

CC5212-1  
 PROCESAMIENTO MASIVO DE DATOS  
 OTOÑO 2015

Lecture 2: Distributed Systems I

Aidan Hogan  
 aidhog@gmail.com

MASSIVE DATA NEEDS  
 DISTRIBUTED SYSTEMS ...

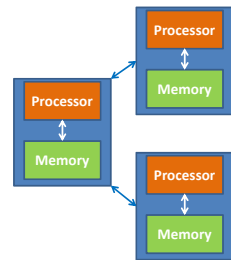
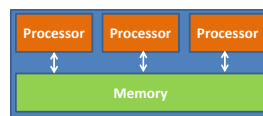
Monolithic vs. Distributed Systems

- One machine that's  $n$  times as powerful?
- vs.**
- $n$  machines that are equally as powerful?



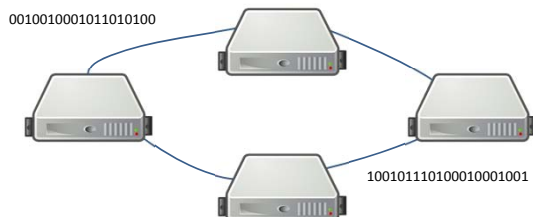
Parallel vs. Distributed Systems

- Parallel System  
 – often = *shared memory*
- Distributed System  
 – often = *shared nothing*



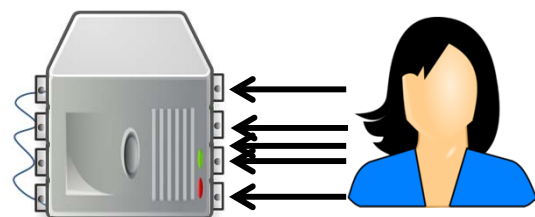
What is a Distributed System?

"A distributed system is a system that enables a collection of independent computers to communicate in order to solve a common goal."



What is a Distributed System?

"An *ideal* distributed system is a system that makes a collection of independent computers look like one computer (to the user)."



## Disadvantages of Distributed Systems

### Advantages

- **Cost**
  - Better performance/price
- **Extensibility**
  - Add another machine!
- **Reliability**
  - No central point of failure!
- **Workload**
  - Balance work automatically
- **Sharing**
  - Remote access to services

### Disadvantages

- **Software**
  - Need specialised programs
- **Networking**
  - Can be slow
- **Maintenance**
  - Debugging sw/hw a pain
- **Security**
  - Multiple users
- **Parallelisation**
  - Not always applicable

## WHAT MAKES A GOOD DISTRIBUTED SYSTEM?

## Distributed System Design

*“An **ideal** distributed system is a system that makes a collection of independent computers look like one computer (to the user).”*

- **Transparency: Abstract/hide:**
  - **Access:** How different machines are accessed
  - **Location:** What machines have what/if they move
  - **Concurrency:** Access by several users
  - **Failure:** Keep it a secret from the user

