# Week 1:   Part 1

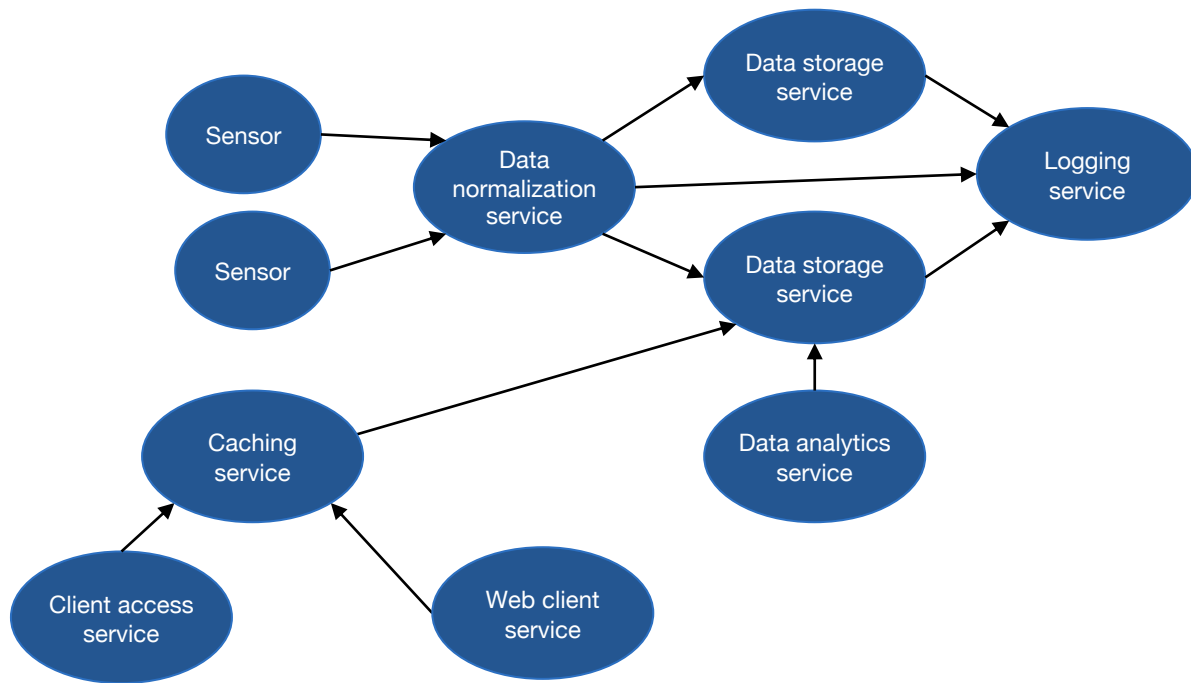## Introduction to distributed systems

**Paul Krzyzanowski**

# What is a Distributed System?

*A collection of independent computers connected through a communication network that work together to accomplish some goal*

– No shared operating system

– No shared memory

– No shared clock

# What is a Distributed System?

A distributed system is a collection of services accessed via network interfaces

*Collection of independent computers that appears as a single system to the user(s)*

*Independent* = autonomous, self-contained
*Single system* = user not aware of distribution

# Classifying parallel and distributed systems

# Flynn's Taxonomy (1966)

Classify computer architectures by looking at the number of instruction streams and number of data streams

1. **SISD** — Single Instruction, Single Data stream
   - Traditional uniprocessor systems

2. **SIMD** — Single Instruction, Multiple Data streams
   - Array (vector) processors
   - Examples:
     - GPUs – Graphical Processing Units for computer graphics, GPGPU (General Purpose GPU): AMD/ATI, NVIDIA
     - AVX: Intel's Advanced Vector Extensions

3. **MISD** — Multiple Instructions, Single Data stream
   - Sometimes (rarely!) applied to classifying fault-tolerant redundant systems

4. **MIMD** — Multiple Instruction, Multiple Data streams
   - Multiple computers, each with a program counter, program (instructions), data
   - **Parallel and distributed systems**

# Subclassifying MIMD

**Memory**

- Shared memory systems: **multiprocessors**
- No shared memory: networks of computers, **multicomputers**

**Interconnect**

- Bus
- Switch

**Delay/bandwidth**

- Tightly coupled systems
- Loosely coupled systems

# Multiprocessors & Multicomputers

**Multiprocessors**

- Shared memory
- Shared clock
- Shared operating system
- All-or-nothing failure

**Multicomputers** (networks of computers) ⇒ *distributed systems*

- No shared memory
- No shared clock
- Partial failures
- Inter-computer communication mechanism needed: the network
  - Traffic volume much lower than memory access

