

# System design

# Concepts

- Scaling: Vertical vs Horizontal
- CAP Theorem
- ACID vs BASE
- Partitioning = Data Sharding
- Consistent Hashing
- Consensus Solving (e.g. Leader election)
- Optimistic vs Pessimistic locking
- Strong vs Eventual consistency
- No-SQLs DBs
- Caching
- Load Balancing: L4 vs L7
- Publishers-Subscribers vs Queues
- Map reduce
- Bloom filters, count-min sketch
- Upload vs Download (write-read)
- Metrics: CPU/RAM/Hard drive/Network
- Alerting and notification

# Tools

- **No-SQL**
    - Cassandra
    - MongoDB
  - **Blob Storage**
    - AWS S3
  - **Load balancers**
    - Nginx
    - HAProxy
  - **Cache**
    - Memcached
    - Redis
  - **Config management**
    - Zookeeper
- **Queue publisher-subscriber**
    - Kafka
  - **Search**
    - Elastic Search
    - Apache Solr
  - **MapReduce**
    - Hadoop
    - Spark

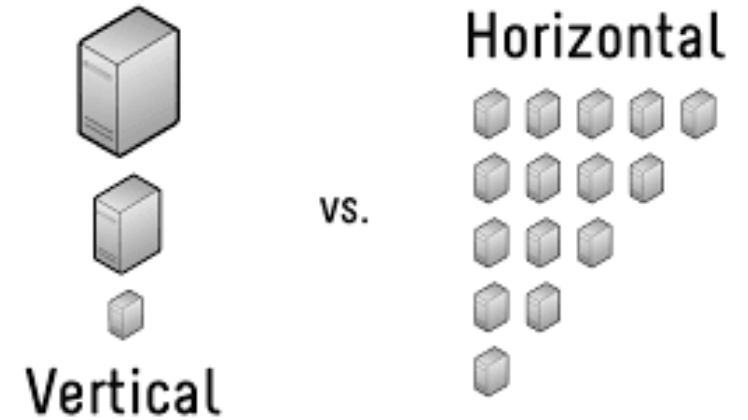
# Technical specifics

- TCP IP model: the layers
- TCP vs UDP
- https vs TLS
- IPV4 vs IPV6
- http vs http2 vs WebSockets
- Long-pooling vs WebSockets vs Server-Send Events
- REST vs SOAP
  - Get/put/post/delete/patch
- CDN, Edge
- DNS lookup
- Public key infrastructure and CA
- Symmetric vs Asymmetric

# Concepts

# Scaling

- Vertical: more RAM/CPU to existing host
  - Can be expensive
  - Limitations per max memory per single host
  - Less problem for distributed systems
- Horizontal : add another host
  - preferred over vertical scaling
  - Should handle issues with distributed systems



