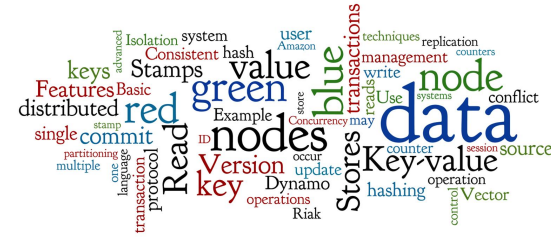




# Agenda



- **Fundamentals** of Key-value Stores
  - Basic Example: Riak
- Key Techniques of Many Key-value Stores
  - Data Sharding: **Consistent hashing** + virtual nodes
  - **Replica** management & consistency (version stamps)
  - **Gossip protocols** (distributed management of nodes)
  - Transactions: Two-phase commit protocol (2PC); MVCC
- **Comparison** of K-V Stores and Applicability
  - **Features** to consider - basic, advanced
  - Modes of **communication** with the database
  - When **(not) to use** Key-value Stores

# Key-value Stores: Basics



- A simple **hash table** (map), primarily used when all accesses to the database are via **primary key**
  - **key-value** mapping
- In RDBMS world: A table with two columns:
  - ID column (**primary key**)
  - DATA column storing the value (unstructured BLOB)
- Basic **operations**:
  - **Put** a value for a key `put(key, value)`
  - **Get** the value for the key `value := get(key)`
  - **Delete** a key-value `delete(key)`

# Querying



- We can **query by the key**
- To query using some **attribute** of the value is **not possible** (in general)
  - We need to read the value to test any query condition
- What if we **do not know** the key?
  - Some systems support additional functionality
    - Using some kind of additional **index** (e.g., full text)
    - The data must be indexed first
    - Example later: Riak search

















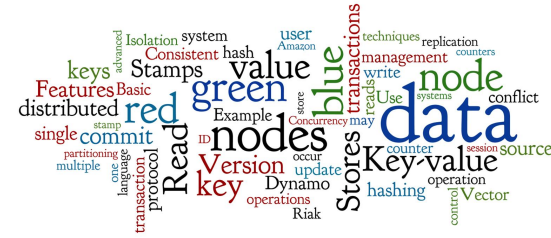




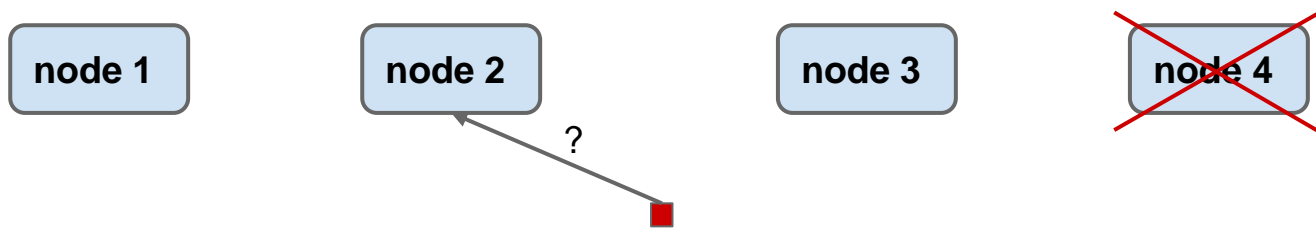




# Sharding: Modulo Hashing



- We want to use  $\text{hash}(\text{key})$  for partitioning of key-value pairs to nodes (**auto sharding**)
- Standard **modulo-based** hashing:



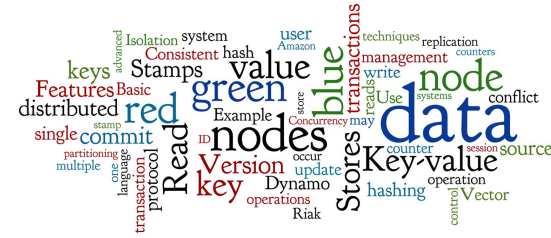
$$\begin{aligned} \text{node}(\text{key}) &= \text{hash}(\text{key}) \bmod (\#\_of\_nodes) \\ \text{node}(\text{key}) &= \text{hash}(\text{key}) \bmod (\#\_of\_nodes - 1) \end{aligned}$$

- **Recalculate** the hashes of all objects, if  $\#\_of\_nodes$  changes
  - and **migrate** practically **all** data objects to different nodes





# Sharding by Hashing



- **Consistent** hashing
  - is used in massively **distributed** systems (like Riak)
- **Modulo**-based hashing
  - is also used, e.g., in Solr and Lucene
- **Modulo** hashing is good for keeping data **balanced**
  - **consistent** hashing **cannot** guarantee balanced data
    - especially for low number of nodes
  - it must use **different** techniques to achieve **balancing**

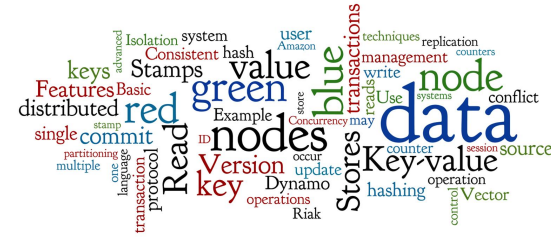








# Quora per Operation



- The **R/W** values can be often set per operation
  - Riak: all / one / quorum / an integer value
  - it is a way to tune efficiency/availability vs. consistency
- example: **N** = 3, quora **W** = 2, **R** = 2
  - value for `key1` is stored on `nodeA`, `nodeB`, `nodeC`
  - at least **two of them** always have the **newest** value
    - and operation `get(key1)` will **always get** the newest value
- we can set **R** = 1 for operation `value := get(key1)`
  - meaning: get the value the from **any** replica, e.g., `nodeB`
  - **even though** `nodeA` and `nodeC` may have a newer value