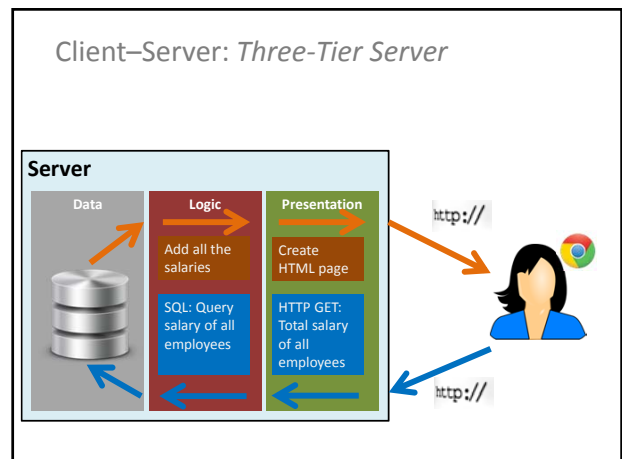
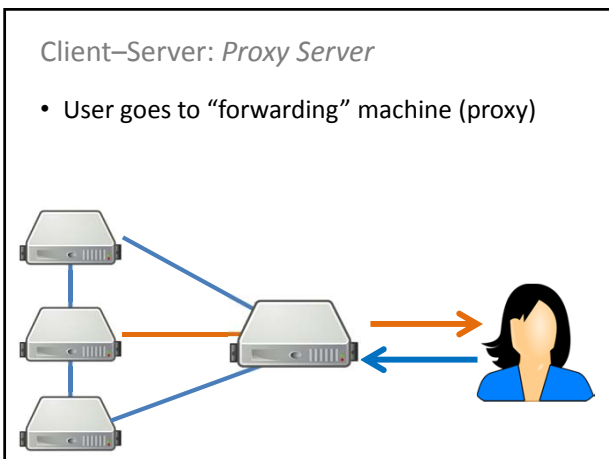
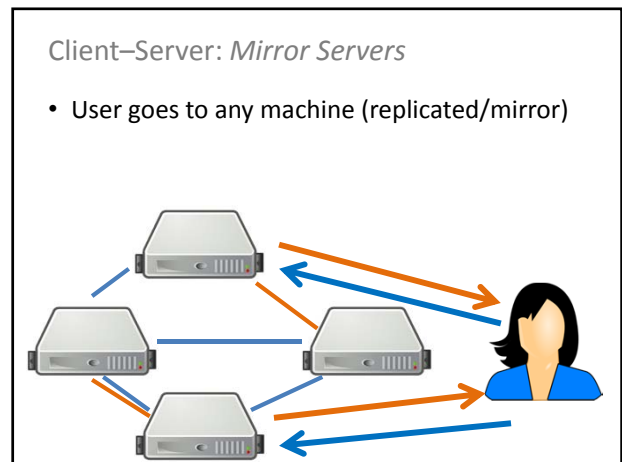
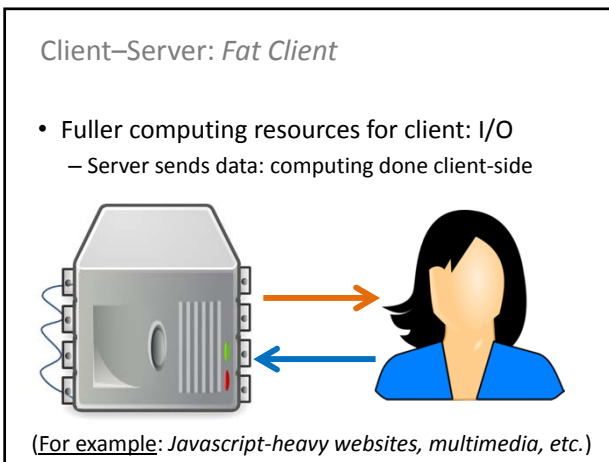
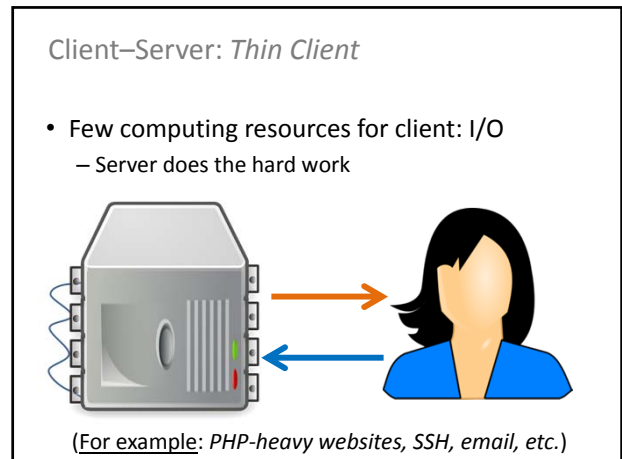
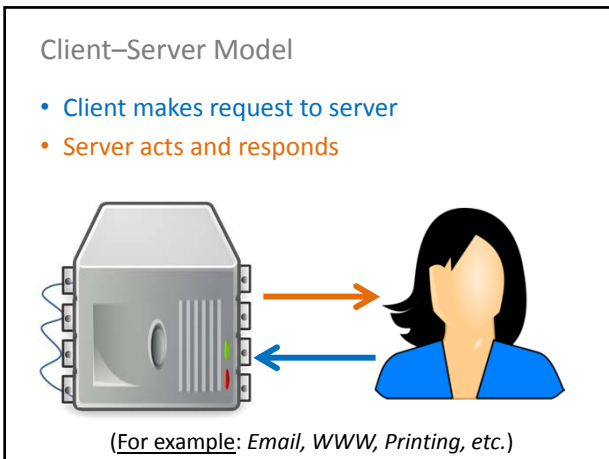
## Distributed System Design

