



# Data Engineering: Introduction

# Computers



# Mark I

**Release date:** August 7, 1944; 77 years ago

**Power:** 5 horsepower (3.7 kW)

**Dimensions:** 816 cubic feet ( $23 \text{ m}^3$ ) – 51 feet (16 m) in length, 8 feet (2.4 m) in height, and 2 feet (0.61 m) deep

**Mass:** 9,445 pounds (4.7 short tons; 4.3 t)

- Store 72 numbers, each 23 decimal digits long
- Do 3 additions or subtractions in a second
- Multiplication took 6 seconds
- Division took 15.3 seconds



# Fagaku

**Memory:** HBM2 32 GiB/node

**Storage:** 1.6 TB NVMe SSD/16 nodes (L1)

- 150 PB shared Lustre FS (L2)
- Cloud storage services (L3)

**Speed:** 442 PFLOPS

**Cost:** US\$1 billion (total programme cost)



Data Engineering  
Introduction

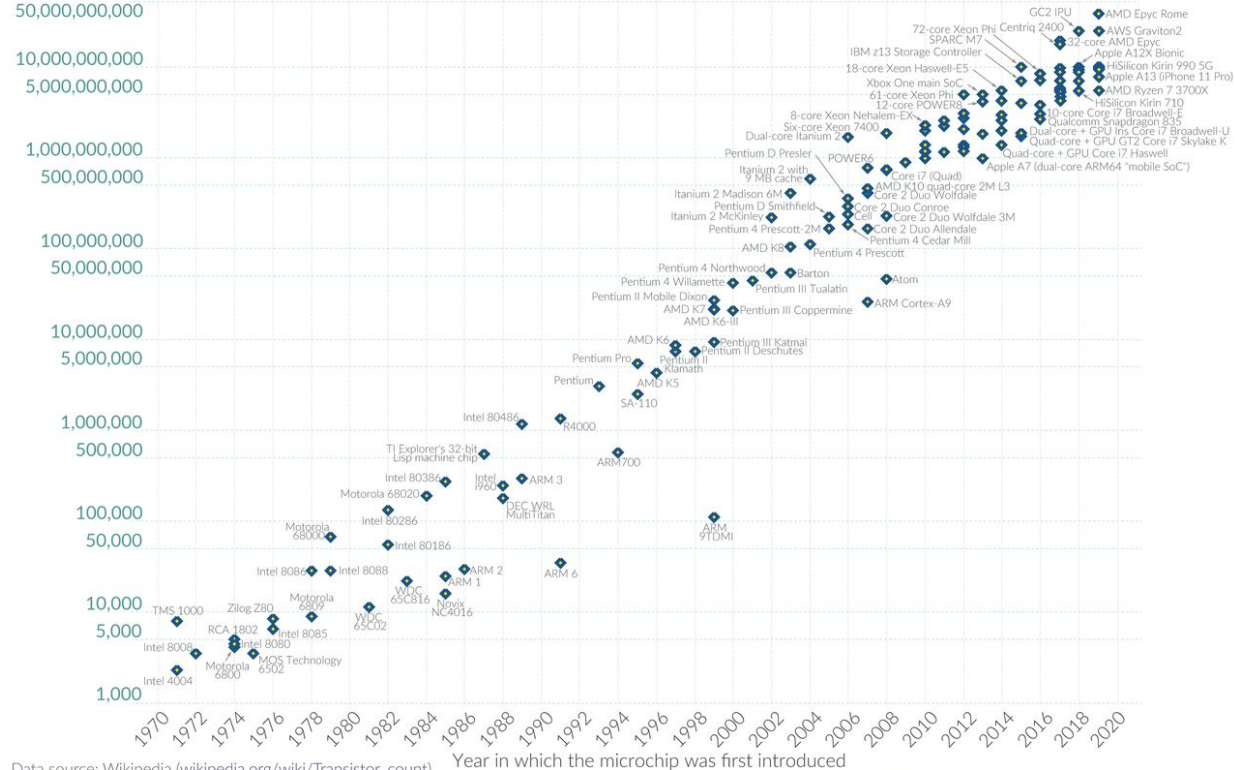


# Moore's Law: The number of transistors on microchips doubles every two years

Our World  
in Data

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

## Transistor count



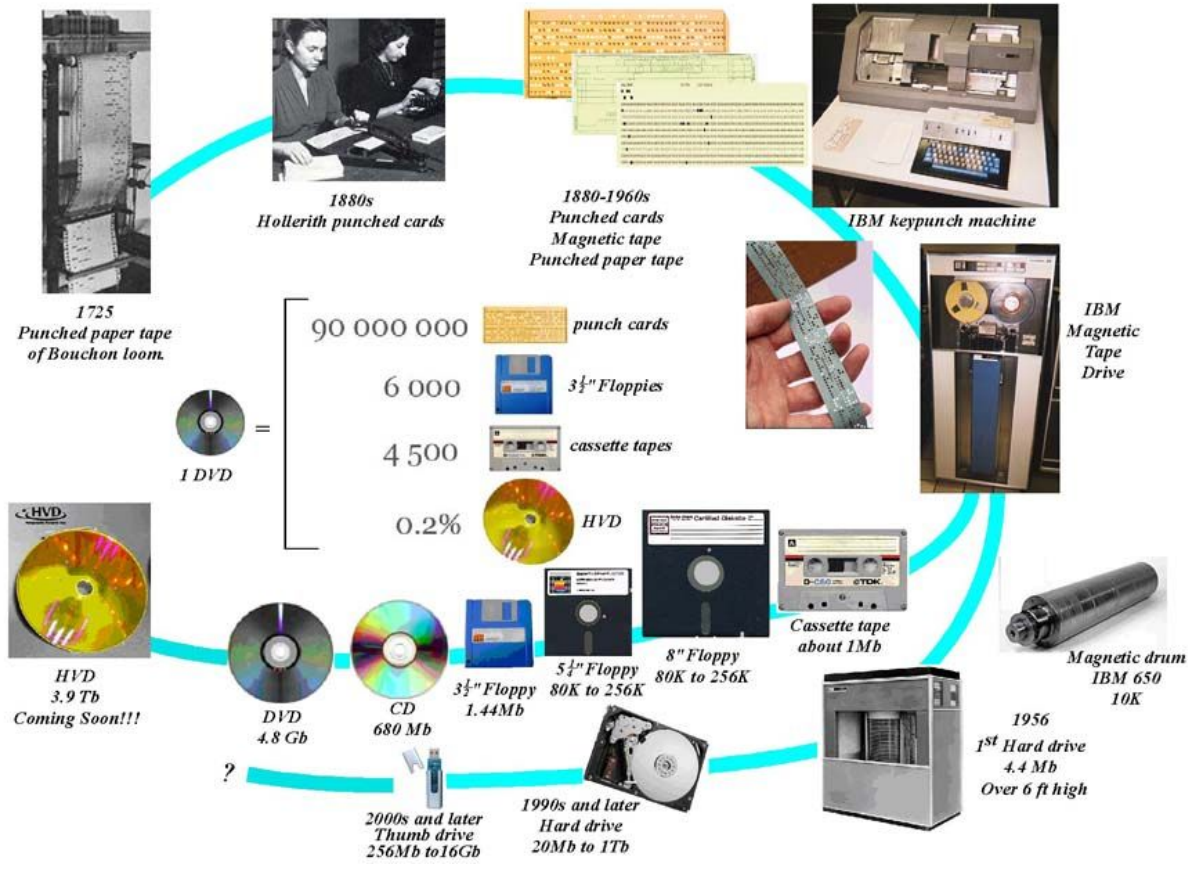
Data source: Wikipedia ([wikipedia.org/wiki/Transistor\\_count](https://wikipedia.org/wiki/Transistor_count))