# Vector Stamps Algorithms



## Vector stamps

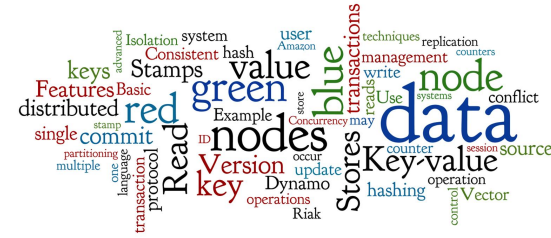
- Family of algorithms for generating a **partial ordering** of events in a **distributed system** and detecting “**conflicts**”.
- **Each node** has its own **counter**
  - for each data item
  - The node’s **counter increments** when its value is updated
- Each node **keeps** a **counter vector** with counters of **all nodes**
  - The nodes **exchange** their values
- Each node uses the **counter vectors** to determine
  - which value is new
  - if there is a conflict







# Conflict Resolution



- There are three general ways to **resolve conflicts**
  - (**reconcile differences** between copies of distributed data)
  - this process is often known as **anti-entropy**
- 1. **Write repair**
  - The correction takes place during a write operation
- 2. **Read repair**
  - The correction is done when a **read finds an inconsistency**
    - Optimistic strategy, read operation is slowed down
- 3. **Asynchronous repair**
  - The correction is done as separate operations
  - AKA **active** “anti-entropy”

# Gossip Protocols



A set of **distributed** protocols

- Each node **periodically sends** its current info
  - To a **randomly**-selected peer
  - The peers keep the newer info

In distributed NoSQL databases, gossip is used for

- **Spreading** information about **current** state
  - of the entering/leaving/failing **nodes**
  - asynchronous **reconciling** of conflicts (anti-entropy)
  - other properties, ...













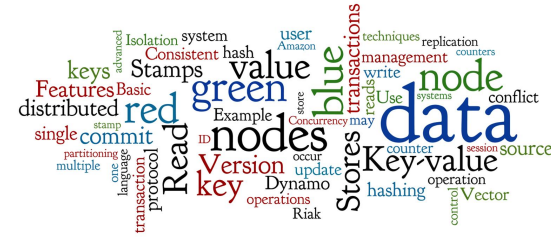






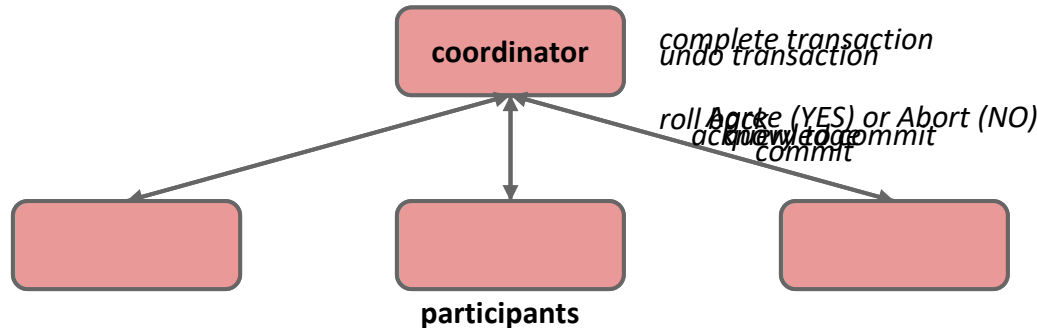


# Two-phase Commit Protocol



- 2PC: Distributed algorithm

- coordinating all participants in a distributed transaction
- on whether to commit or abort (roll back) the transaction
  - it's a special type of consensus protocol



1. Commit request phase (voting phase)
2. Commit phase
  - a. SUCCESS (agreement from all)
  - b. FAILURE (abort from any)

1. Execute transactions up to COMMIT phase
2. Write entry to UNDO and REDO logs

1. Complete operation
  2. Release locks
1. Undo operation
  2. Release locks













# Internal Features (1)



- **Durability**

- if the system supports data **persistence**
- and how is it done (storage models)

- **Data Partitioning (Sharding)**

- if the system supports (semi)-**automatic** data sharding

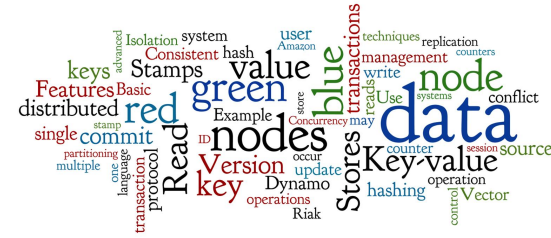
- **Data Replication**

- replication can **speed up** read/write and provide reliability
- **master-slave**, P2P, R/W quora, version control, etc.





# Advanced Features (2)



- **Distributed Transactions Management**
  - is implemented any concept of **transaction management**
  - like X/Open XA (eXtended Architecture)
- **Map-Reduce processing**
  - is available some **distributed** operation **execution** like M-R
- **Triggers**
  - procedures started **automatically** when something happens