# Why do we want distributed systems?

1. Scale

2. Collaboration

3. Reduced latency

4. Mobility

5. High availability & Fault tolerance

6. Incremental cost

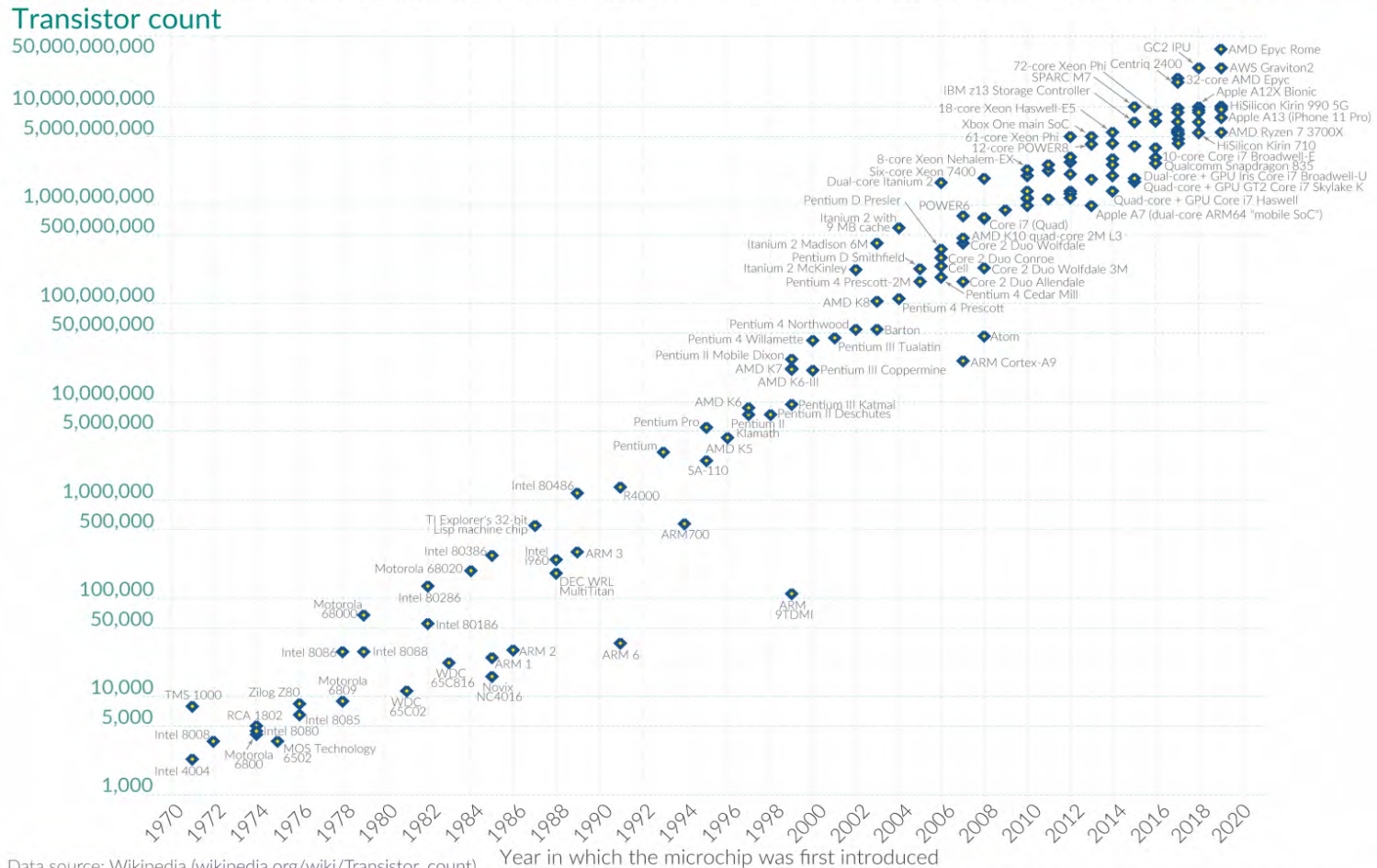7. Delegated infrastructure & operations

# 1. Scale

# Scale: Increased Performance

Computers are getting faster

## Moore's Law

Prediction: performance doubles approximately every 18 months because of faster transistors and more transistors per chip

*Not a real law* – just an observation from the mid 1970s

Transistor count — Year in which the microchip was first introduced

Data source: Wikipedia (wikipedia.org/wiki/Transistor_count)

# Scaling a single system has limits

Getting harder for technology to keep up with Moore's law

- More cores per chip
    - → requires multithreaded programming

- There are limits to the die size and # of transistors
    - Intel Xeon W-3175X CPU: 28 cores per chip ($2,999/chip!)
        - 8 billion transistors, 255 watts @ 3.1-4.3 GHz
    - AMD EPYC 7601 CPU: 32 cores per chip ($4,200/chip)
        - 19.2 billion transistors, 180 watts
    - NVIDIA GeForce RTX 2080 Ti: 4,352 CUDA cores per chip
        - Special purpose apps: Graphics rendering, neural networks

# More performance

What if we need more performance than a single CPU?

Combine them ⇒ multiprocessors

Distributed systems allow us to achieve massive performance

# Our computing needs exceed CPU advances

**Movie rendering**

– *Toy Story* (1995) – 117 computers; 45 mins — 30 hours to render a frame

– *Toy Story 4* (2019) – 60-160 hours to render a frame

**Google**

– Over 63,000 search queries per second on average

– Over 130 trillion pages indexed

– Uses hundreds of thousands of servers to do this

**Facebook**

– Approximately 100M requests per second with 4B users

# Example: Google

- In 1999, it took Google one month to crawl and build an index of about 50 million pages

- In 2012, the same task was accomplished in less than one minute.