OurWorldInData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

Data Engineering  
Introduction

Wikipedia

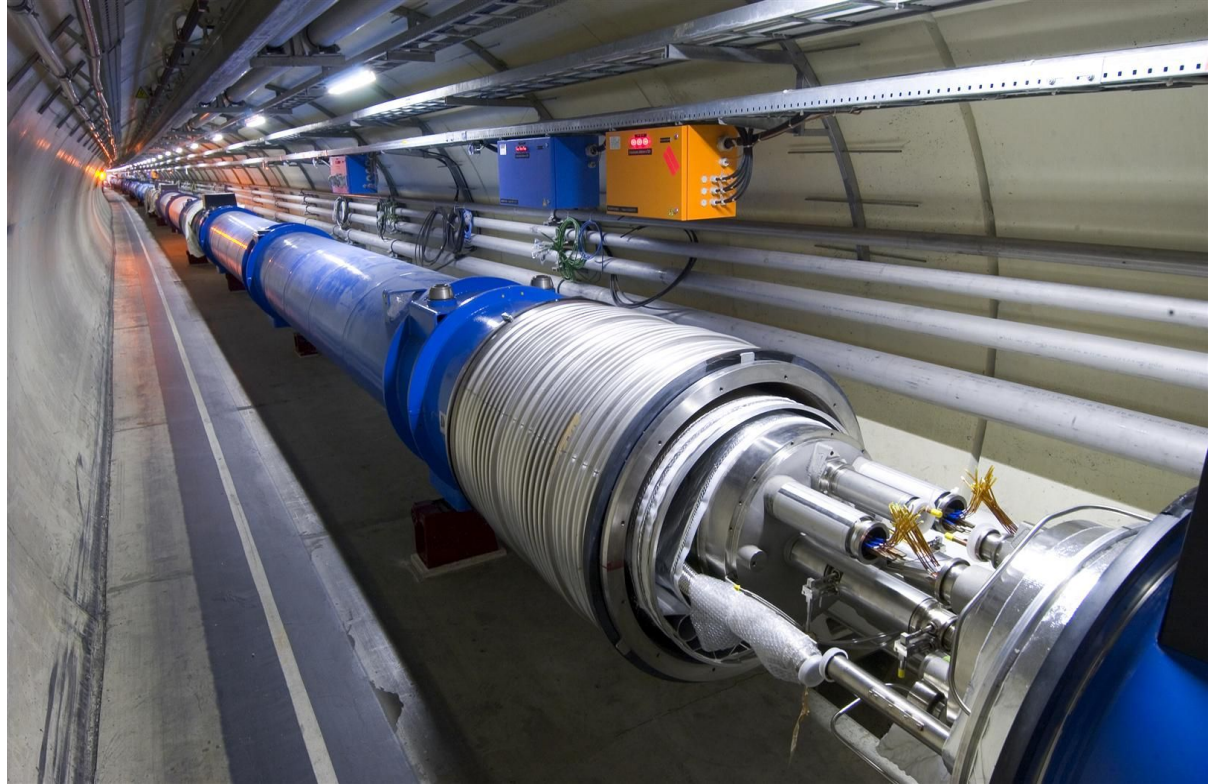




# CERN

- 50 PB/Year
- Storage: 1 ZB
- 15000 servers

<https://monit-grafana-open.cern.ch/d/00000884/it-overview?orgId=16>



# Why make a system distributed?

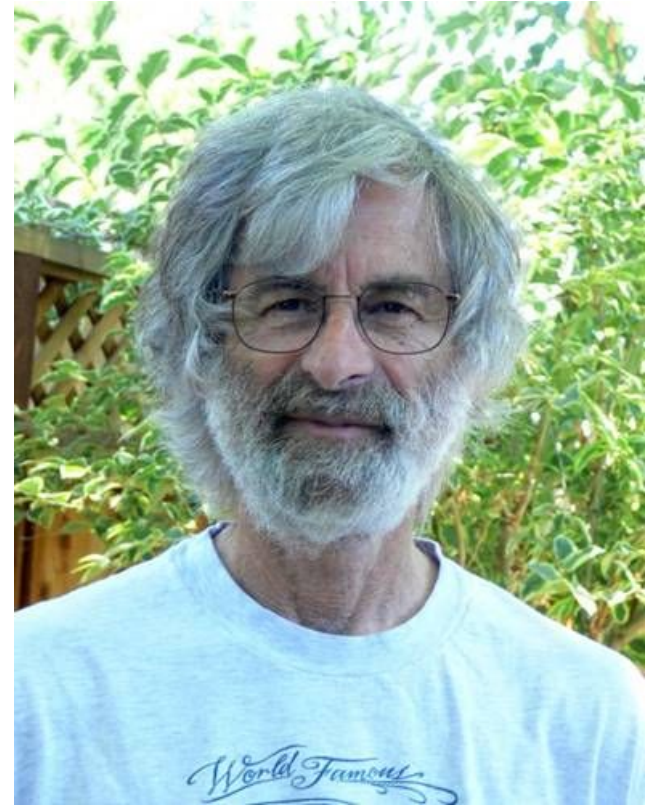
- It's inherently distributed
- For better reliability
- For better performance
- To solve bigger problems





# Distributed System

“... a system in which the failure of a computer you didn't even know existed can render your own computer unusable.” — Leslie Lamport



# Scalability



# Instagram

## 1 Billion Monthly Active Users

1B  
June 2018

800M  
Sep 2017

700M  
April 2017

600M  
Dec 2016

500M  
June 2016

400M  
Sep 2015

300M  
Dec 2014

200M  
March 2014

150M  
Sep 2013

100M  
Feb 2013

50M  
May 2012

30M iPhone users /  
Android Launch  
April 3, 2012

1M Android signups  
in 24 hours  
April 4, 2012

5M Android users /  
Acquired by Facebook  
April 9, 2012

Launch on  
iPhone  
Oct 2010

1M  
Dec 2010

5M  
June 2011

10M  
Sep 2011

2010

2011

2012

2013

2014

2015

2016

2017

2018



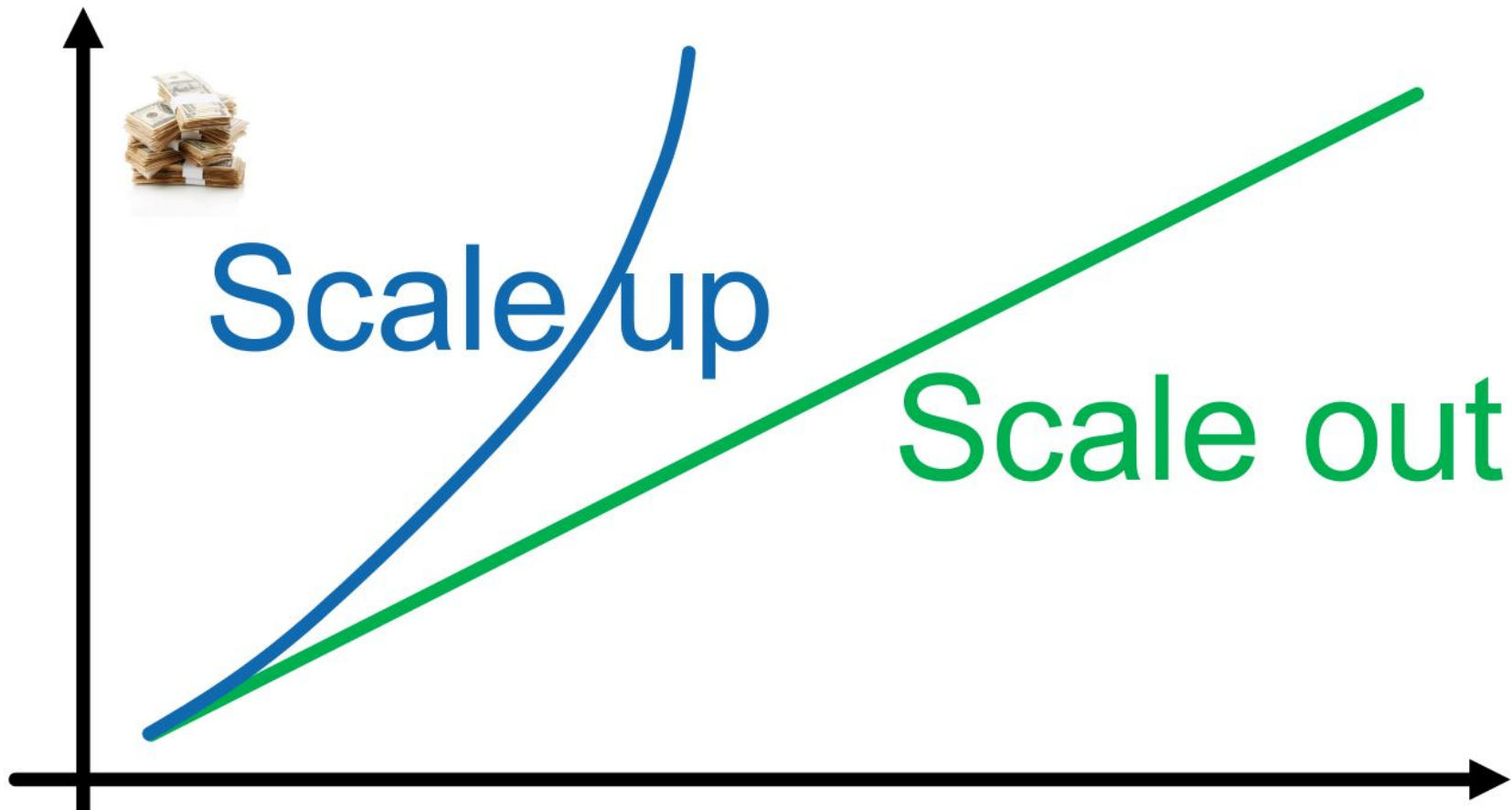


Scale Up



Scale Out

# Hardware price comparison



# What's the problems?

Application Clients (End Users)