# CAP

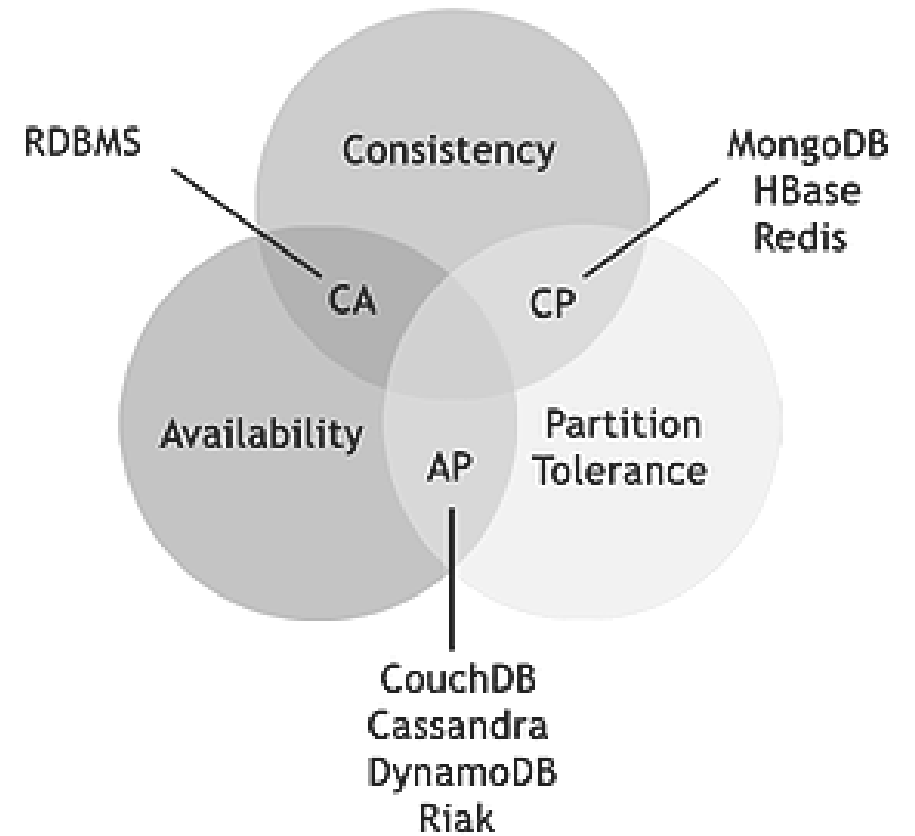
**Consistency**



**Availability**

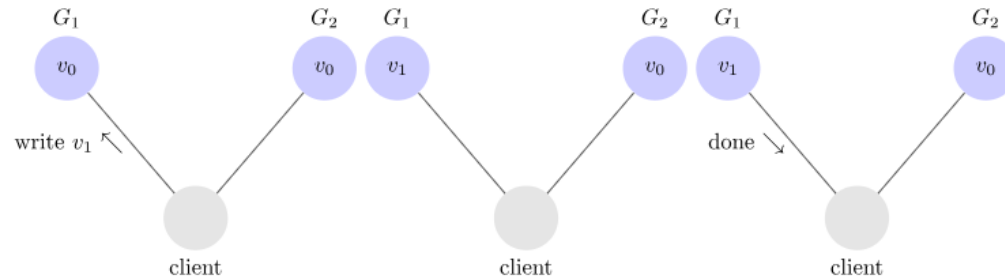


**Partition Tolerance**

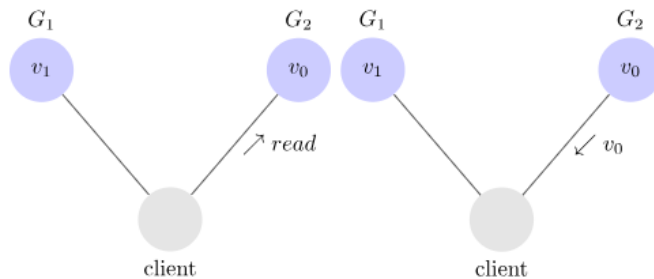


# CAP

Next, we have our client request that  $v_1$  be written to  $G_1$ . Since our system is available,  $G_1$  must respond. Since the network is partitioned, however,  $G_1$  cannot replicate its data to  $G_2$ . Gilbert and Lynch call this phase of execution  $\alpha_1$ .



Next, we have our client issue a read request to  $G_2$ . Again, since our system is available,  $G_2$  must respond. And since the network is partitioned,  $G_2$  cannot update its value from  $G_1$ . It returns  $v_0$ . Gilbert and Lynch call this phase of execution  $\alpha_2$ .



$G_2$  returns  $v_0$  to our client after the client had already written  $v_1$  to  $G_1$ . This is inconsistent.

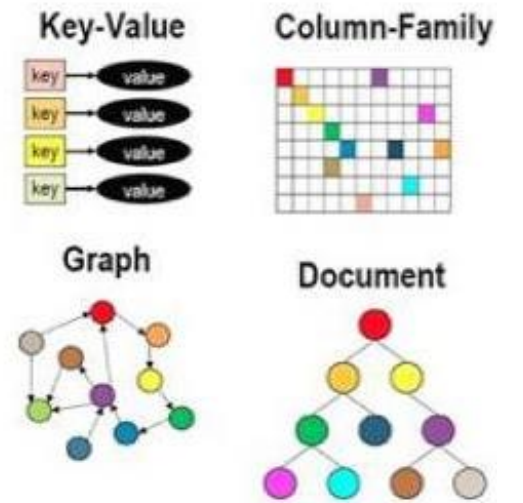
We assumed a consistent, available, partition tolerant system existed, but we just showed that there exists an execution for any such system in which the system acts inconsistently. Thus, no such system exists.

[https://mwhittaker.github.io/blog/an\\_illustrated\\_proof\\_of\\_the\\_cap\\_theorem/](https://mwhittaker.github.io/blog/an_illustrated_proof_of_the_cap_theorem/)



# No-SQL

Key-Value Persistent	Key-Value Volatile	Column	Graph	Document
Redis (CP)	memcached	Cassandra (AP)	FlockDB (twitter)	MongoDB (CP)
Membase (memcached)	Hazelcast	BigTable (Google)		CouchDB (AP)
Dynamo (AWS)		SimpleDB (AWS)		



# ACID vs BASE

## BASE (no-SQL)

- **Basically Available:** system guarantee availability in terms of the CAP
- **Soft state:** state of the system may change over time, even without input
- **Eventual consistency:** updates will eventually ripple through to all servers, given enough time.

## ACID (traditional relational DB)

- **Atomicity:** each transaction is a "unit" which either succeeds completely, or fails completely
- **Consistency:** ensures transaction can only bring the DB from one valid state to another
- **Isolation** state(concurrent transactions) = state(executed sequentially transactions)
- **Durability** once transaction is committed, it will remain committed even in the case of a system failure

# Data sharding = Partitioning

- Horizontal sharding = Range based sharding
- Vertical sharding = Different features of an entity are placed in different shards
- Key or hash based sharding = This hash value determines which database server(shard) to use.
  - if we want to add X more servers, keys would need to be remapped and migrated to new servers
  - Both new and old hash function are not valid. So requests cannot be serviced till the migration completes
  - Consistent hashing can solve this

## Drawbacks

- Database Joins become more expensive and not feasible in certain cases
- application layer needs additional level of asynchronous code and exception handling
- cross machine joins may not be an option for high availability SLA