- **Flexibility:**
  - Add/remove/move machines
  - Generic interfaces
- **Reliability:**
  - **Fault-tolerant:** recover from errors
  - **Security:** user authentication
  - **Availability:** uptime/total-time

## Distributed System Design

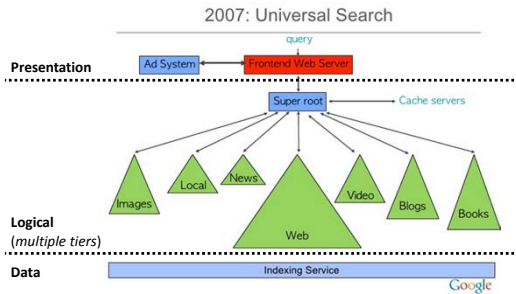
- **Performance:**
  - Runtimes (processing)
  - Latency, throughput and bandwidth (data)
- **Scalability**
  - Network and infrastructure scales
  - Applications scale
  - Minimise global knowledge/bottlenecks!

## DISTRIBUTED SYSTEMS: CLIENT–SERVER ARCHITECTURE



Client-Server: *n-Tier Server*

- Slide from Google's Jeff Dean:

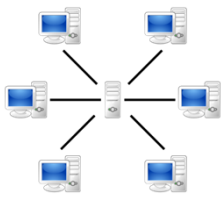


**DISTRIBUTED SYSTEMS:  
PEER-TO-PEER ARCHITECTURE**

Peer-to-Peer (P2P)

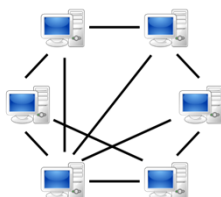
**Client-Server**

- Clients interact directly with a "central" server

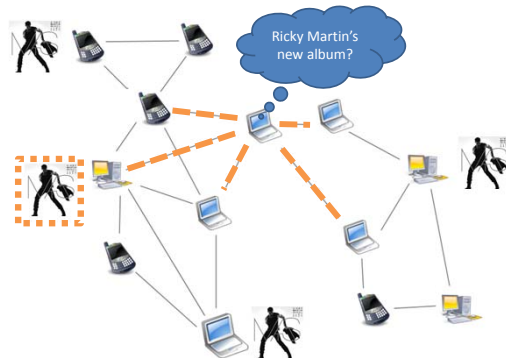


**Peer-to-Peer**

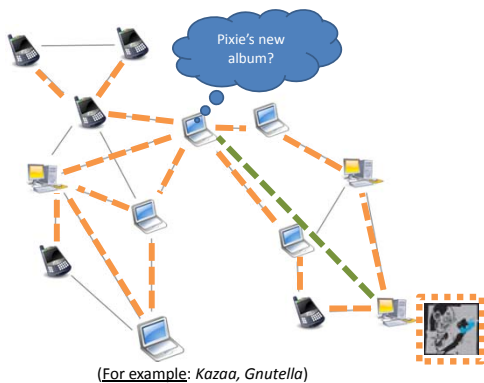
- Peers interact directly amongst themselves



Peer-to-Peer: *Unstructured*

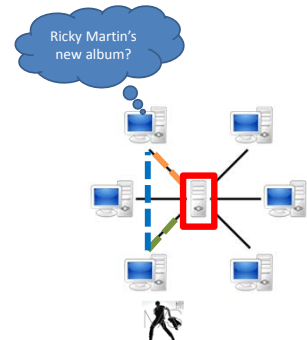


Peer-to-Peer: *Unstructured*



Peer-to-Peer: *Structured (Central)*

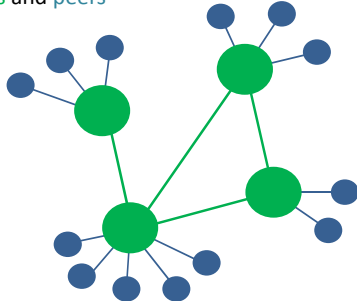
- In central server, each peer registers
  - Content
  - Address
- Peer requests content from server
- Peers connect directly
- Central point-of-failure