Internet



Software Load Balancer



Hardware Load Balancer

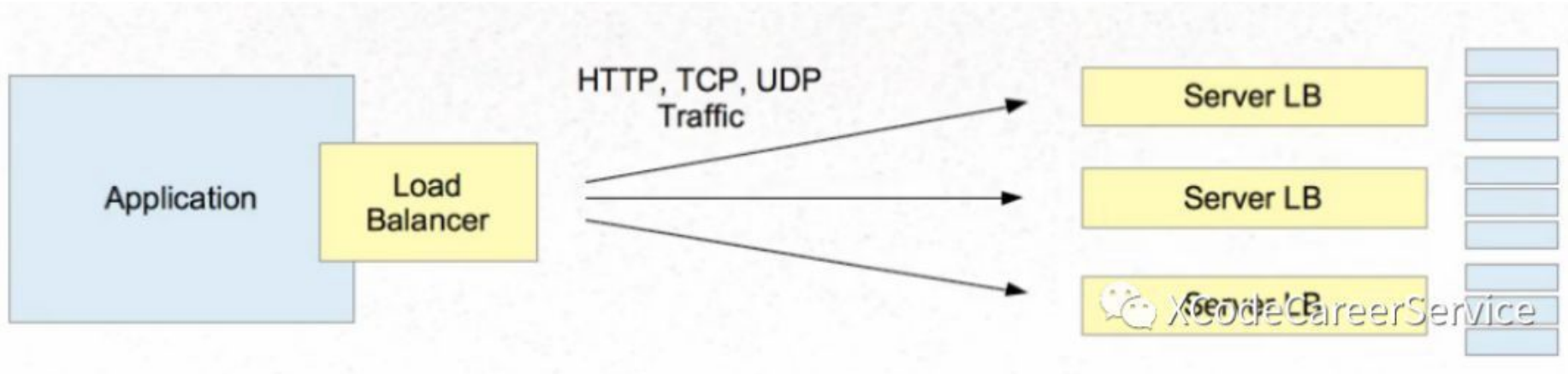


Application Servers

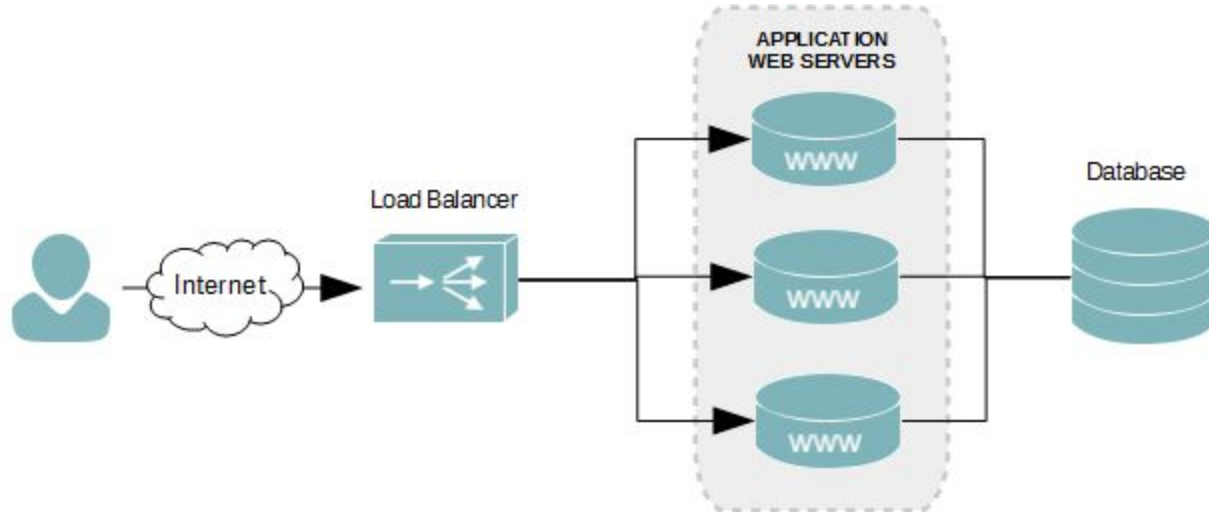




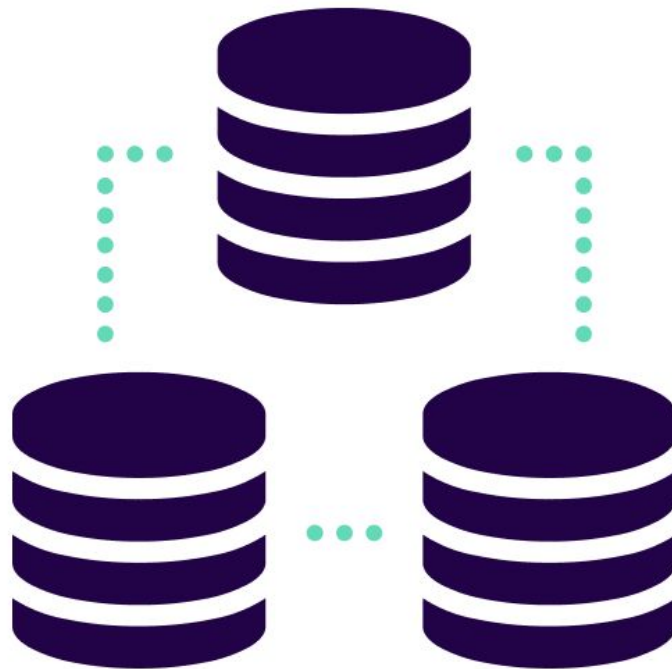
# Smart Client-Side Load Balancing

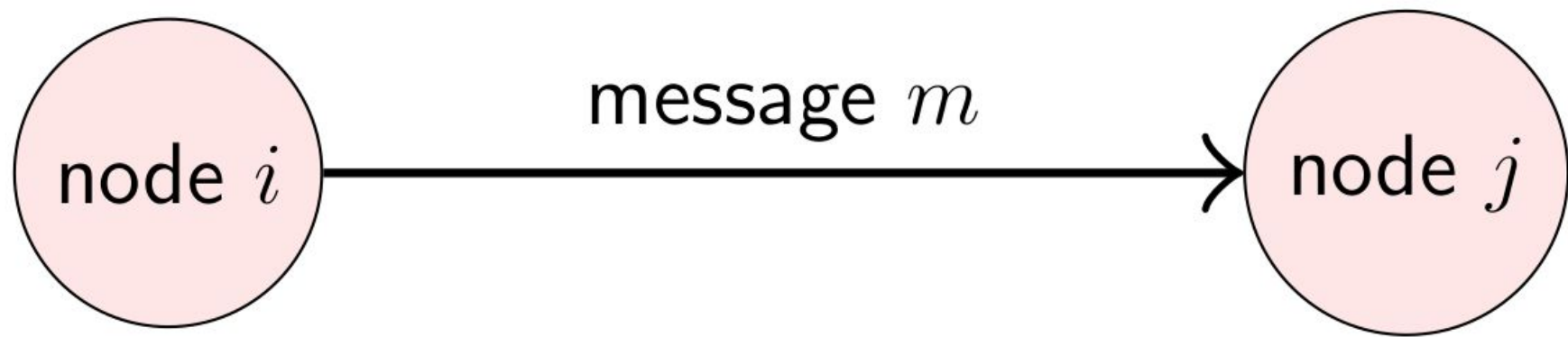


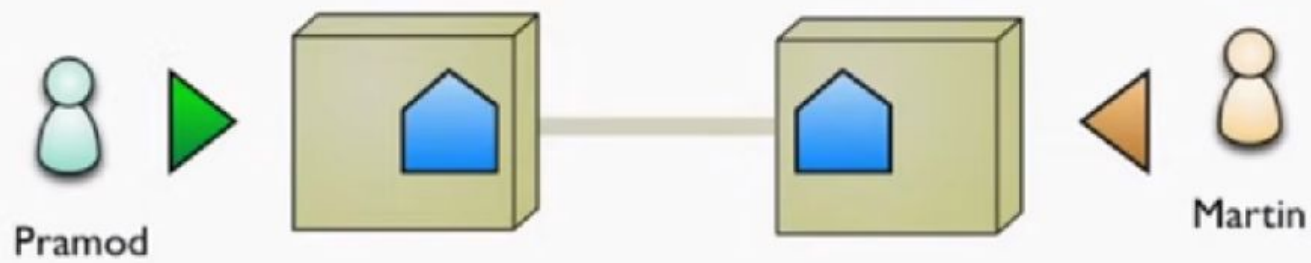
# Shared Database



- Sharding
- Replication

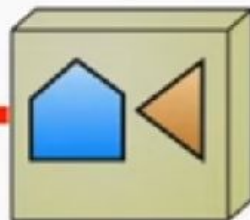
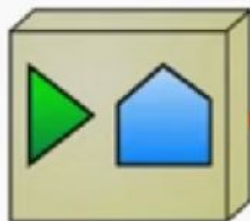