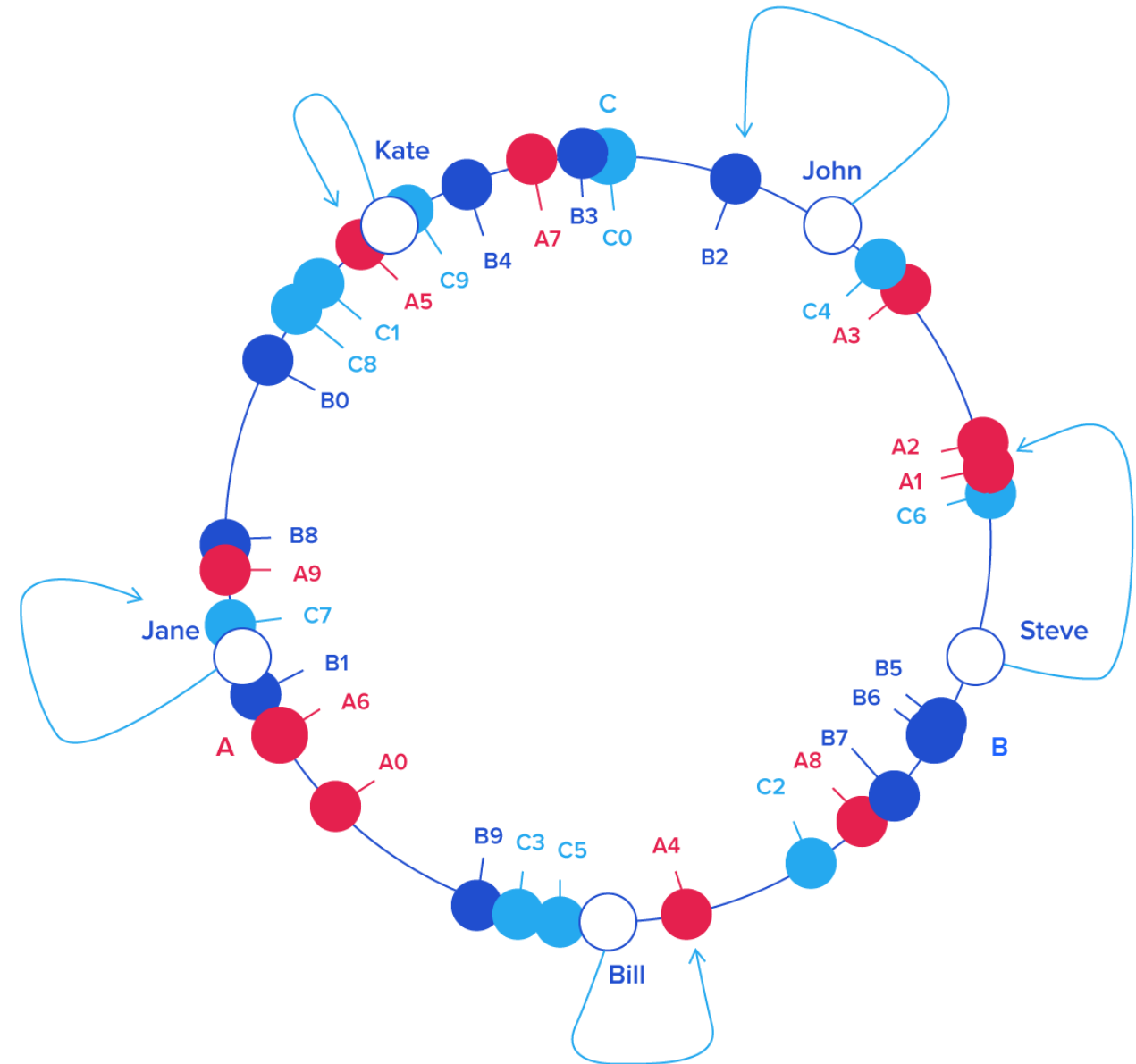
## Use sharding when

- data need to scale beyond a single storage node
- improve performance by reducing contention in a data store

# Consistent hashing

special kind of hashing: when a hash table is resized, only  $K/n$  keys need to be remapped on average

- $K$  is the number of keys
- $n$  is the number of slots.



# Pessimistic vs Optimistic locking

## **Pessimistic Locking**

- lock the record for your exclusive use until you have finished with it

This strategy is most applicable

- direct connection to the database or
- an externally available transaction ID that can be used independently of the connection
- cases when a collision is anticipated

# Pessimistic vs Optimistic locking

## Optimistic Locking

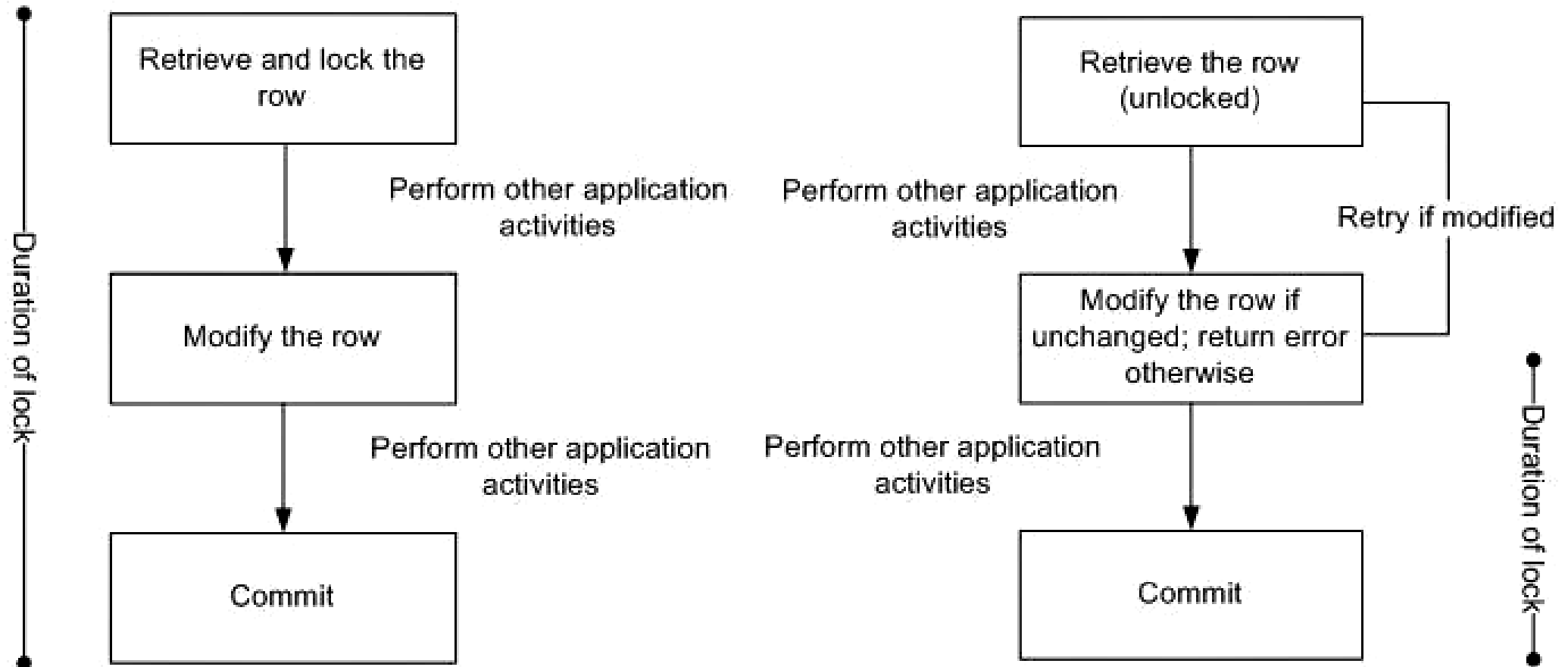
- read a record,
- take note of a version number (dates/timestamps/checksums/hashes)
- check that the version hasn't changed before you do a write operation
- Write the record
- filter the update on the version to make sure it's atomic and update the version in one hit.
  - Record should not be updated between when you check the version and write the record

If the record is dirty (i.e. different version to yours) you abort the transaction and the user can re-start it.