(For example: *Napster ... central directory was shut down*)

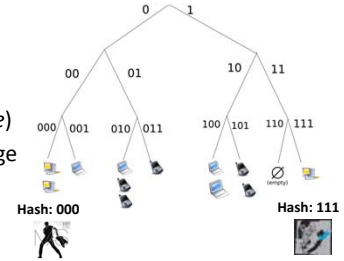
### Peer-to-Peer: Structured (Hierarchical)

- Super-peers and peers



### Peer-to-Peer: Structured (DHT)

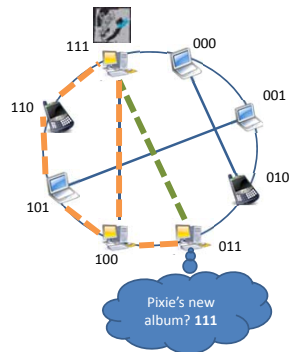
- Distributed Hash Table
- (key,value) pairs
- key based on hash
- Query with key
- Insert with (key,value)
- Peer indexes key range



(For example: BitTorrent's Tracker)

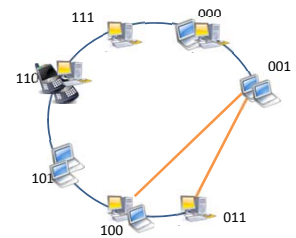
### Peer-to-Peer: Structured (DHT)

- Circular DHT:
  - Only aware of neighbours
  - $O(n)$  lookups
- Implement shortcuts
  - Skips ahead
  - Enables binary-search-like behaviour
  - $O(\log(n))$  lookups



### Peer-to-Peer: Structured (DHT)

- Handle peers leaving (churn)
  - Keep  $n$  successors
- New peers
  - Fill gaps
  - Replicate



### Comparison of P2P Systems

For Peer-to-Peer, what are the benefits of (1) central directory vs. (2) unstructured, vs. (3) structured?

#### 1) Central Directory

- Search follows directory (1 lookup)
- Connections  $\rightarrow 1$
- Central point of failure
- Peers control their data
- No neighbours

#### 2) Unstructured

- Search requires flooding ( $n$  lookups)
- Connections  $\rightarrow n^2$
- No central point of failure
- Peers control their data
- Peers control neighbours

#### 3) Structured

- Search follows structure ( $\log(n)$  lookups)
- Connections  $\rightarrow \log(n)$
- No central point of failure
- Peers assigned data
- Peers assigned neighbours

### P2P vs. Client-Server

What are the benefits of Peer-to-Peer vs. Client-Server?

#### Client-Server

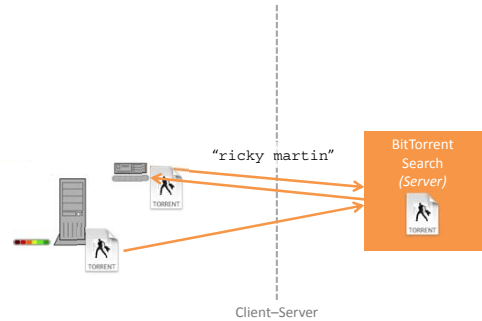
- Data lost in failure/deletes
- Search easier/faster
- Network often faster (to websites on backbones)
- Often central host
  - Data centralised
  - Remote hosts control data
  - Bandwidth centralised
  - Dictatorial
  - Can be taken off-line

#### Peer-to-Peer

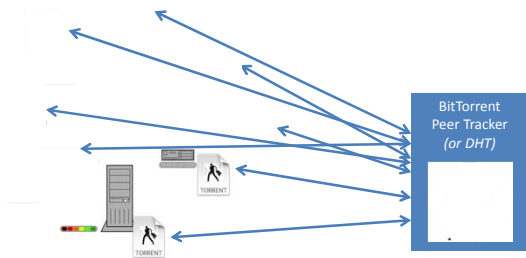
- May lose rare data (churn)
- Search difficult (churn)
- Network often slower (to conventional users)
- Multiple hosts
  - Data decentralised
  - Users (often) control data
  - Bandwidth decentralised
  - Democratic
  - Difficult to take off-line

