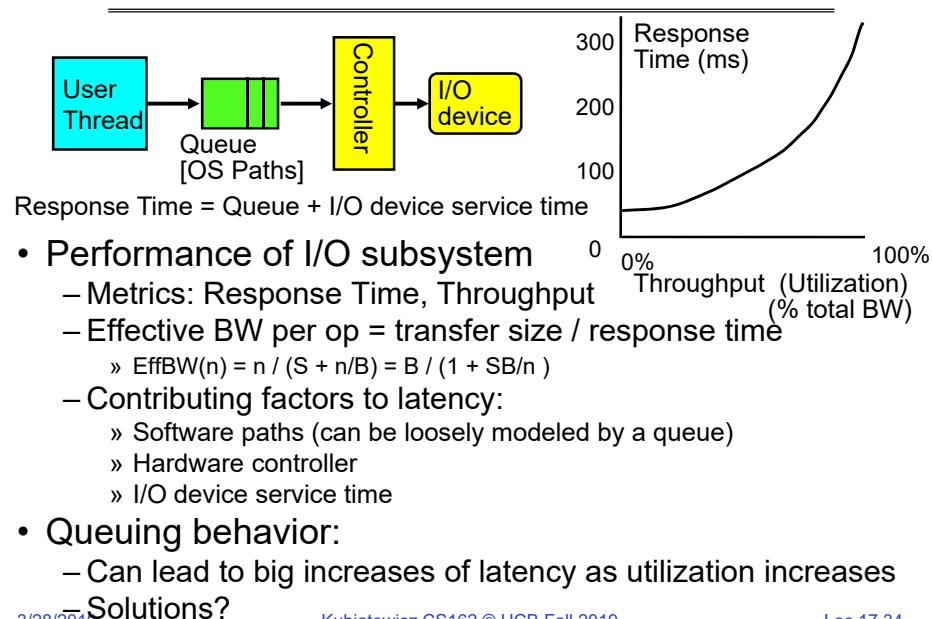Response Time = Queue + I/O device service time

- **Performance of I/O subsystem**
  - Metrics: Response Time, Throughput
  - Effective BW per op = transfer size / response time
    - » $EffBW(n) = n / (S + n/B) = B / (1 + SB/n)$

| # of ops | time per op |

Fixed overhead

## 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
    - » $EffBW(n) = n / (S + n/B) = B / (1 + SB/n)$
  - 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?
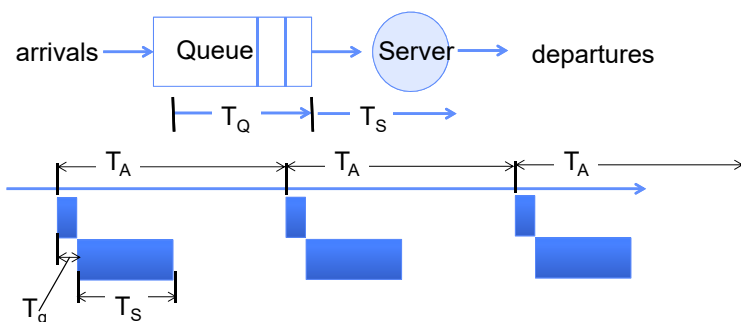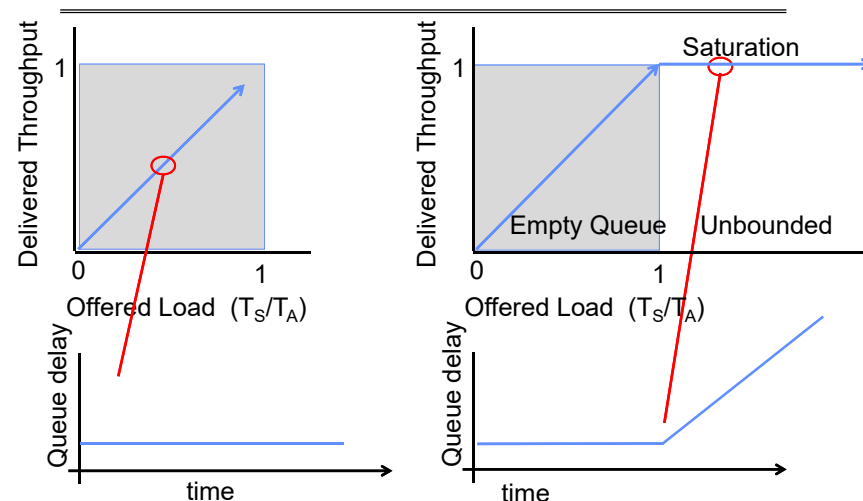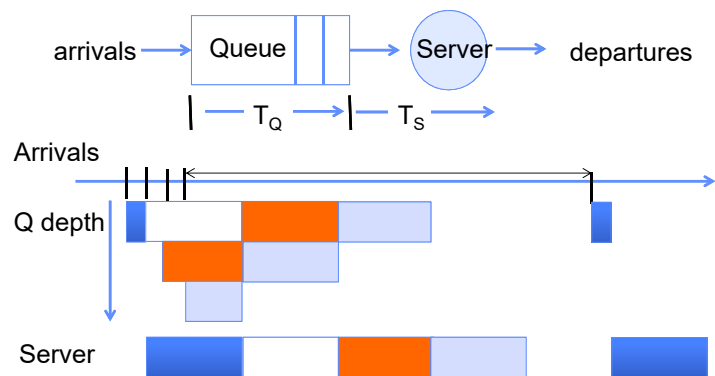
## A Simple Deterministic World