This strategy is most applicable for

- high-volume systems where you do not necessarily maintain a connection to the DB
- cases when you don't expect many collisions.

# Pessimistic vs Optimistic locking



# Consistency: Eventual, Strong Eventual, Strong

Conflicts arise because each node can update its own copy. If we read the data from different nodes we will see different values

## **Strong consistent**

Data will get passed on to all the replicas as soon as a write request comes to one of the replicas of the database. But during the time these replicas are being updated with new data, response to any subsequent read/write requests by any of the replicas will get delayed as all replicas are busy in keeping each other consistent

## **Eventual consistency**

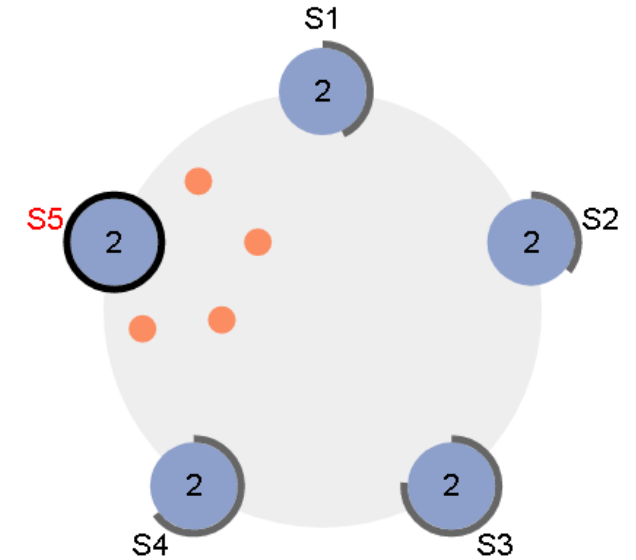
Used in distributed computing to achieve high availability that informally guarantees that, if no new updates are made to a given data item, **eventually** all accesses to that item will return the last updated value



# Consensus Solving

Leader election  
Log replication

Paxos vs RAFT

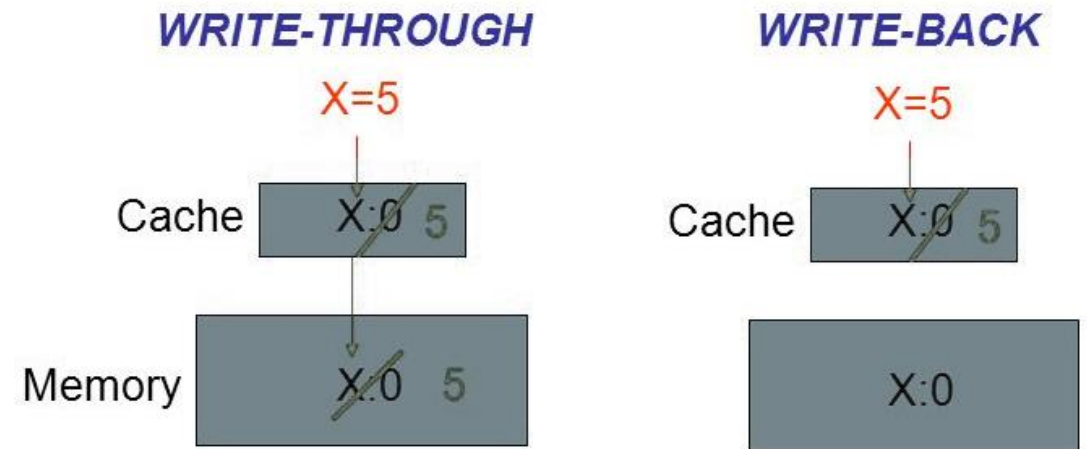


<http://thesecretlivesofdata.com/raft/>

<https://www.youtube.com/watch?v=IE3r12vBx6I>

# Caching

- Should be small to fit in memory
- Is not a source of Truth
- Should comply with caching eviction policy  
FIFO, LIFO, Least Recently Used (LRU), Most Recently Used (MRU), Least Frequently Used (LFU), Random Replacement (RR).
- **write-through:** write to the cache causes a synchronous write to the backing store
- **write-back:** writes are not immediately mirrored to DB store. Instead, the cache tracks which of its locations have been written over and marks these locations as dirty. The data in these locations is written back to the backing store when those data are evicted from the cache, an effect referred to as a lazy write.
- Reduce DB/network load (get recent data from cache)
- Reduce recalculations (e.g. aggregations)



# Load balancing: methods

- Round robin: an incoming request is routed to each available server in a sequential manner.
- Weighted round robin: a static weight is preassigned to each server
- Least connection: This method reduces the overload of a server by assigning an incoming request to a server with the lowest number of connections currently maintained.
- Weighted least connection: In this method, a weight is added to a server depending on its capacity. This weight is used with the least connection method to determine the load allocated to each server.
- Least connection slow start time -- Here, a ramp-up time is specified for a server using least connection scheduling to ensure that the server is not overloaded on startup.
- Agent-based adaptive balancing -- This is an adaptive method that regularly checks a server irrespective of its weight to schedule the traffic in real time.
- Fixed weighted -- In this method, the weight of each server is preassigned and most of the requests are routed to the server with the highest priority. If the server with the highest priority fails, the server that has the second highest priority takes over the services.
- Weighted response -- Here, the response time from each server is used to calculate its weight.
- Source IP hash -- In this method, an IP hash is used to find the server that must attend to a request.