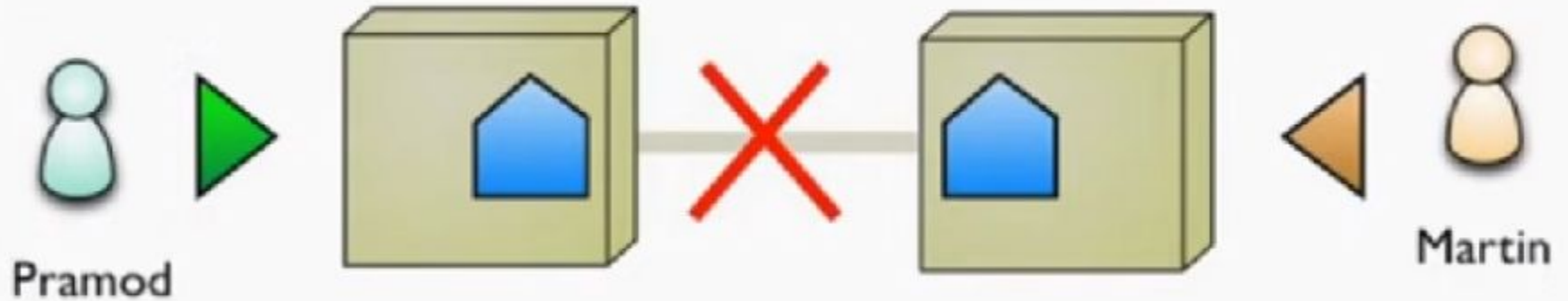
Pramod

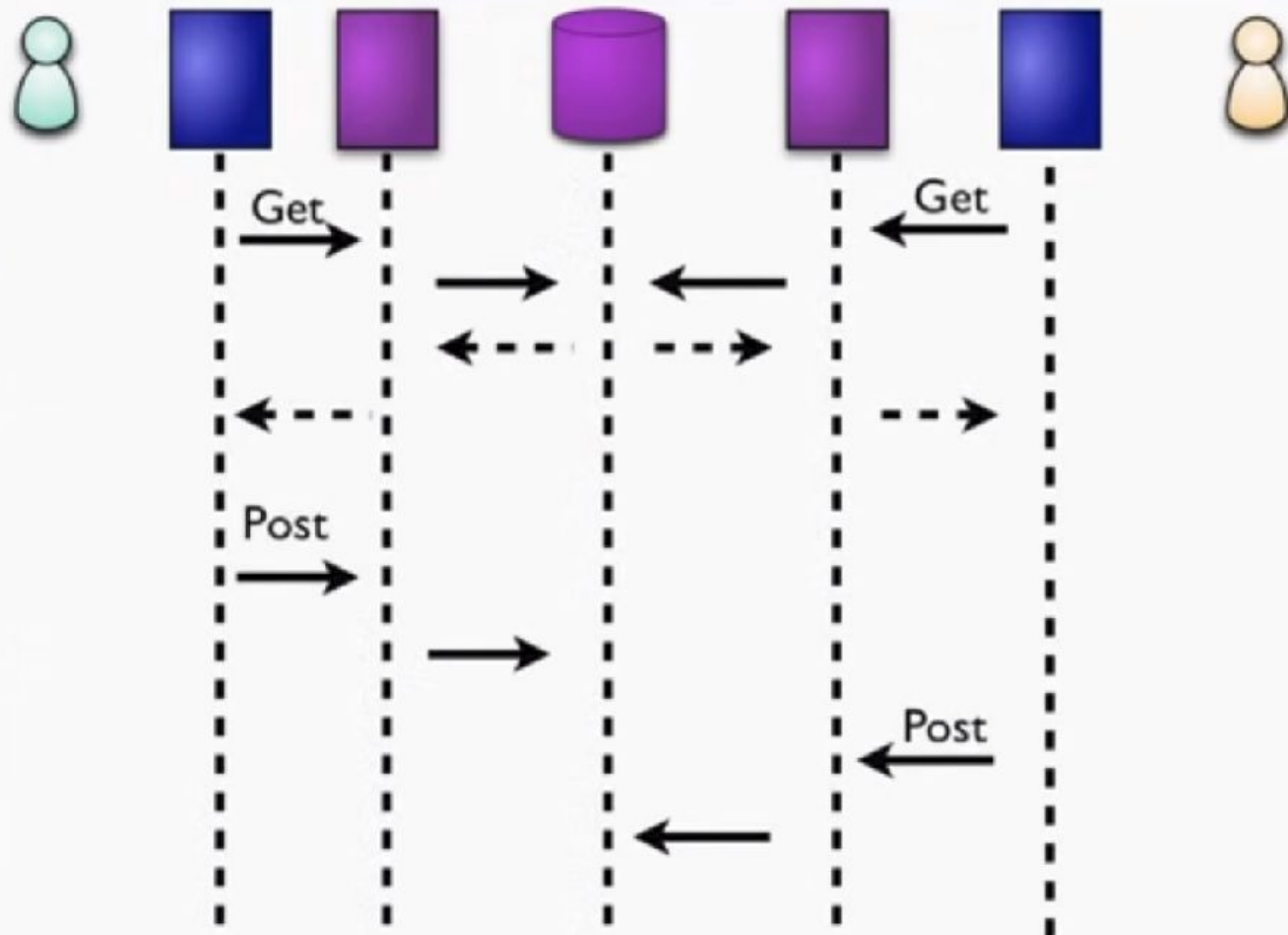


Martin

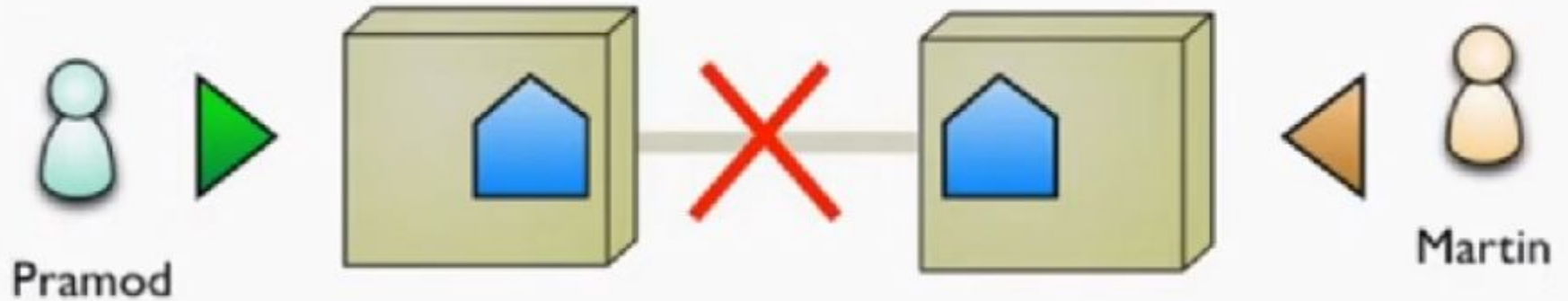


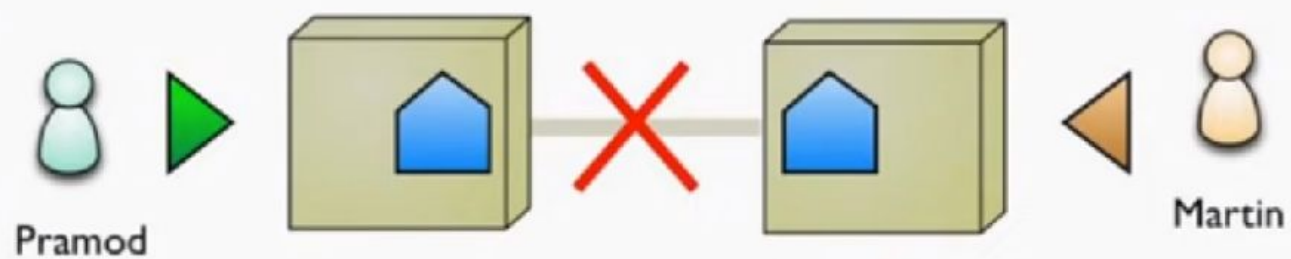
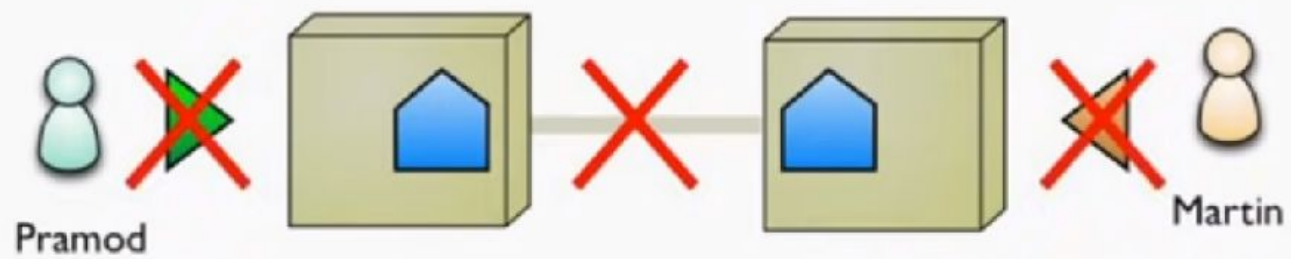
# Partition Tolerance