**DISTRIBUTED SYSTEMS:  
HYBRID EXAMPLE (BITTORRENT)**

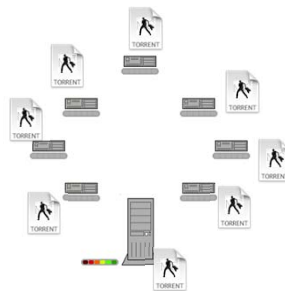
BitTorrent: Search Server



BitTorrent: Tracker



BitTorrent: File-Sharing



BitTorrent: Hybrid

**Uploader**

1. Creates torrent file
2. Uploads torrent file
3. Announces on tracker
4. Monitors for downloaders
5. Connects to downloaders
6. Sends file parts

**Downloader**

1. Searches torrent file
2. Downloads torrent file
3. Announces to tracker
4. Monitors for peers/seeds
5. Connects to peers/seeds
6. Sends & receives file parts
7. Watches illegal movie

Local / Client-Server / Structured P2P / Direct P2P  
(Torrent Search Engines target of law-suits)

**DISTRIBUTED SYSTEMS:  
IN THE REAL WORLD**

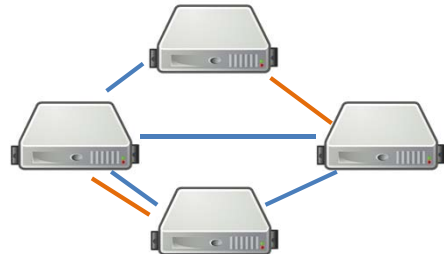
Real-World Architectures: Hybrid

• **Often hybrid!**

- Architectures herein are simplified/idealised
- No clear black-and-white (just good software!)
- For example, *BitTorrent* mixes different paradigms
- But good to know the paradigms

Physical Location: Cluster Computing

- Machines (typically) in a central, local location; e.g., a local LAN in a server room

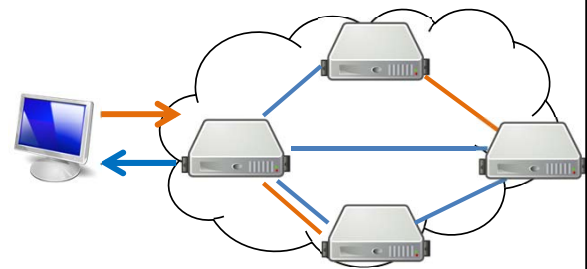


Physical Location: Cluster Computing

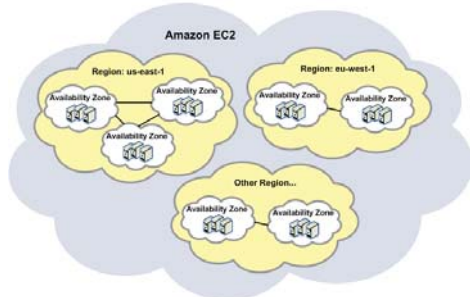


Physical Location: Cloud Computing

- Machines (typically) in a central, remote location; e.g., a server farm like Amazon EC2

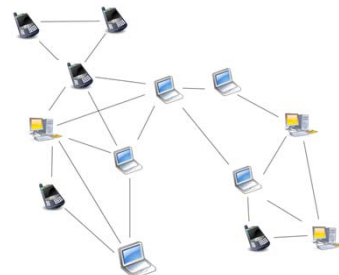


Physical Location: Cloud Computing



Physical Location: Grid Computing

- Machines in diverse locations



Physical Location: Grid Computing



Physical Locations

- Cluster computing:
  - Typically centralised, local
- Cloud computing:
  - Typically centralised, remote
- Grid computing:
  - Typically decentralised, remote

LIMITATIONS OF DISTRIBUTED SYSTEMS: EIGHT FALLACIES

Eight Fallacies

- By L. Peter Deutsch (1994)
  - James Gosling (1997)

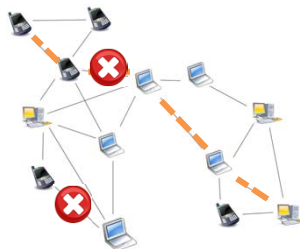
*“Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause big trouble and painful learning experiences.” — L. Peter Deutsch*

- Each fallacy is a **false statement!**

1. The network is reliable

Machines fail, connections fail, firewall eats messages

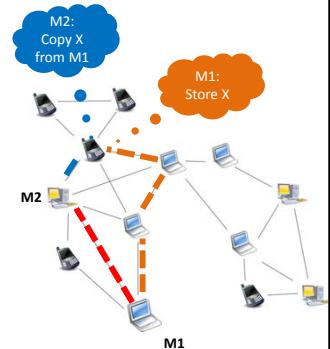
- flexible routing
- retry messages
- acknowledgements!