# Load balancing: L4 vs L7

Open systems interconnection - Transmission Control Protocol - Internet protocol

	OSI Layer	TCP/IP	Datagrams are called
Software	Layer 7 Application	HTTP, SMTP, IMAP, SNMP, POP3, FTP	Upper Layer Data
	Layer 6 Presentation	ASCII Characters, MPEG, SSL, TSL, Compression (Encryption & Decryption)	
	Layer 5 Session	NetBIOS, SAP, Handshaking connection	
	Layer 4 Transport	TCP, UDP	Segment
	Layer 3 Network	IPv4, IPv6, ICMP, IP <u>Sec</u> , MPLS, ARP	Packet
Hardware	Layer 2 Data Link	Ethernet, 802.1x, PPP, ATM, <u>Fiber</u> Channel, MPLS, FDDI, MAC Addresses	Frame
	Layer 1 Physical	Cables, Connectors, Hubs (DLS, RS232, 10BaseT, 100BaseTX, ISDN, T1)	Bits

# Load balancing: L4 vs L7

**L4:** directs traffic based on data from network and transport layer protocols, such as IP address and TCP port.

**L7:** adds content switching to load balancing. This allows routing decisions based on attributes like HTTP header, uniform resource identifier, SSL session ID and HTML form data. Most of LBs are working like this

<https://www.educative.io/collection/page/5668639101419520/5649050225344512/5747976207073280>

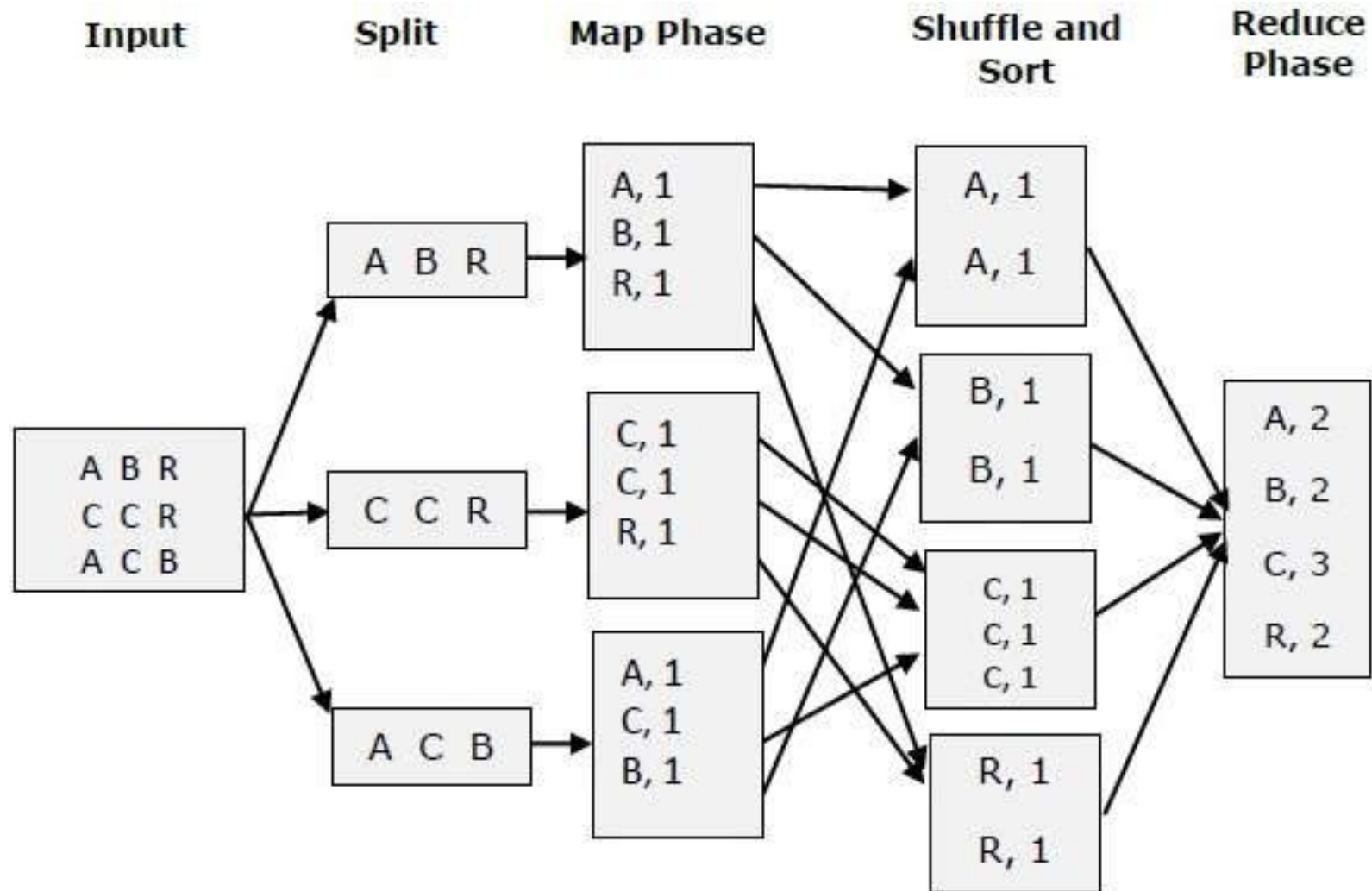
<https://blog.envoyproxy.io/introduction-to-modern-network-load-balancing-and-proxying-a57f6ff80236>