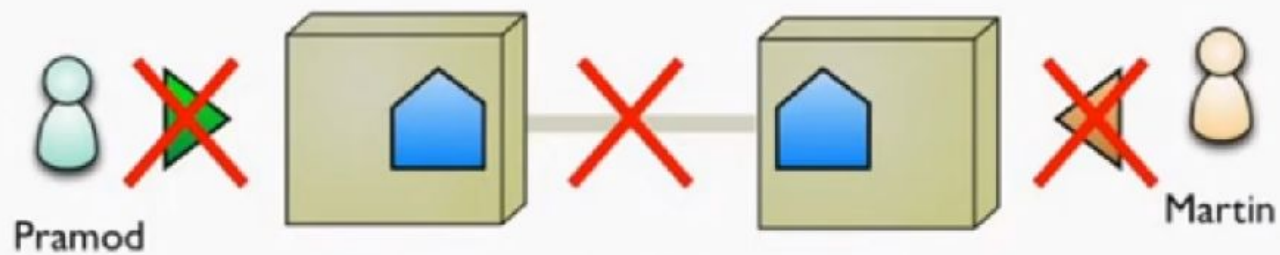




# Partition Tolerance



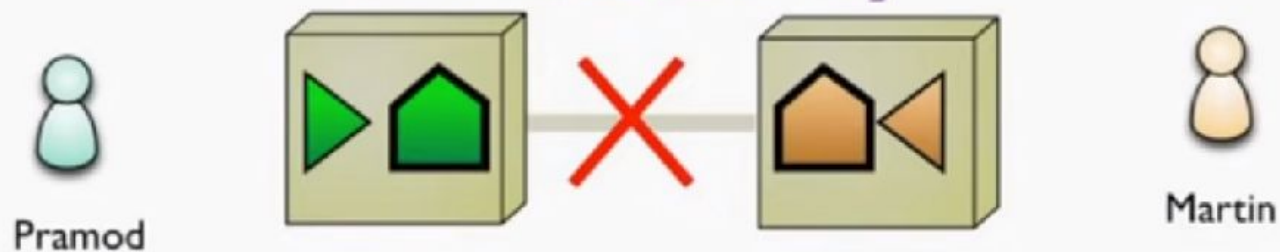




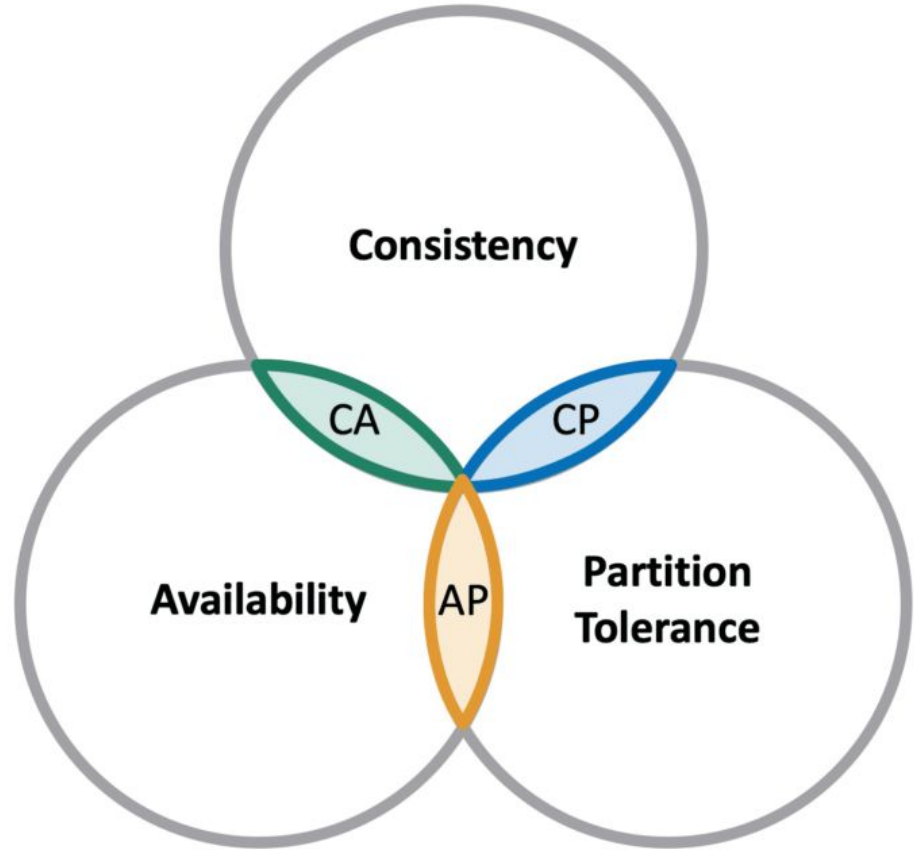
## Consistency



## Availability



# Cap Theorem





# Assumptions

- The network is reliable
- The network is secure
- The network is homogeneous
- The topology does not change
- Latency is zero
- Bandwidth is infinite
- Transport cost is zero
- There is one administrator

# An Example: Distributed Cache System



## HDD

CrystalDiskMark 3.0.1 x64

File Edit Theme Help Language

5 1000MB C: 13% (88/699GB)

	Read [MB/s]	Write [MB/s]
All		
Seq	112.3	109.3
512K	41.69	48.05
4K	0.543	0.693
4K Qb32	1.004	0.698

## SSD

CrystalDiskMark 3.0.1 x64

File Edit Theme Help Language

5 1000MB C: 45% (101/224GB)

	Read [MB/s]	Write [MB/s]
All		
Seq	477.9	235.7
512K	402.7	248.9
4K	30.49	64.67
4K Qb32	200.3	233.2

## RAM

CrystalDiskMark 3.0.1 x64

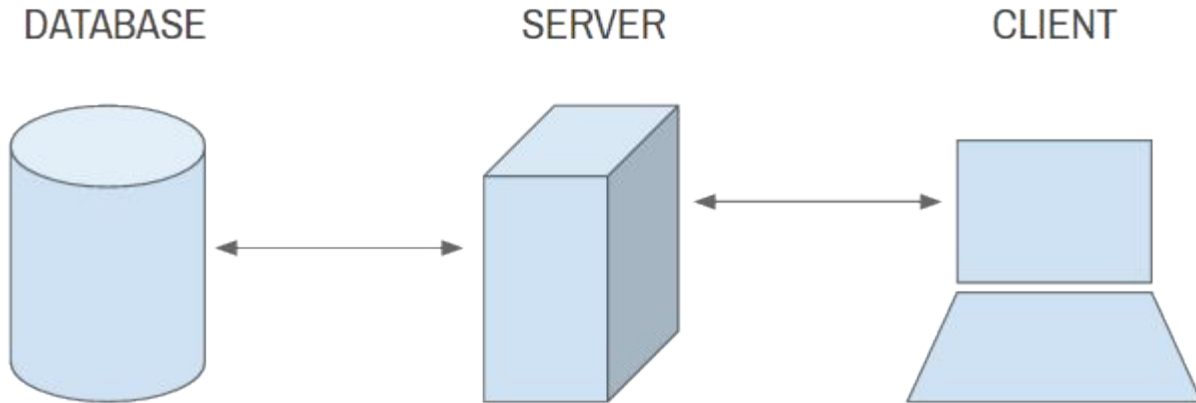
File Edit Theme Help Language

5 1000MB R: 1% (48/4089MB)

	Read [MB/s]	Write [MB/s]
All		
Seq	5766	7760
512K	5649	7172
4K	657.0	554.8
4K Qb32	631.9	544.7