- 16% to 20% of queries that get asked every day have never been asked before

- Every query has to travel on average 1,500 miles to a data center and back to return the answer to the user

- A single Google query uses 1,000 computers in 0.2 seconds to retrieve an answer

Source: http://www.internetlivestats.com/google-search-statistics/

# 2. Collaboration

# Collaboration & Content

- Collaborative work & play
- Social connectivity
- Commerce
- News & media

# Metcalfe's Law

*The value of a telecommunications network is proportional to the square of the number of connected users of the system.*

The Network Effect ⇒ This makes networking interesting to us … and to investors!

# 3. Reduced latency

# Reduced Latency

- Cache data close to where it is needed

- *Caching* vs. *replication*
  - Replication: multiple copies of data for increased fault tolerance
  - Caching: temporary copies of frequently accessed data closer to where it's needed


- Some caching services:
  Akamai, Cloudflare, Amazon Cloudfront, Apache Ignite
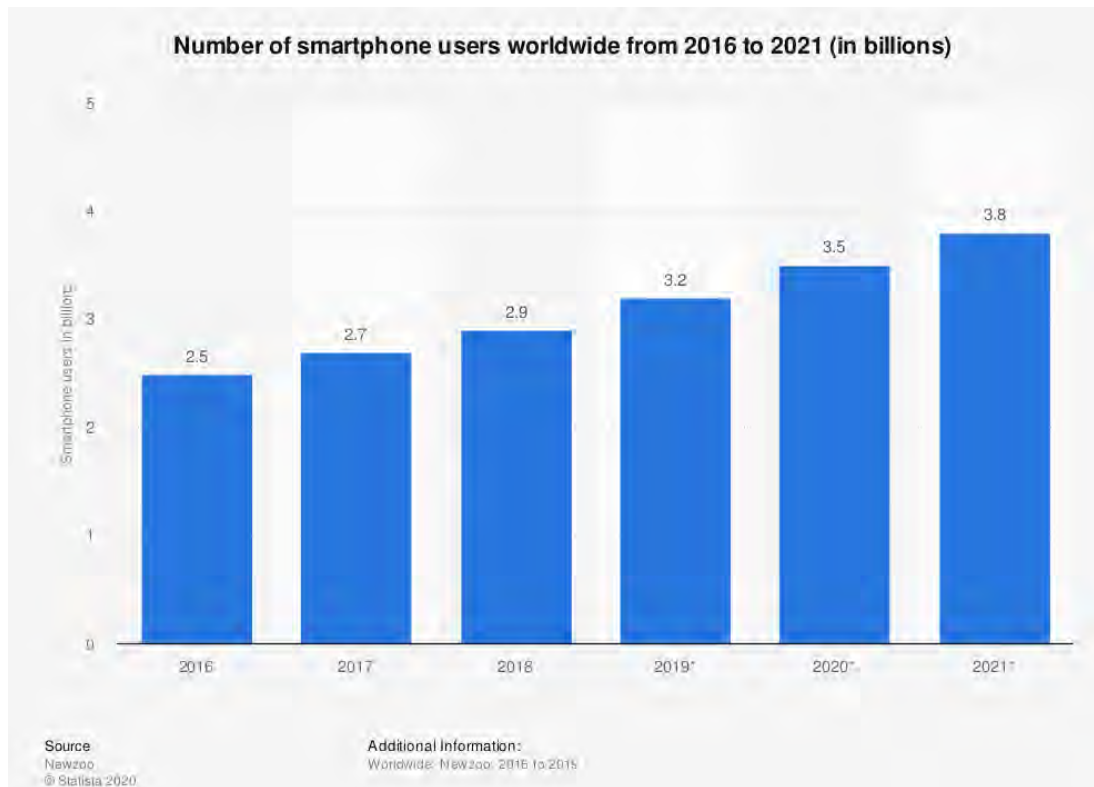
# 4. Mobility

3.5 billion smartphone users

Remote sensors
- Cars
- Traffic cameras
- Toll collection
- Shipping containers
- Vending machines

IoT = Internet of Things
- 2017: more IoT devices than humans



Number of smartphone users worldwide from 2016 to 2021 (in billions)

| Year | Smartphone users in billions |
| --- | --- |
| 2016 | 2.5 |
| 2017 | 2.7 |
| 2018 | 2.9 |
| 2019 | 3.2 |
| 2020 | 3.5 |
| 2021 | 3.8 |

Source
Newzoo
© Statista 2020

Additional Information:
Worldwide; Newzoo; 2016 to 2015

# 5. High availability & Fault tolerance

# High availability

Redundancy = replicated components

Service can run even if some systems die

> Reminder:
>
> $$P(\text{A and B}) = P(\text{A}) \times P(\text{B})$$

If $P(\text{any one system down}) = 5\%$

$P(\text{two systems down at the same time}) = 5\% \times 5\% = 0.25\%$

**Uptime = 1 – downtime = 1 – 0.0025 = 99.75%**

We get 99.7% uptime instead of 95% because we need _both_ replicated components to fail instead of just one.

No redundancy = dependence on <u>all</u> components

If we need **_all_** systems running to provide a service

$P$(any system down) = 1 - $P$( A is up <u>AND</u> B is up )

$\quad$ = 1 - (1-5%) × (1-5%) = 1 - 0.95 × 0.95 = 9.75%

$\quad$ ⇒ 39x greater than a single component failure with redundancy!