# Publish-Subscribe vs Queues

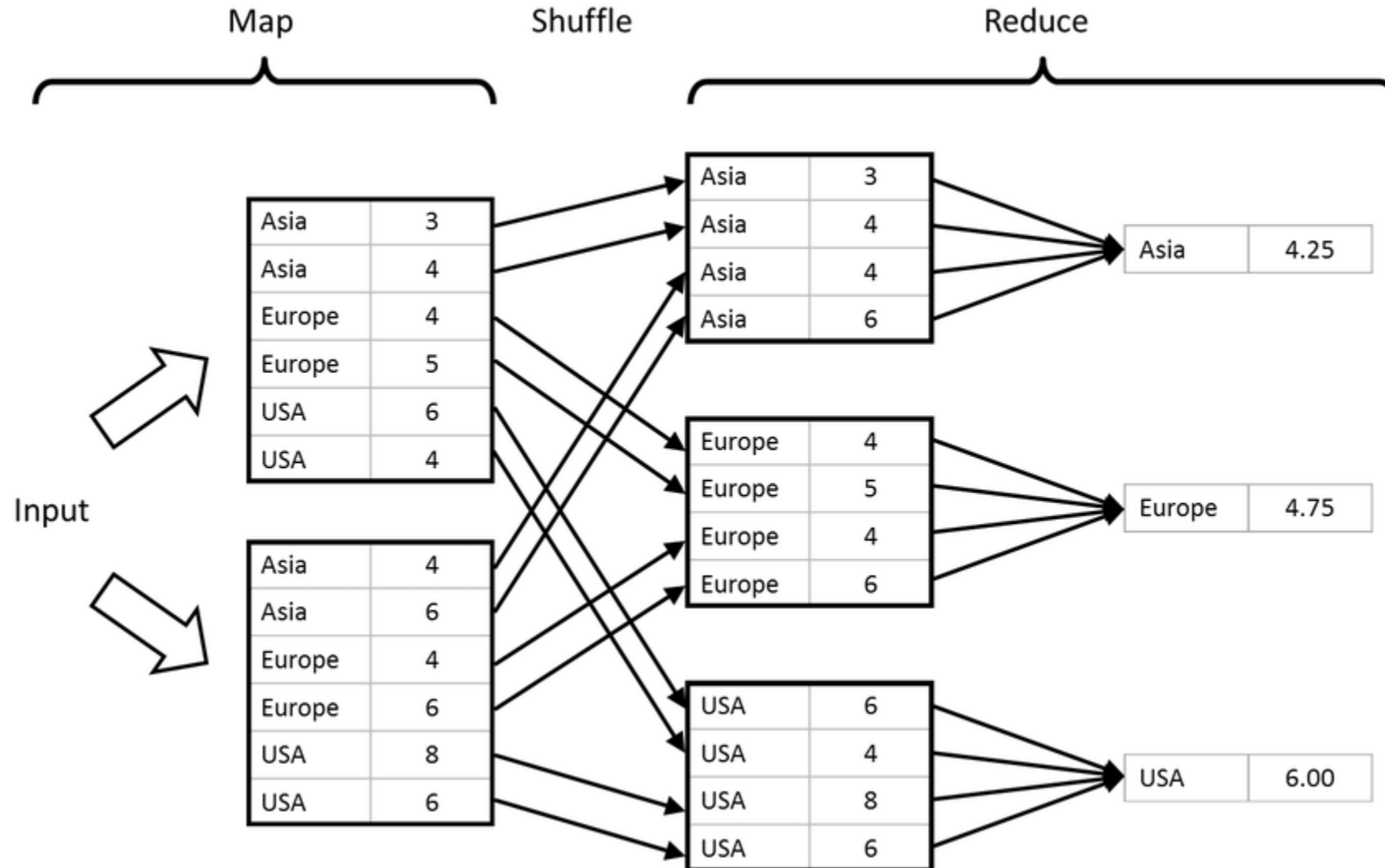
A message queue receives incoming messages and ensures that each message for a given topic or channel is delivered to and processed by exactly one consumer. Message queues can support high rates of consumption by adding multiple consumers for each topic, but only one consumer will receive each message on the topic. Which consumer receives which message is determined by the implementation of the message queue. To ensure that a message is only processed by one consumer, each message is deleted from the queue once it has been received and processed by a consumer (i.e. once a consumer has acknowledged consumption of the message to the messaging system).

in contrast to message queuing, publish-subscribe messaging allows multiple consumers to receive each message in a topic. Further, pub-sub messaging ensures that each consumer receives messages in a topic in the exact order in which they were received by the messaging system