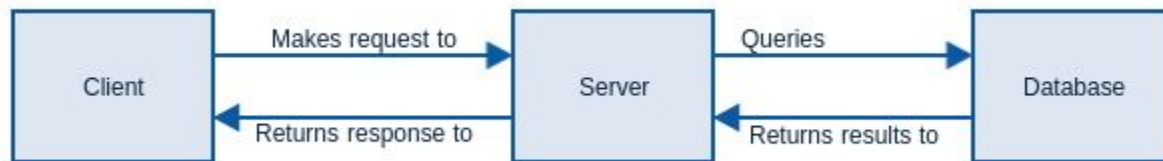


# Latency

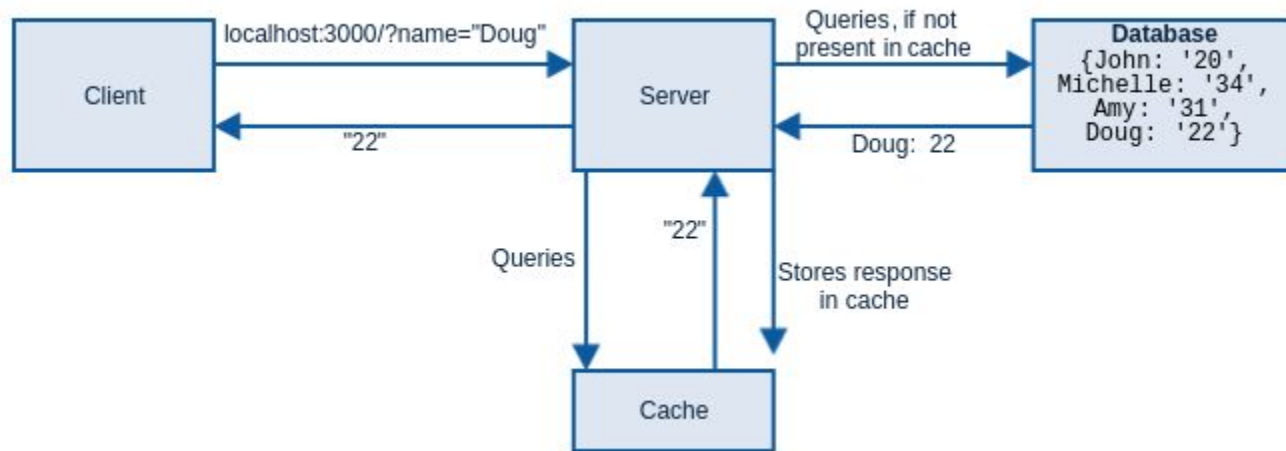


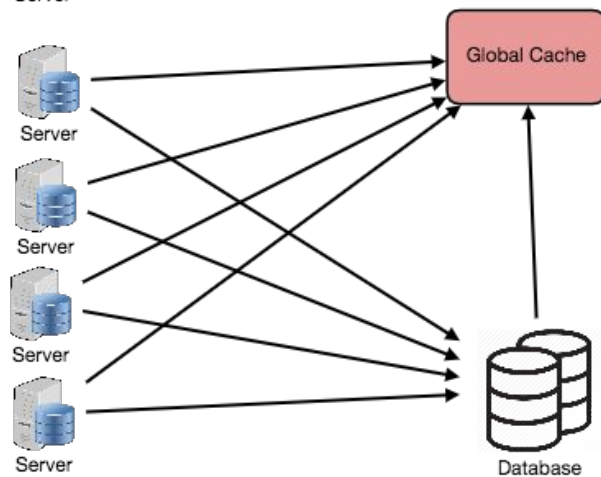
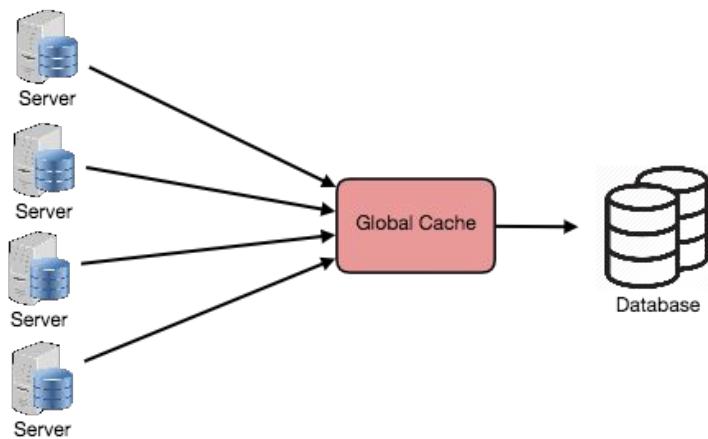
# A Cache is like short-term memory.