2. Latency is zero

There are significant communication delays

- avoid “races”
- local order ≠ remote order
- acknowledgements
- minimise remote calls
  - batch data!
- avoid waiting
  - multiple-threads

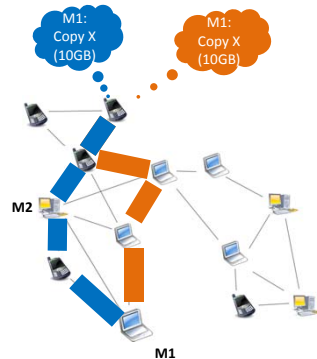




### 3. Bandwidth is infinite

Limited in amount of data that can be transferred

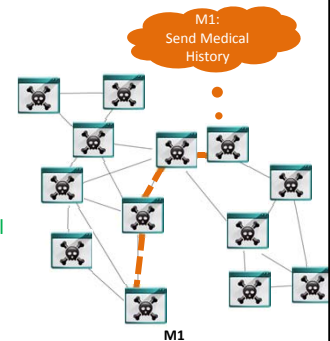
- avoid resending data
- direct connections
- caching!!



### 4. The network is secure

Network is vulnerable to hackers, eavesdropping, viruses, etc.

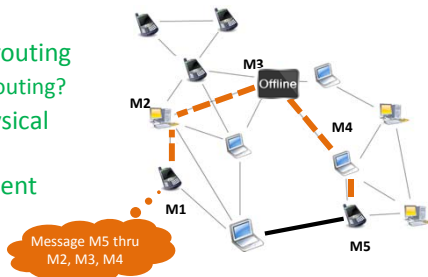
- send sensitive data directly
- isolate hacked nodes
  - hack one node ≠ hack all nodes
- authenticate messages
- secure connections



### 5. Topology doesn't change

How machines are physically connected may change ("churn")!

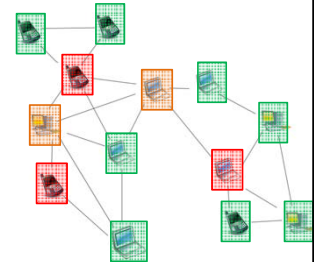
- avoid fixed routing
  - next-hop routing?
- abstract physical addresses
- flexible content structure



### 6. There is one administrator

Different machines have different policies!

- Beware of firewalls!
- Don't assume most recent version
  - Backwards compat.

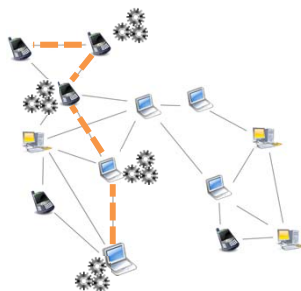


### 7. Transport cost is zero

It costs time/money to transport data: not just bandwidth

(Again)

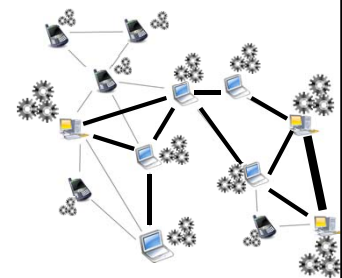
- minimise redundant data transfer
  - avoid shuffling data
  - caching
- direct connection
- compression?



### 8. The network is homogeneous

Devices and connections are not uniform

- interoperability!
  - Java vs. .NET?
- route for speed
  - not hops
- load-balancing



### Eight Fallacies (to avoid)

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous

**Severity of fallacies vary in different scenarios!**

Which fallacies apply/do not apply for:

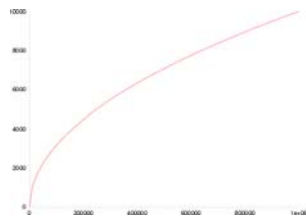
- Gigabit ethernet LAN?
- BitTorrent
- The Web

### LABS REVIEW/PREP

### Why did it work in memory?

We processed a lot of data. Why did it work in memory?