# Map-Reduce

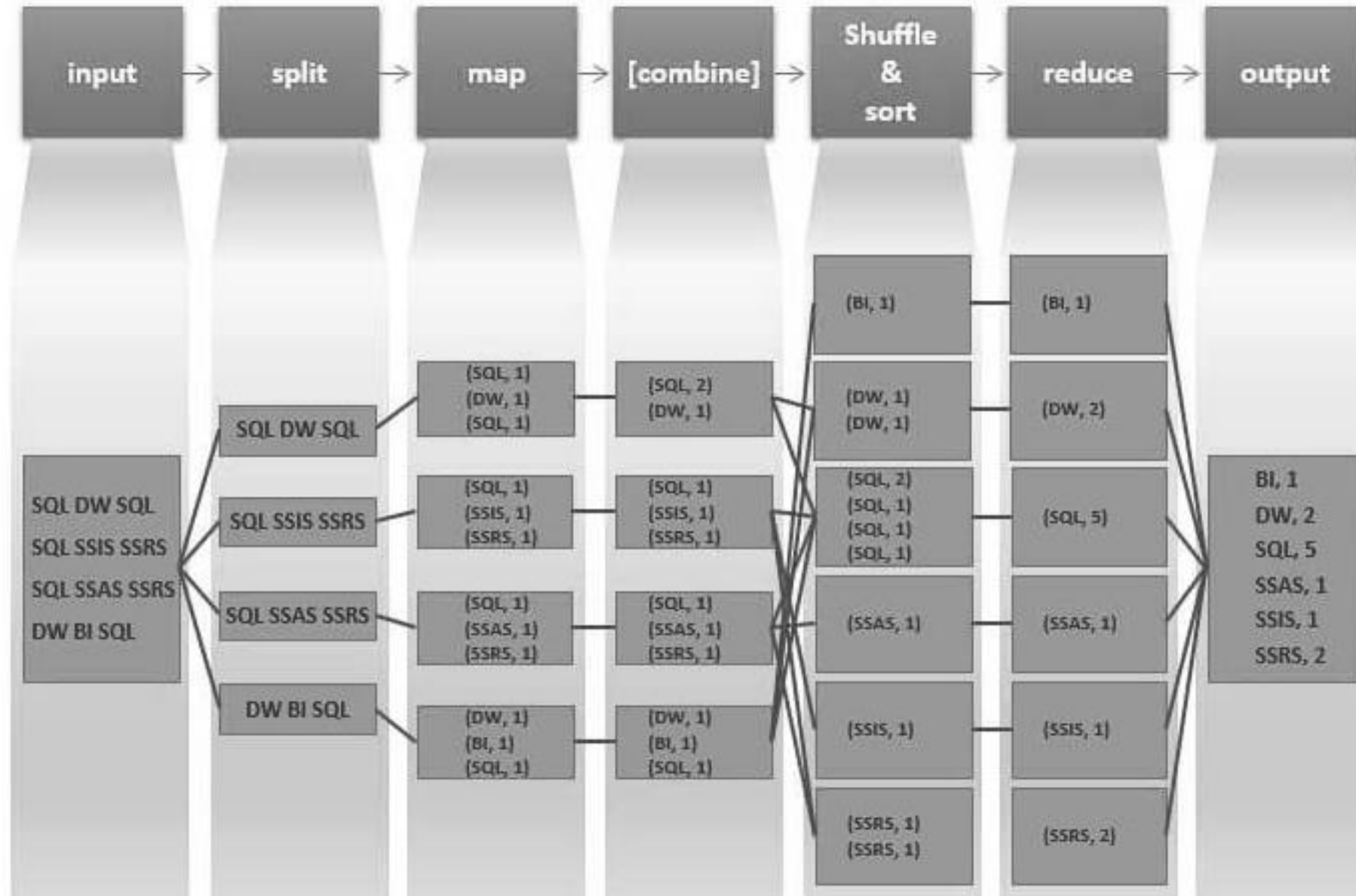


# Map-Reduce





# Map-Reduce



# Tools

# No-SQL Cassandra: column based

## Cassandra

- column/tab based (goes well with the historical RDBMS)
- **Availability-Partition**
- support HiveQL (SQL like syntax)
- Supports sharding
- Preferable option when # of writes > # of reads
- Great for both key-value and time series data
- Can provide both eventual and strong consistency



# No-SQL: Redis key-value

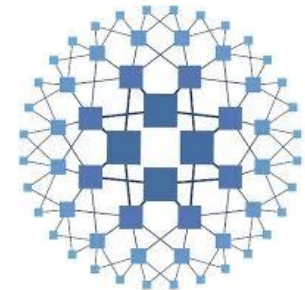
## Redis

- key-value based
- **Consistency-Partition**
- great for rapid changing data



# Load balancing: nginx vs HAProxy

1000+ connections



# Cache

- Memcached vs Redis vs Hazelcast



# Search:

ElasticSearch

Solr



# Blob Storage:

Amazon S3

S3 vs EC2





# Zookeeper

## Centralized configuration management

- Distributed locking
- Leader election
- Scales great for reads
- Scales worse for writes



Apache  
Zookeeper

# Messaging

## **Fault-tolerant highly available queue for publisher-subscriber**

- Can deliver message exactly once
- Keep messages ordered inside the partition

### Queues

[Apache ActiveMQ](#), [Amazon SQS](#), [RabbitMQ](#)

### Pub-Sub

[Apache Kafka](#), [Google Cloud Pub/Sub](#)



# Map-Reduce

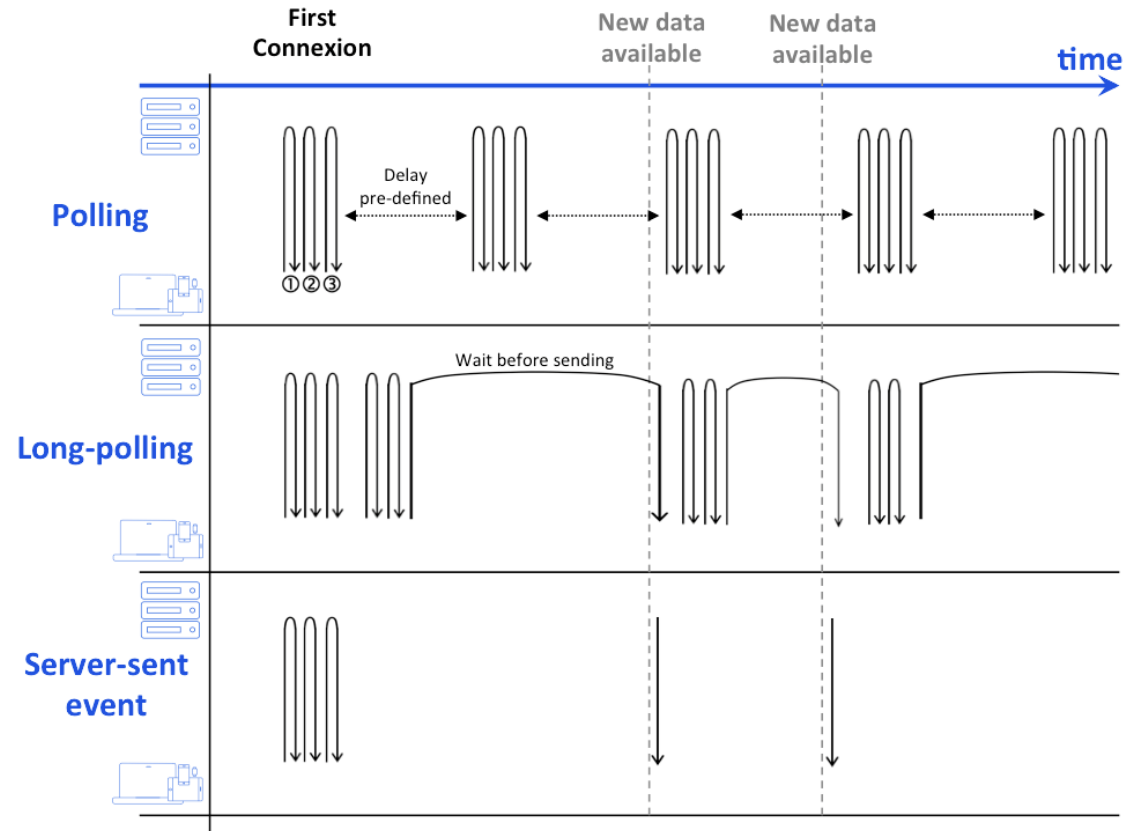
- Hadoop
- Spark: in-memory



# Technical specifics

# Long-pooling vs WebSockets vs Server-Send Events

- **Long/short polling (client pull)**
- client asking server for updates at certain regular intervals
- *Short polling* is an AJAX-based timer that calls at fixed delays
- *Long polling* is based on [Comet](#)
- **WebSockets** (server push)
- server is proactively pushing updates to the client (reverse of client pull)
- **Server-Sent Events** (server push)



# REST: Representational State Transfer

## Key principles:

- this is a client-server architecture
- **It's stateless:** *communication between the client and the server always contains all the information needed to perform the request.*
- There is **no session state** in the server, it is kept entirely on the client's side. *(e.g. if access to a resource requires authentication, then the client needs to authenticate itself with every request)*
- Cacheable
- provides a uniform interface between components.