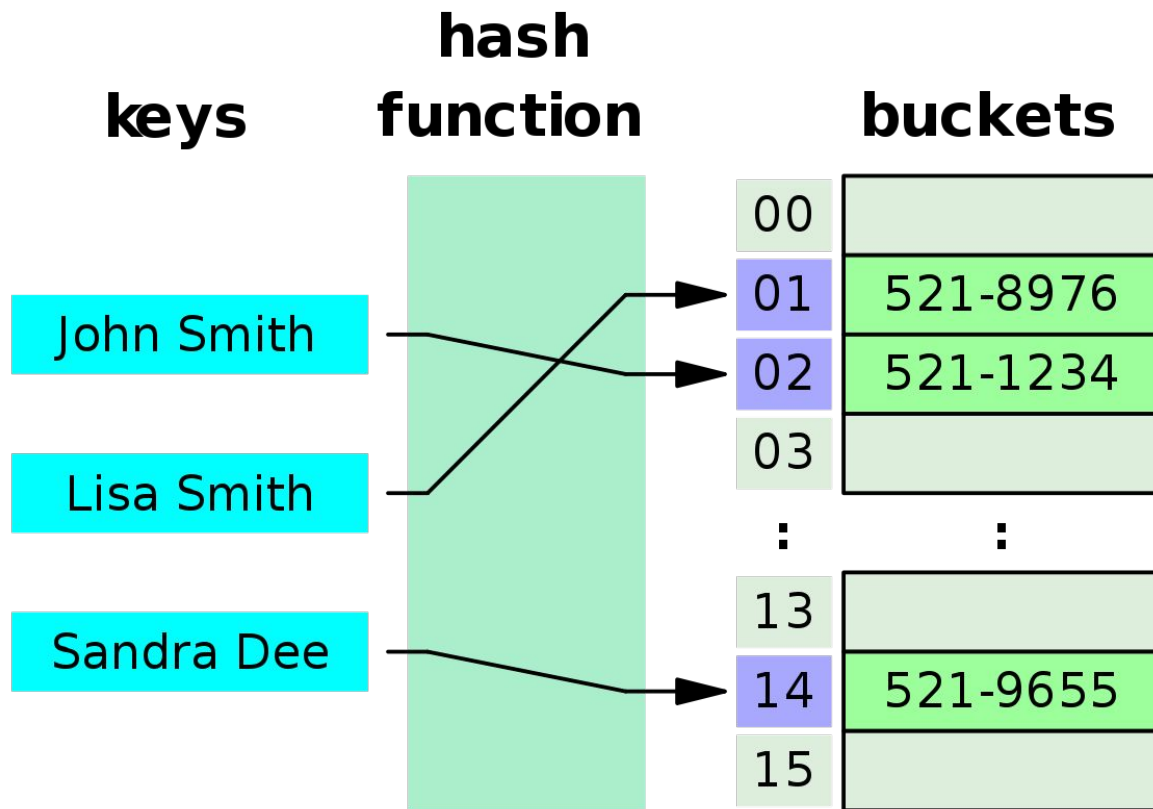




# Hash Table

Data structure used to implement cache



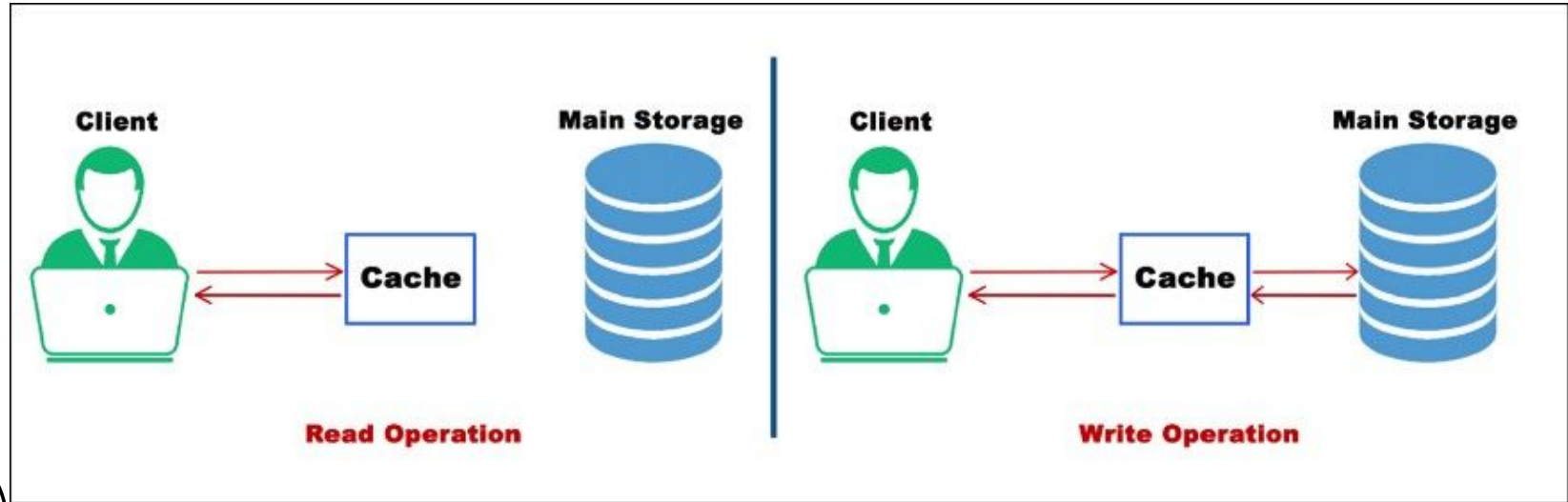


# Cache Invalidation and Modification



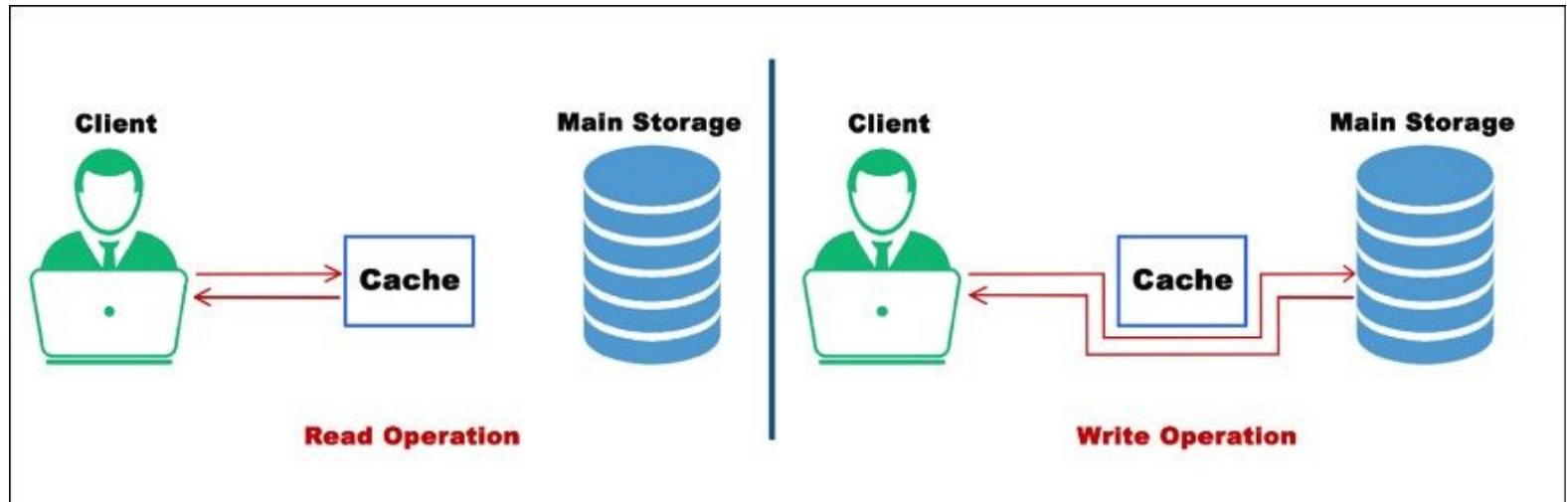
# Read-Through

The data is written into the cache and the corresponding database at the same time. This scheme maintains the complete data **consistency** between the cache and the main storage.



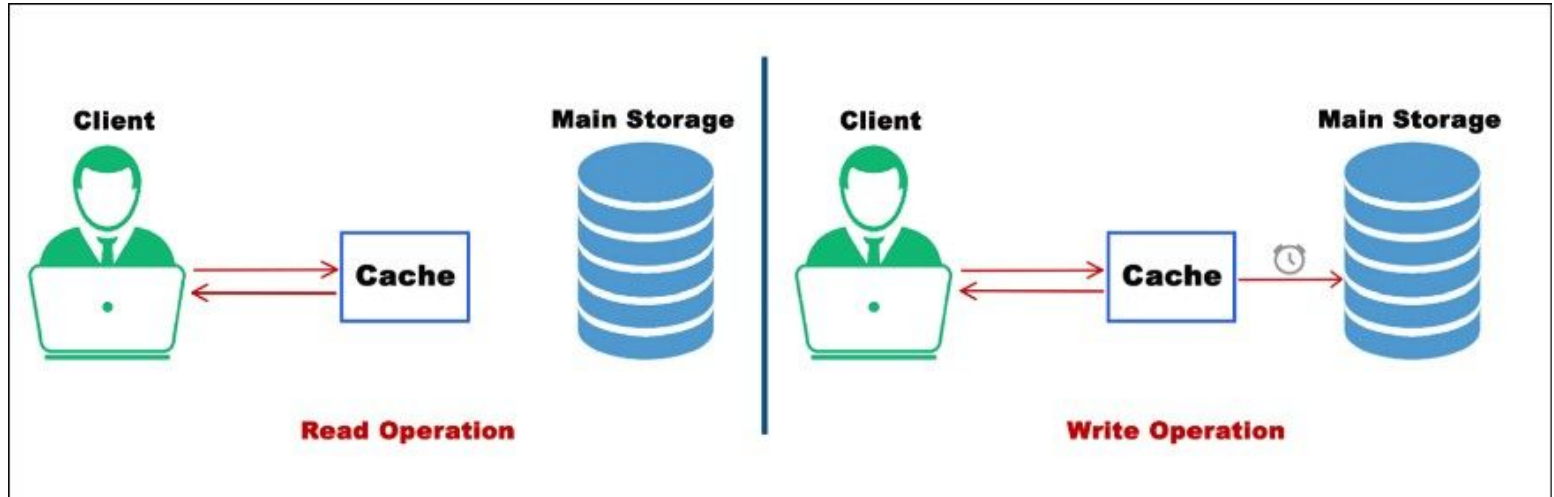
# Write-Through

In this strategy, every information directly written to the database just bypassing the cache.



# Write-Back

The data is written to cache only and completion is immediately confirmed to the client. The write to the permanent storage is done after specified intervals or under certain conditions. This results in low latency and high throughput for write-intensive applications.





# Fault Tolerance

As we know all of our hash table and data are stored in RAM, what happens if there is a power loss, all of our data goes for the toss. This means our cache system is not persistent so to make it persistent we have to do something.

- **Regular interval snapshot**
- **Log reconstruction**



# RAM is smaller than a hard disk!



# Cache Eviction Policies

- **First In First Out (FIFO):** The cache evicts the first block accessed first without any regard to how often or how many times it was accessed before.
- **Last In First Out (LIFO):** The cache evicts the block accessed most recently first without any regard to how often or how many times it was accessed before.
- **Least Recently Used (LRU):** Discards the least recently used items first.
- **Most Recently Used (MRU):** Discards, in contrast to LRU, the most recently used items first.
- **Least Frequently Used (LFU):** Counts how often an item is needed. Those that are used least often are discarded first.



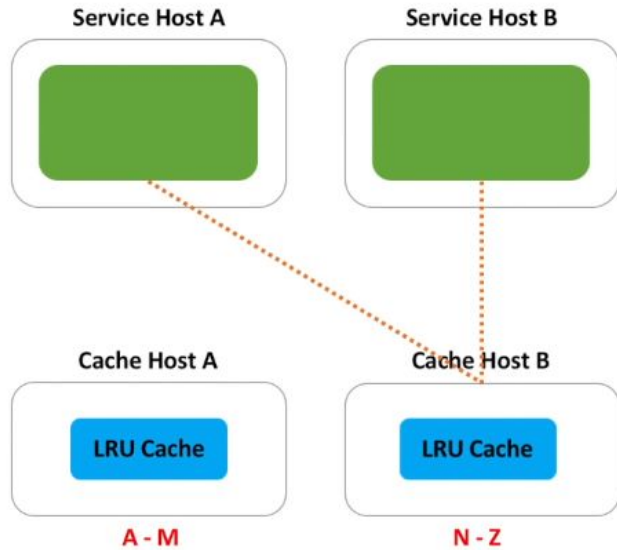
# Make it Distributed!



# Distributed Cache

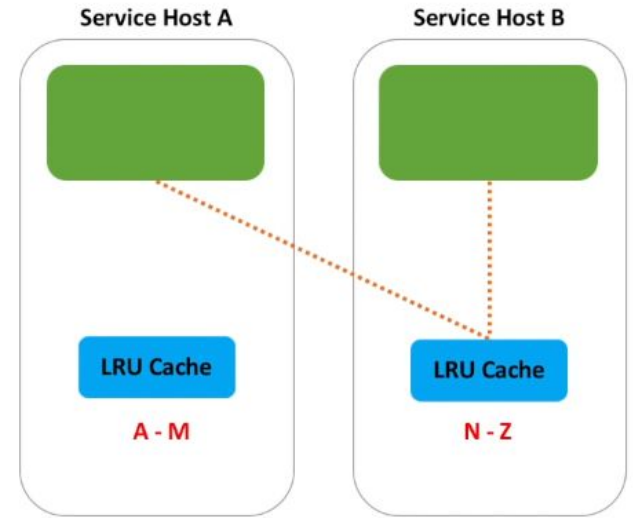
- Scalability
- High Availability
- Fault Tolerance
- Consistency





### Dedicated cache cluster

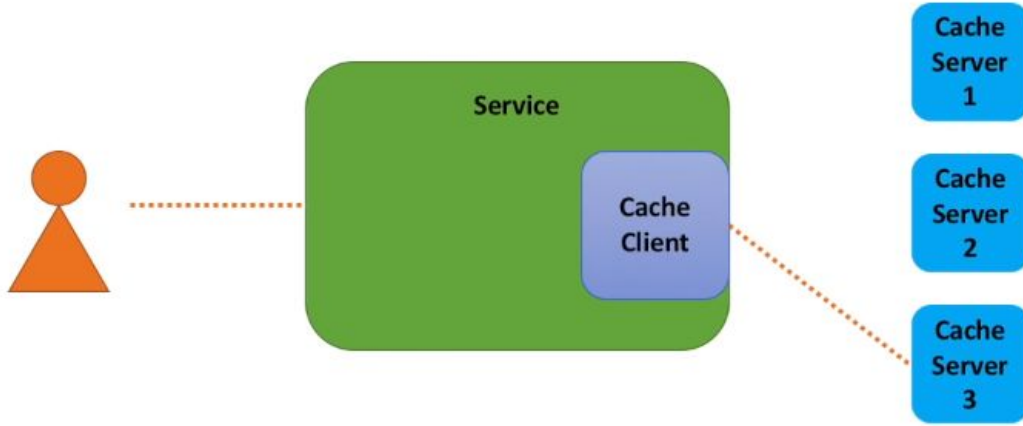
- Isolation of resources between service and cache
- Can be used by multiple services
- Flexibility in choosing hardware