**Uptime = 1 – downtime = 1 – 0.0975 = 90.25%**

*With a large # of systems, P*(any system down) approaches 100% !

Requiring a lot of components to be up & running is a losing proposition.
With large enough systems, something is always breaking!

# Availability: series & parallel systems

**Series system**:  The system fails if ANY of its components fail

*P(system failure) = 1 - P(system survival)*

If $P_i$ = P(component *i* fails) then for *n* components:

$$P(system\ failure) = 1 - \prod_{i}^{n}(1 - P_i)$$

**Parallel system**:  The system fails only if ALL of its components fail

*P(system failure) = P(component$_1$ fails) × P(component$_2$ fails)  …*
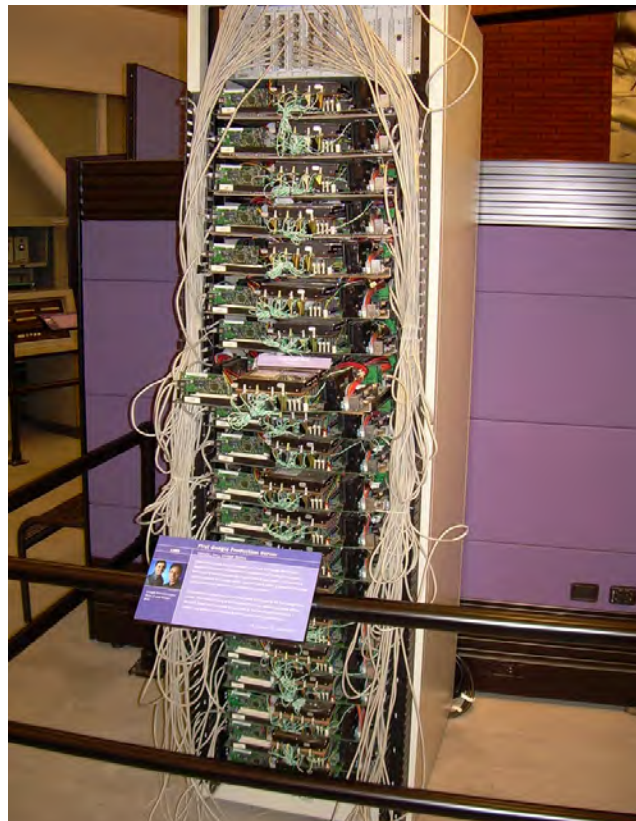
$$P(system\ failure) = \prod_{i}^{n} P_i$$

# Availability requires fault tolerance

- Fault tolerance
  - Identify & recover from component failures

- Recoverability
  - Software can restart and function
  - May involve restoring state

# 6. Incremental growth & cost

- Version 1 does not have to be the full system

  – Add more servers & storage over time

  – Scale also implies cost – you don't need millions of $ for v1.0

# 7. Delegated infrastructure & operations

# Delegated operations

- **Offload responsibility**
  - Let someone else manage systems
  - Use third-party services

- **Speed deployment**
  - Don't buy & configure your own systems
  - Don't build your own data center

- **Modularize services on different systems**
  - Dedicated systems for storage, email, etc.

- **Use cloud, network attached storage**
  - Let someone else figure out how to expand storage and do backups

# Transparency as a Design Goal

# Transparency

High level: hide distribution from users

Low level: hide distribution from software

- **Location transparency**
  Users don't care where resources are

- **Migration transparency**
  Resources move at will

- **Replication transparency**
  Users cannot tell whether there are copies of resources

- **Concurrency transparency**
  Users share resources transparently

- **Parallelism transparency**
  Operations take place in parallel without user's knowledge

# Core challenges in distributed systems design

1. Concurrency

2. Latency

3. Partial Failure

# Concurrency

# Concurrency

- Lots of requests may occur at the same time

- Need to deal with concurrent requests
  - Need to ensure consistency of all data
  - Understand critical sections & mutual exclusion
  - Beware: mutual exclusion (locking) can affect performance

- Replication adds complexity
  - All operations must appear to occur in the same order on all replicas

# Latency

# Latency

Network messages may take a long time to arrive

- **Synchronous network model**
  - There is some upper bound, $T$, between when a node sends a message and another node receives it
  - Knowing $T$ enables a node to distinguish between a node that has failed and a node that is taking a long time to respond
- **Partially synchronous network model**
  - There's an upper bound for message communication but the programmer doesn't know it – it has to be discovered
  - Protocols will operate correctly only if all messages are received within some time, $T$
    - We cannot make assumptions on the delay time distribution
- **Asynchronous network model**
  - Messages can take arbitrarily long to reach a peer node
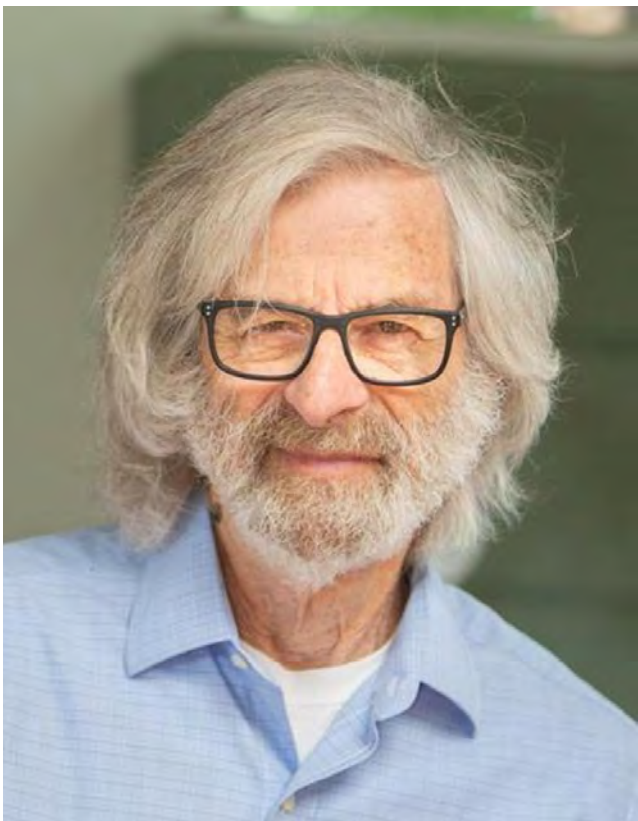  - **This is what we get from the Internet!**