- 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 unbounded at a rate ~ $(T_s/T_A)$ 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
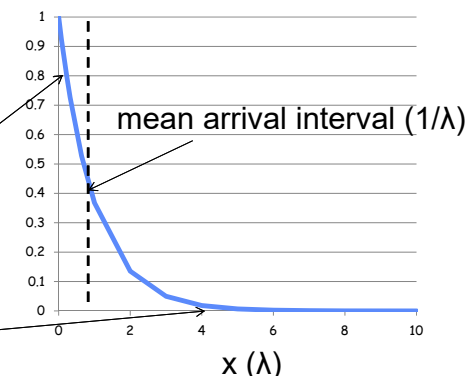- Even though average utilization is low

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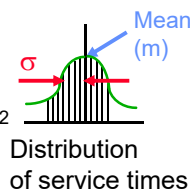
Few long gaps (i.e., low instantaneous rate)
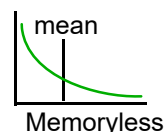


mean arrival interval $(1/\lambda)$

$x (\lambda)$

## Background: General Use of Random Distributions

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



Distribution of service times

- 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
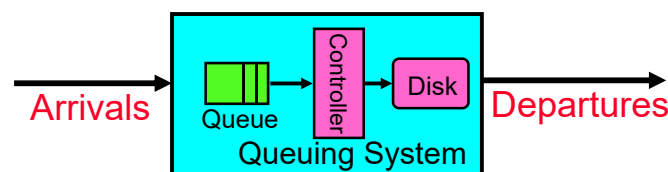  - Disk response times $C \approx 1.5$ (majority seeks < average)

mean

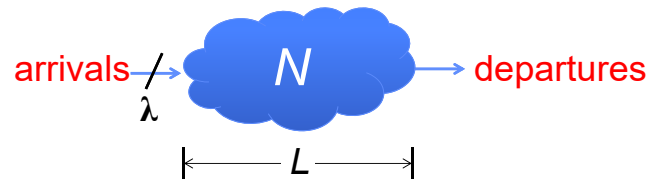Memoryless

## 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


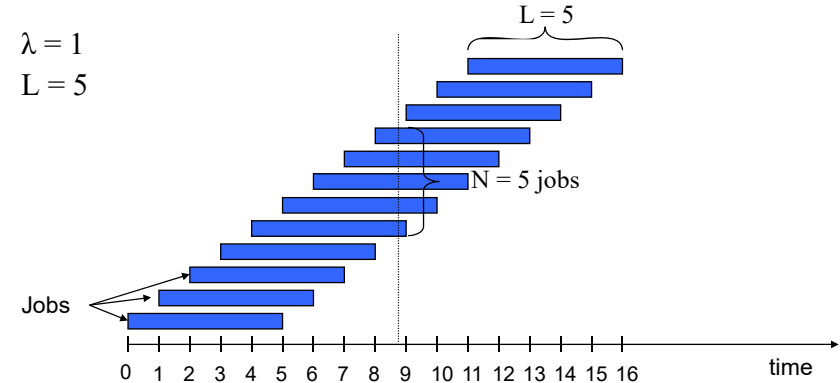
arrivals → **N** → departures
λ

L

- 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) x 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

λ = 1
L = 5

L = 5

N = 5 jobs

Jobs

0  1  2  3  4  5  6  7  8  9  10 11 12 13 14 15 16    time

**A:** N = λ x L

- E.g., N = λ x L = 5

---

## Little's Theorem: Proof Sketch

arrivals → **N** → departures
λ

L

Job i

L(i) = response time of job *i*
N(t) = number of jobs in system at time *t*

N(t)

L(1)          T          time

---

## Little's Theorem: Proof Sketch

arrivals → **N** → departures
λ

L

Job i

L(i) = response time of job *i*
N(t) = number of jobs in system at time *t*

N(t)

T          time

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