### Co-located cache

- No extra hardware and operational cost
- Scales together with the service

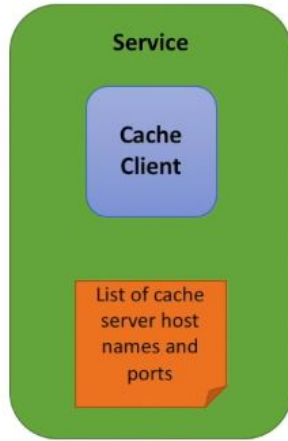




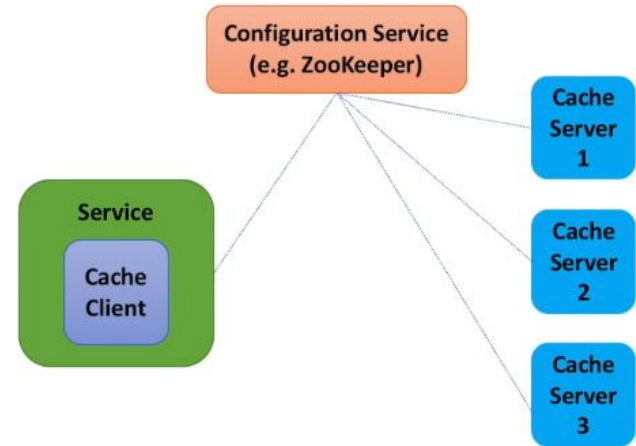
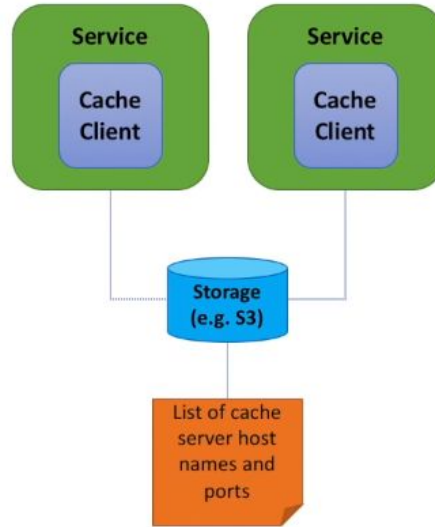
How do clients get the list of servers? How is this list updated when server dies or a new one is added?

- Cache client knows about all cache servers
- All cache clients should have the same list of servers
- Client stores list of servers in sorted order (by hash value, e.g. TreeMap in Java)
- Binary search is used to identify the server  $O(\log n)$
- Cache client uses TCP or UDP protocol to talk to servers
- If server is unavailable, client proceeds as though it was a cache miss





Use configuration management tools  
(e.g. Chef, Puppet) to deploy modified  
file to every service host



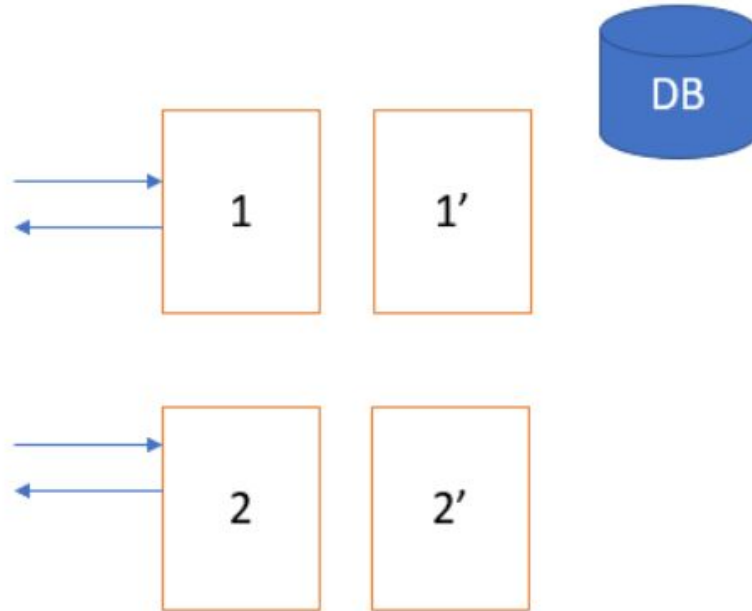


# Distributed Cache

- Scalability
- High Availability
- Fault Tolerance
- Consistency

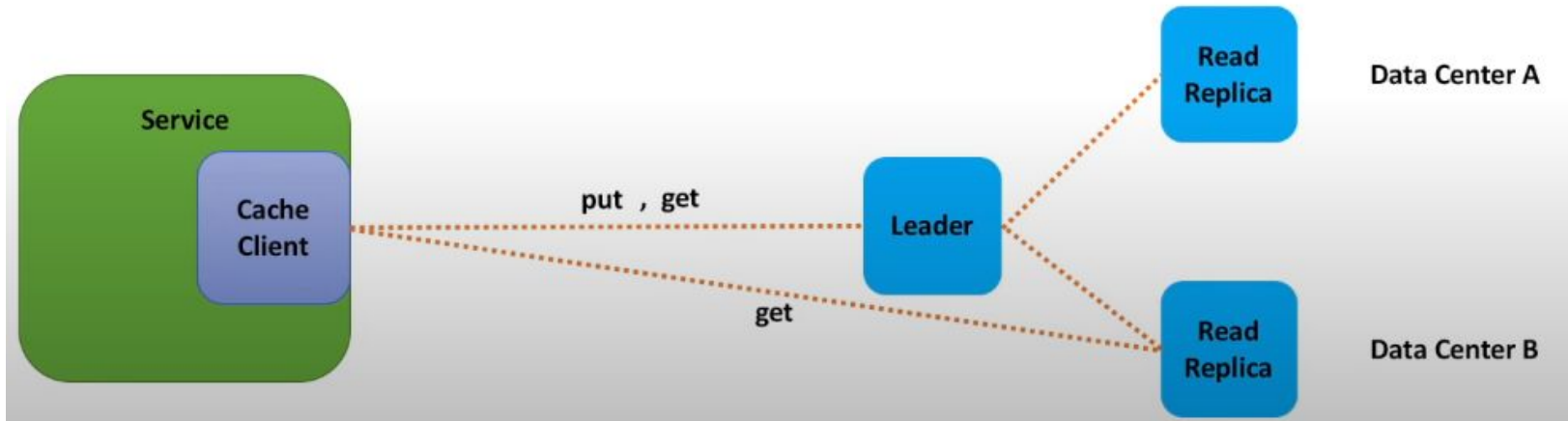
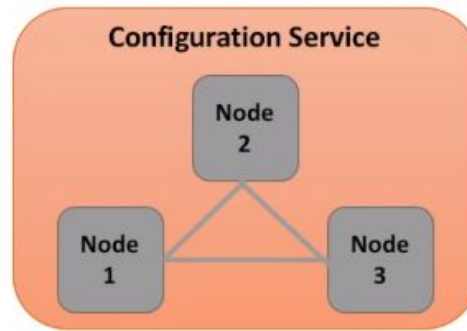


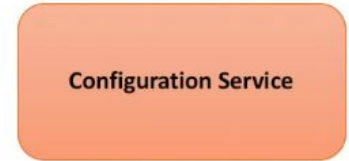
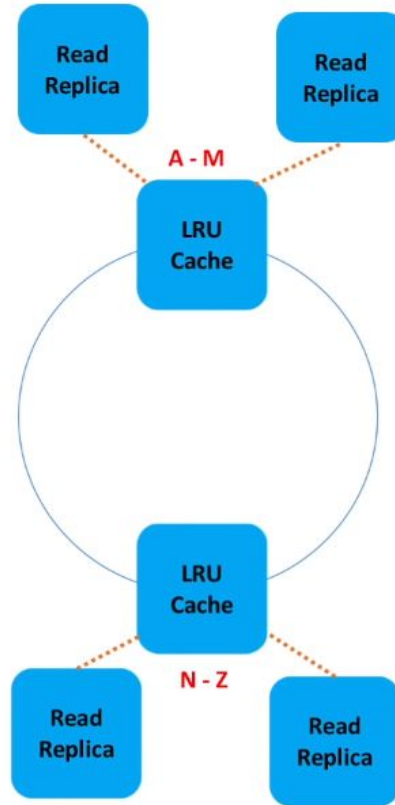
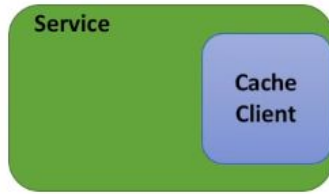
# Availability



# Consistency! 🤔







# Distributed Cache

- Scalability
- High Availability
- Fault Tolerance
- Consistency



# Split to Smaller Shards!



# Popular Distributed Cache

- Memcache
- Redis





# Reference

- <https://serhatgiydiren.github.io/system-design-interview-distributed-cache>
- <https://medium.com/system-design-concepts/distributed-cache-system-design-9560f7dd07f2>
- <https://medium.com/rtkal/distributed-cache-design-348cbe334df1>
- [https://www.youtube.com/watch?v=ql\\_g07C\\_Q5I&t=2562s](https://www.youtube.com/watch?v=ql_g07C_Q5I&t=2562s)
- 