# Latency & asynchronous networks

- Asynchronous networks can be a pain

- Messages may take an unpredictable amount of time

  - We may think a message is lost but it's really delayed

  - May lead to retransmissions → duplicate messages

  - May lead us to assume a service is dead when it isn't

  - May mess with our perception of time

  - May cause messages to arrive in a different order
        … or a different order on different systems

# Latency

- Speed up data access via **caching** – temporary copies of data

- Keep data close to where it's processed to maximize efficiency
  - Memory vs. disk
  - Local disk vs. remote server
  - Remote memory vs. remote disk
  - **Cache coherence**: cached data can become **stale**
    - Underlying data can change → cache needs to be invalidated
    - System using the cache may change the data → propagate results
      - *Write-through cache*
      - But updates take time ⇒ can lead to **inconsistencies** (**incoherent views**)

# Partial Failure

You know you have a distributed system when the crash of a computer you've never heard of stops you from getting any work done.

*— Leslie Lamport*

**Failure is a fact of life in distributed systems!**

In local systems, failure is usually **total** (*all-or-nothing*)

In distributed systems, we get **partial failure**

- A component can fail while others continue to work

- Failure of a network link is indistinguishable from a remote server failure

- Send a request but don't get a response ⇒ what happened?

No **global state**

- There is no global state that can be examined to determine errors
- There is no agent that can determine which components failed and inform everyone else

Need to ensure the state of the entire system is consistent after a failure

# Handling failure

Handle **detection**, **recovery**, and **restart**

**Availability** = fraction of time system is usable

– Achieve with redundancy

– But then consistency is an issue!

**Reliability**: data must not get lost

– Includes security

# System Failure Types

- **Fail-stop**

    - Failed component stops functioning

    - **Halting** = stop without notice

    - Detect failed components via **timeouts**

        - But you can't count on timeouts in asynchronous networks

            - And what if the network isn't reliable?

        - Sometimes we guess

- **Fail-restart**

    - Component stops but then restarts

    - Danger: **stale state**
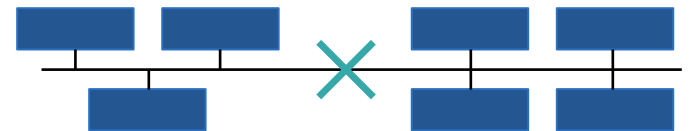
# Network Failure Types

- **Omission**

  – Failure to send or receive messages

    - Due to queue overflow in router, corrupted data, receive buffer overflow

- **Timing**

  – Messages take longer than expected

    - We may assume a system is dead when it isn't

  – Unsynchronized clocks can alter process coordination

- **Partition**

  – Network fragments into two or more
    sub-networks that cannot communicate with each other

# Network & System Failure Types

- **Fail-silent**
  - A failed component (process or hardware) does not produce any output

- **Byzantine failures**
  - Instead of stopping, a component produces faulty data
  - Due to bad hardware, software, network problems, or malicious interference

Goal: avoid single points of failure

# Redundancy

We deal with failures by adding redundancy
- – Replicated components


But this means we need to keep the **state** of those components replicated

# State, replicas, and caches

- **State**

  - Information about some component that cannot be reconstructed
  - Network connection info, process memory, list of clients with open files, lists of which clients finished their tasks

- **Replicas**

  - Redundant copies of data → *used to address fault tolerance*

- **Cache**

  - Local storage of frequently-accessed data to reduce latency
    → *used to address latency*

# No global knowledge

- Nobody has the true **global state** of a system
  - There is no global state that can be examined to determine errors
  - There is no agent that can determine which components failed and inform everyone else
  - No shared memory

- A process knows its current state
  - It may know the *last reported state* of other processes
  - It may periodically report its state to others

*No foolproof way to detect failure in all cases*

# Other design considerations

# Handling Scale

- Need to be able to add and remove components

- Impacts failure handling
  - If failed components are removed, the system should still work
  - If replacements are brought in, the system should integrate them

# Security

- The environment
  - Public networks, remotely-managed services, 3$^{rd}$ party services

- Some issues
  - Malicious interference, bad user input, impersonation of users & services
  - Protocol attacks, input validation attacks, time-based attacks, replay attacks

Rely on authentication, cryptography (hashes, encryption)
                                    … and good defensive programming!