Task	Method	Path
Create a new task	POST	/tasks
Delete an existing task	DELETE	/tasks/{id}
Get a specific task	GET	/tasks/{id}
Search for tasks	GET	/tasks
Update an existing task	PUT	/tasks/{id}

# Practice

- Designing URL Shortening service
- Designing Instagram
- Designing Twitter
- Designing Dropbox
- Designing Youtube
- Design a Parking Lot

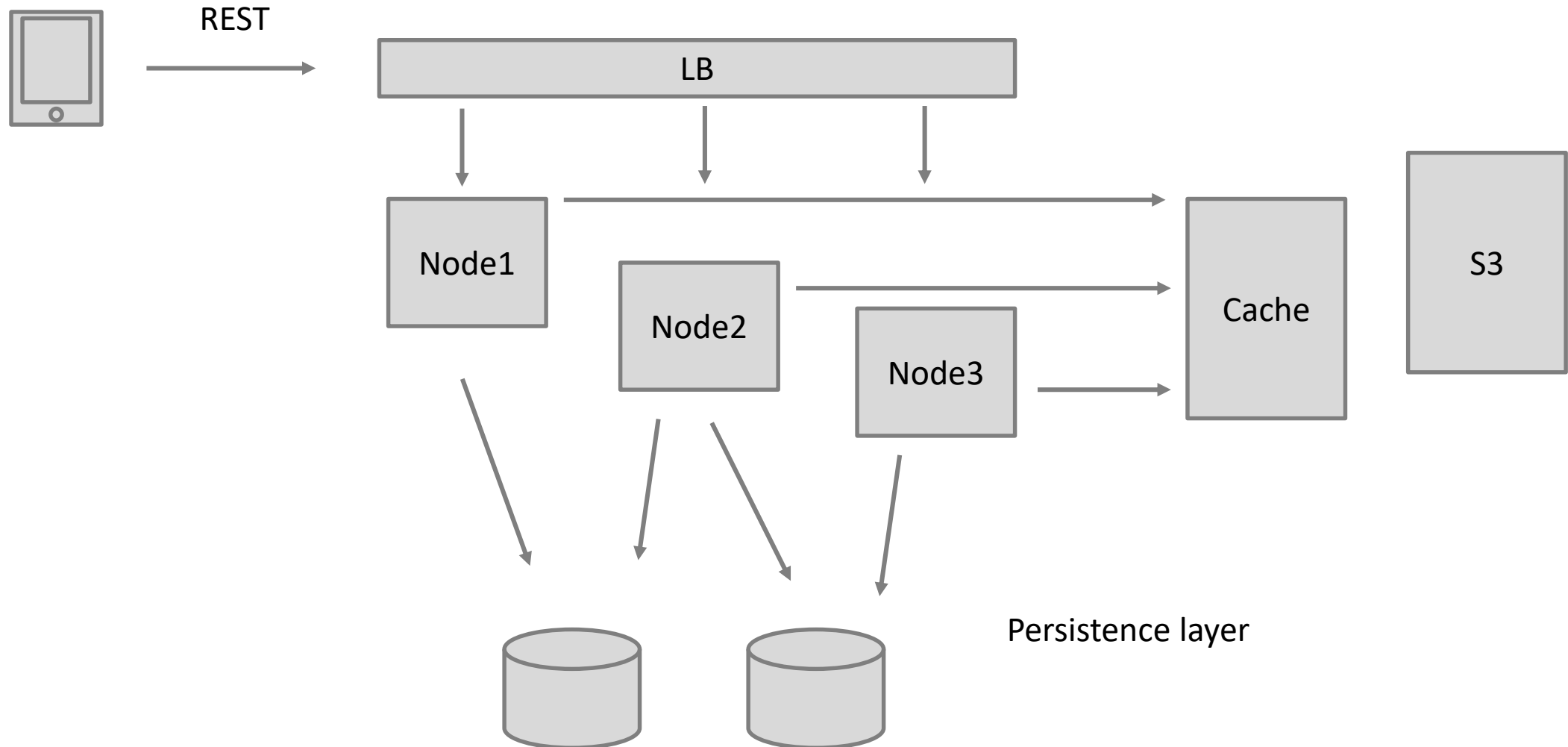


# The approach

- Requirements and Features: functional, non-functional
- APIs
- Durability, Availability, Performance
- Capacity estimation
- Scalability
- High level design, Entity diagram
- Data model
- Detailed design of each component
- Bottlenecks
- Security / Privacy
- Cost



# The approach



# Shortening URL

<https://www.educative.io/collection/page/5668639101419520/5649050225344512/5668600916475904>

<https://www.quora.com/What-is-the-best-way-to-prepare-for-a-System-Design-interview-for-Amazon>

[https://www.youtube.com/user/tusharroy2525/videos?view=0&sort=dd&shelf\\_id=3](https://www.youtube.com/user/tusharroy2525/videos?view=0&sort=dd&shelf_id=3)

# Shortening URL: Requirements

## **Functional**

- given a URL, generate a unique short link
- given a short link, provide the original URL

## **Non-functional**

- The system should be highly available
- Shortened links should not be predictable (not easy to guess)

# Shortening URL: APIs

get  
get

# TODO

---

- xxx