- Not so many unique words ...
  - but lots of new proper nouns
  - Heap's law:
    - $U(n) \approx Kn^\beta$
    - English text
      - $K \approx 10$
      - $\beta \approx 0.6$



### What if it doesn't work in memory?

What if it doesn't work in memory?



How could we implement a word-count (or a bi-gram count) using the hard disk for storage?

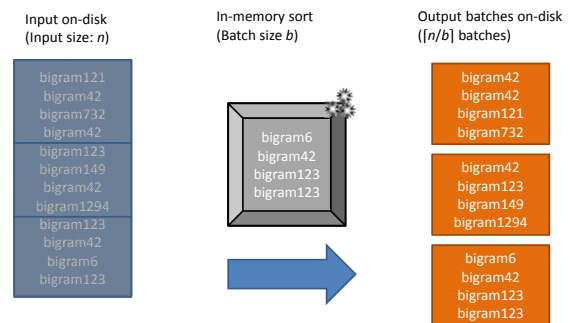
### Most generic method: use sorting

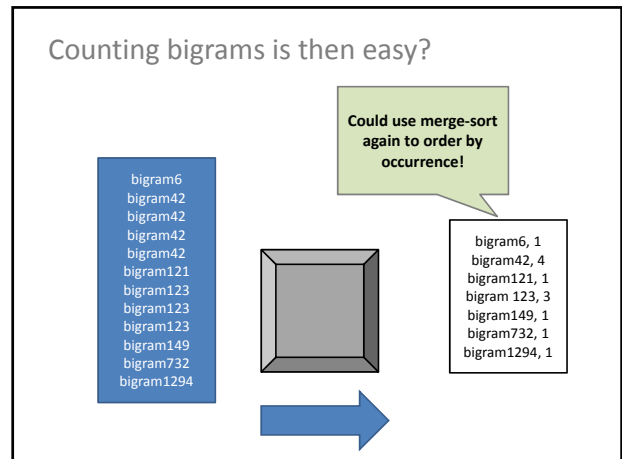
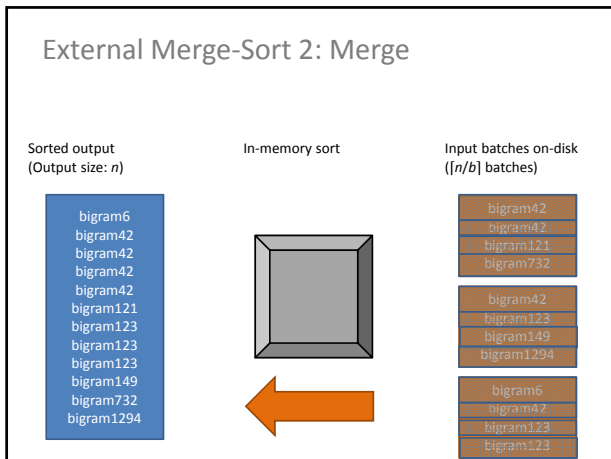


How can we use the disk to sort?

### External Merge-Sort 1: Batch

- Sort in batches





### Does external merge-sorting scale?

Any problem with external merge-sorting as we scale really high?

- If you have too many batches to read simultaneously, disk will go nuts

Any solution(s)?

- Use lots of main-memory to reduce batch count
- Only merge  $k$  at a time

If we have  $n$  batches and merge them  $k$  at a time, how many passes will we need?

### RECAP

- ### Topics Covered
- What is a (good) Distributed System?
  - Client-Server model
    - Fat/thin client
    - Mirror/proxy servers
    - Three-tier
  - Peer-to-Peer (P2P) model
    - Central directory
    - Unstructured
    - Structured (Hierarchical/DHT)
    - BitTorrent

- ### Topics Covered
- Physical locations:
    - Cluster (local, centralised) vs.
    - Cloud (remote, centralised) vs.
    - Grid (remote, decentralised)
  - 8 fallacies
    - Network isn't reliable
    - Latency is not zero
    - Bandwidth not infinite,
    - etc.

### Topics Covered

- External Merge Sorting
  - Split data into batches
  - Sort batches in memory
  - Write batches to disk
  - Merge sorted batches into final output

### Questions?