- Users also want convenience
  - Single sign-on, no repeated entering of login credentials
  - Controlled access to services

# Other design considerations

- Algorithms & environment
  - Distributable vs. centralized algorithms
  - Programming languages
  - APIs and frameworks

# Main themes in distributed systems

- **Availability & fault tolerance**
  - Fraction of time that the system is functioning
  - Dead systems, dead processes, dead communication links, lost messages

- **Scalability**
  - Things are easy on a small scale
  - But on a large scale
    - Geographic latency (multiple data centers), administering many thousands of systems

- **Latency & asynchronous processes**
  - Processes run asynchronously: concurrency
  - Some messages may take longer to arrive than others

- **Security**
  - Authentication, authorization, encryption

# Key approaches in distributed systems

- Divide & conquer
  - Break up data sets (sharding) and have each system work on a small part
  - Merging results is usually the easy & efficient part

- Replication
  - For high availability, caching, and sharing data
  - Challenge: keep replicas consistent even if systems go down and come up

- Quorum/consensus
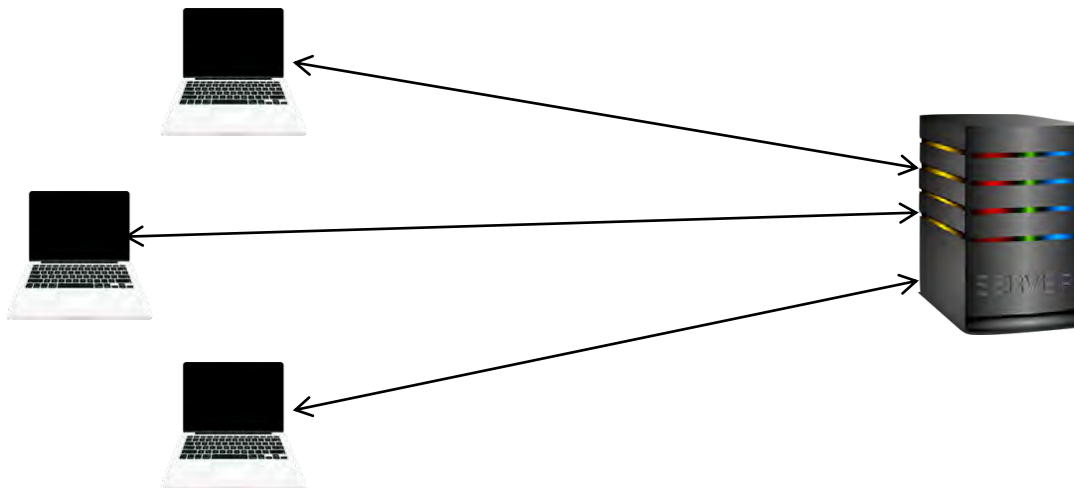  - Enable a group to reach agreement

# Service Models (Application Architectures)

# Centralized model

- No networking

- Traditional time-sharing system

- Single workstation/PC or direct connection of multiple terminals to a computer

- One or several CPUs

- Not easily scalable

- Limiting factor: number of CPUs in system
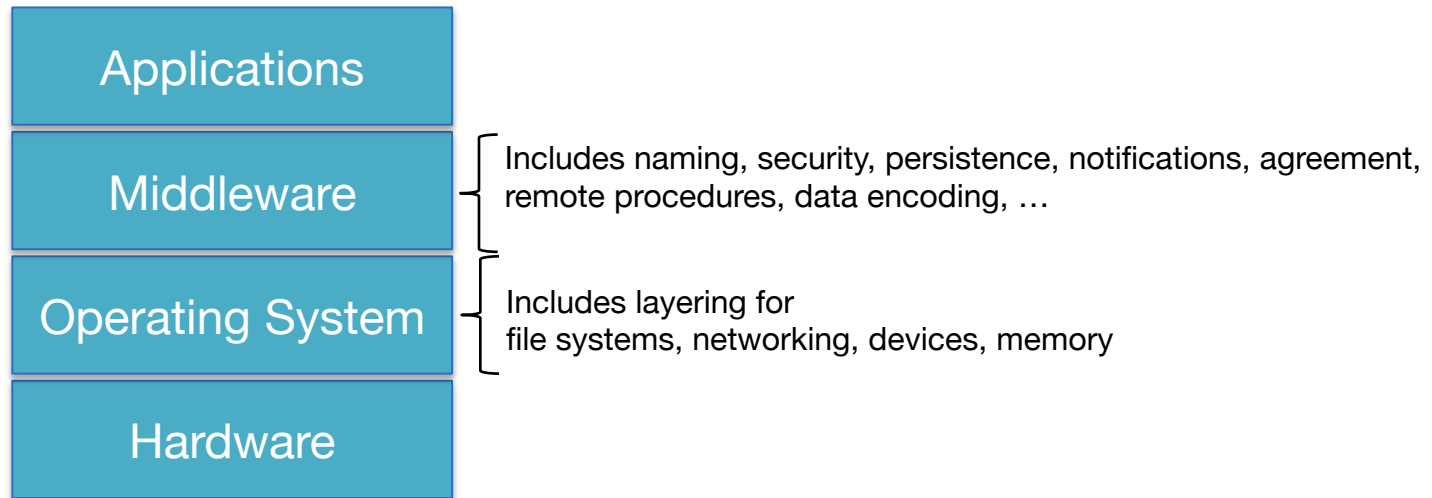  - Contention for same resources (memory, network, devices)

# Client-Server model

- Clients send requests to servers

- A server is a system that runs a service

- Clients do not communicate with other clients

# Layered architectures in software design

- Break functionality into multiple layers

- Each layer handles a specific abstraction
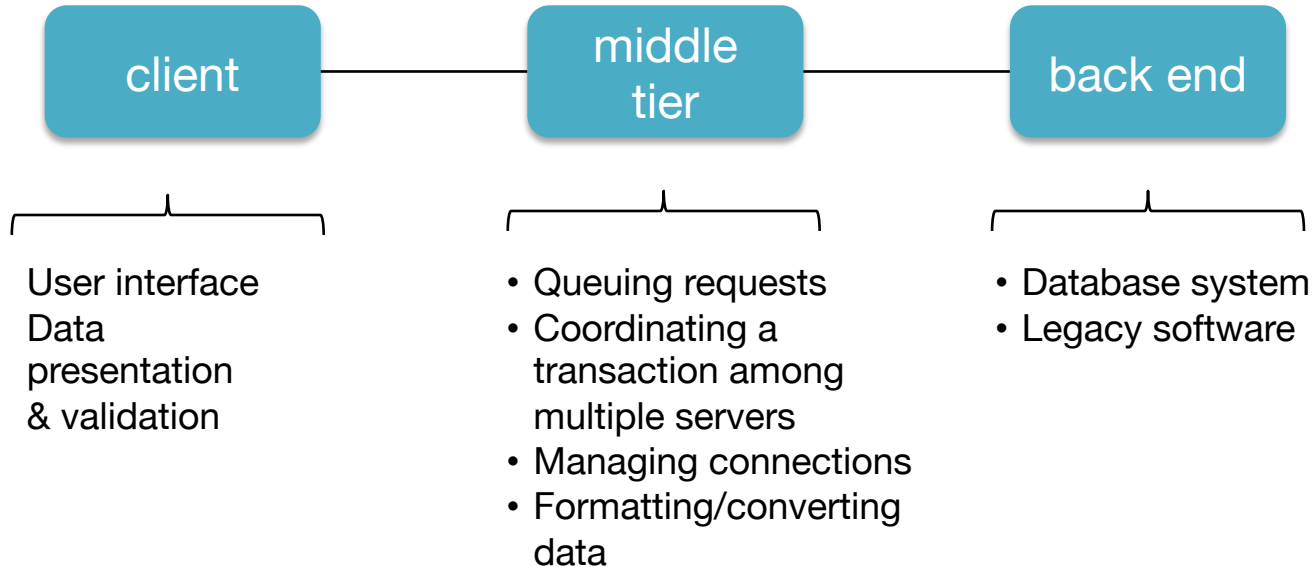  - Hides implementation details and specifics of hardware, OS, network abstractions, data encoding, …

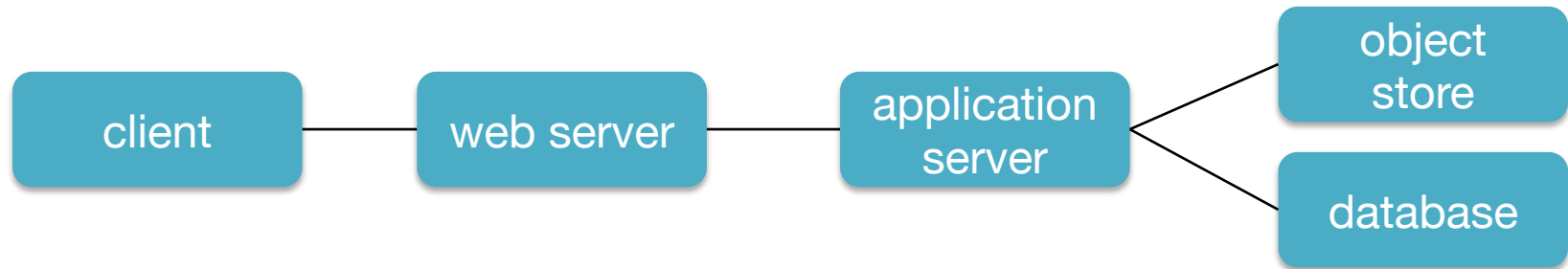| | |
|---|---|
| **Applications** | |
| **Middleware** | Includes naming, security, persistence, notifications, agreement, remote procedures, data encoding, … |
| **Operating System** | Includes layering for file systems, networking, devices, memory |
| **Hardware** | |

# Tiered architectures in networked systems

- Tiered (multi-tier) architectures
  - Distributed systems analogy to a layered architecture

- Each tier (layer)
  - Runs as a network service
  - Is accessed by surrounding layers

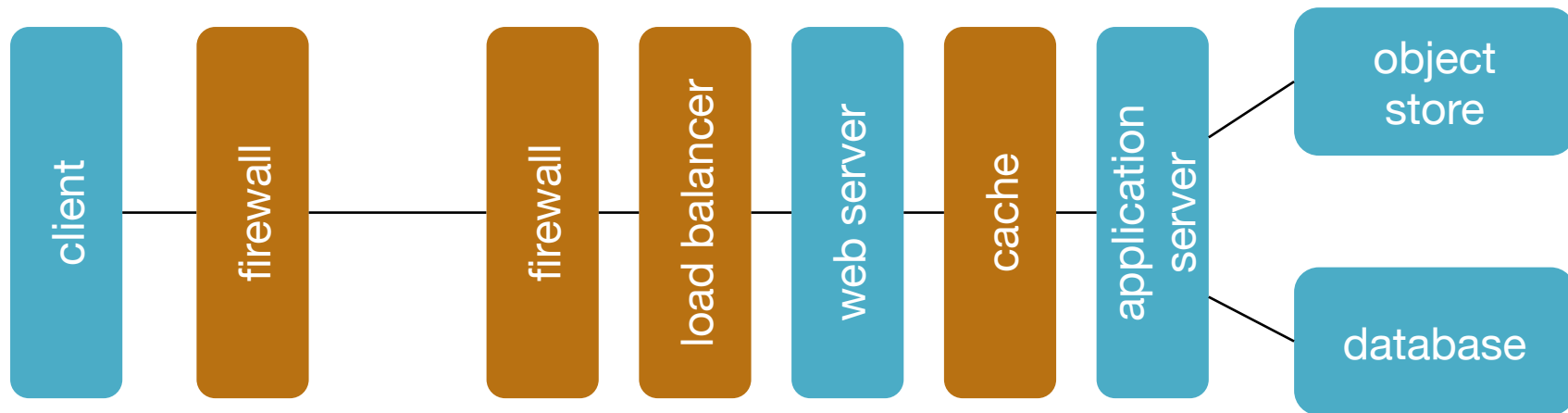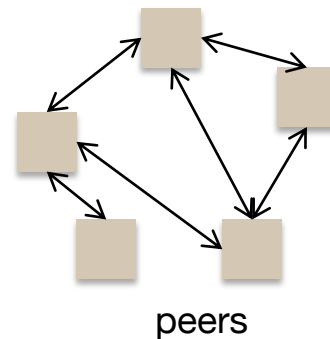The basic client-server architecture is a two-tier model

# Multi-tier example



| client | middle tier | back end |
|--------|-------------|----------|
| User interface<br>Data presentation & validation | • Queuing requests<br>• Coordinating a transaction among multiple servers<br>• Managing connections<br>• Formatting/converting data | • Database system<br>• Legacy software |

# Multi-tier example



client — web server — application server — object store / database

# Multi-tier example

Some tiers may be transparent to the application



client — firewall — firewall — load balancer — web server — cache — application server — object store / database
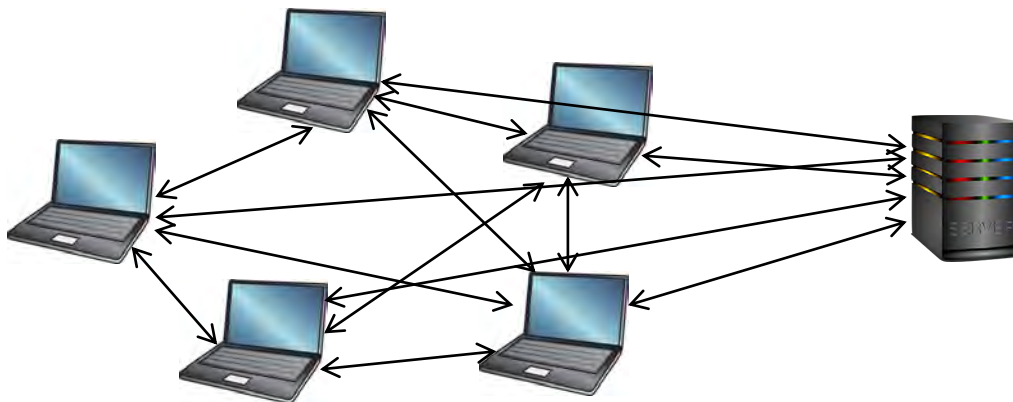
# Peer-to-Peer (P2P) Model

- No reliance on servers

- Machines (peers) communicate with each other

- Goals
  - Robustness
  - Self-scalability

- Examples
  - BitTorrent, Skype

server

clients

peers

# Hybrid model

- Many peer-to-peer architectures still rely on a server
  - Look up, track users
  - Track content
  - Coordinate access

- But traffic-intensive workloads are delegated to peers

Images from: http://clipart-library.com/laptop-cliparts.html

# Processor pool model

- Collection of CPUs that can be assigned processes on demand

- Similar to hybrid model
  - Coordinator dispatches work requests to available processors

- Render farms, big data processing, machine learning

# Cloud Computing

Resources are provided as a network (Internet) service

**Software as a Service** (SaaS)
Remotely hosted software: email, productivity, games, …
   *Salesforce.com, Google Apps, Microsoft 365*

**Platform as a Service** (PaaS)
Execution runtimes, databases, web servers, development environments, …
   *Google App Engine, AWS Elastic Beanstalk*

**Infrastructure as a Service** (IaaS)
Compute + storage + networking: VMs, storage servers, load balancers
   *Microsoft Azure, Google Compute Engine, Amazon Web Services*

**Storage**
Remote file storage
   • *Dropbox, Box, Google Drive, OneDrive, …*

# The End