

NoSQL and Big Data Lecture

Felix Gessert, Norbert Ritter

fg@baqend.com

May 5, 2019



Universität Hamburg



Outline



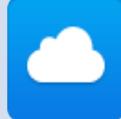
NoSQL Foundations and Motivation



The NoSQL Toolbox: Common Techniques



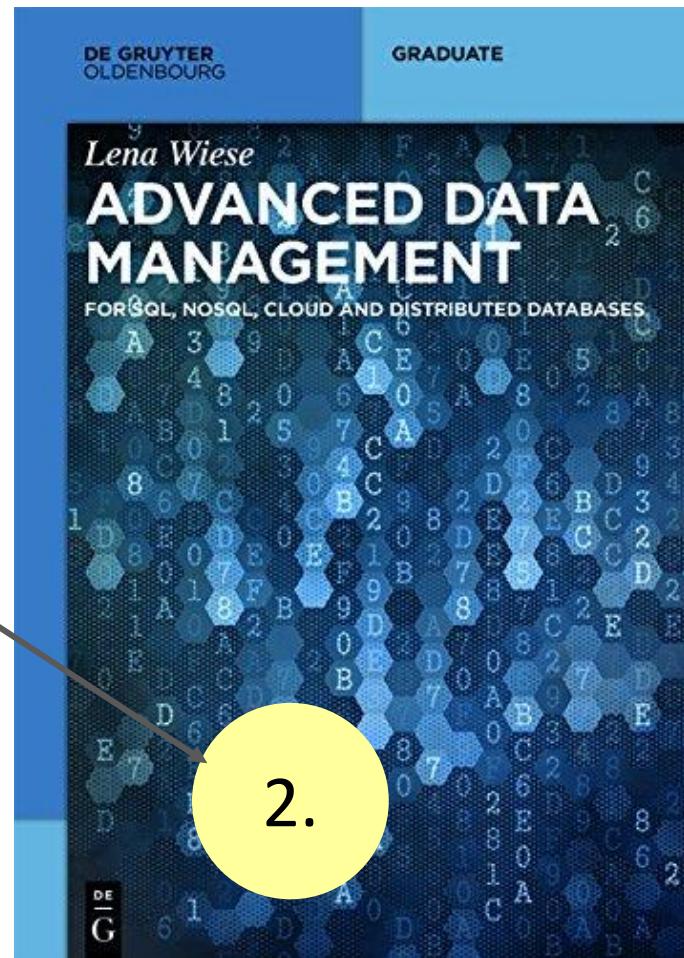
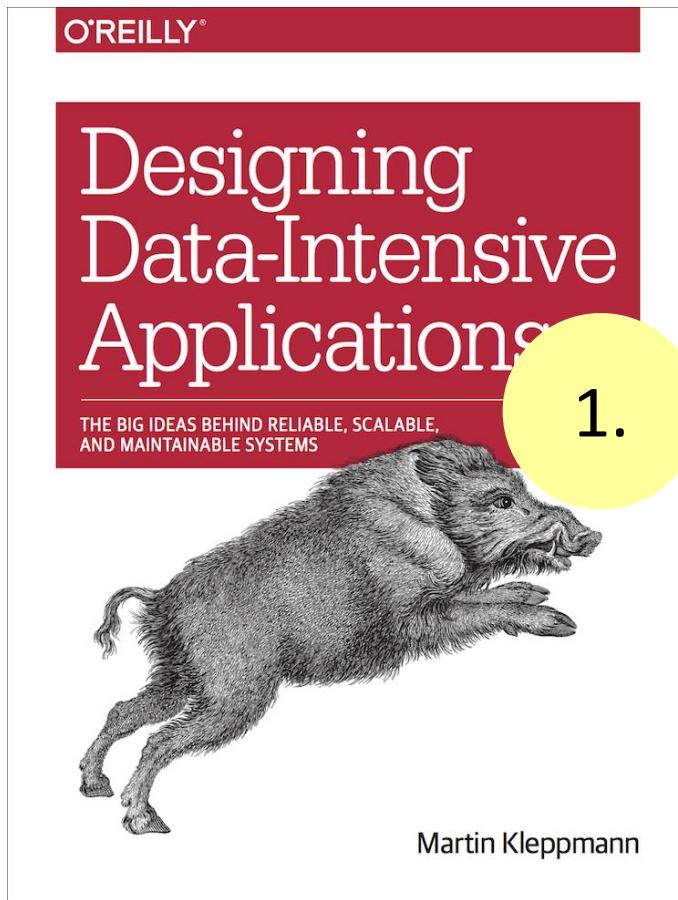
NoSQL Systems & Decision Guidance



Cloud Databases

- Literature
- The Database Explosion
- NoSQL: Motivation and Origins
- The 4 Classes of NoSQL Databases:
 - Key-Value Stores
 - Wide-Column Stores
 - Document Stores
 - Graph Databases
- CAP Theorem

Recommended Literature: NoSQL Books



Recommended Literature: Blogs



<http://www.infoq.com/nosql/>



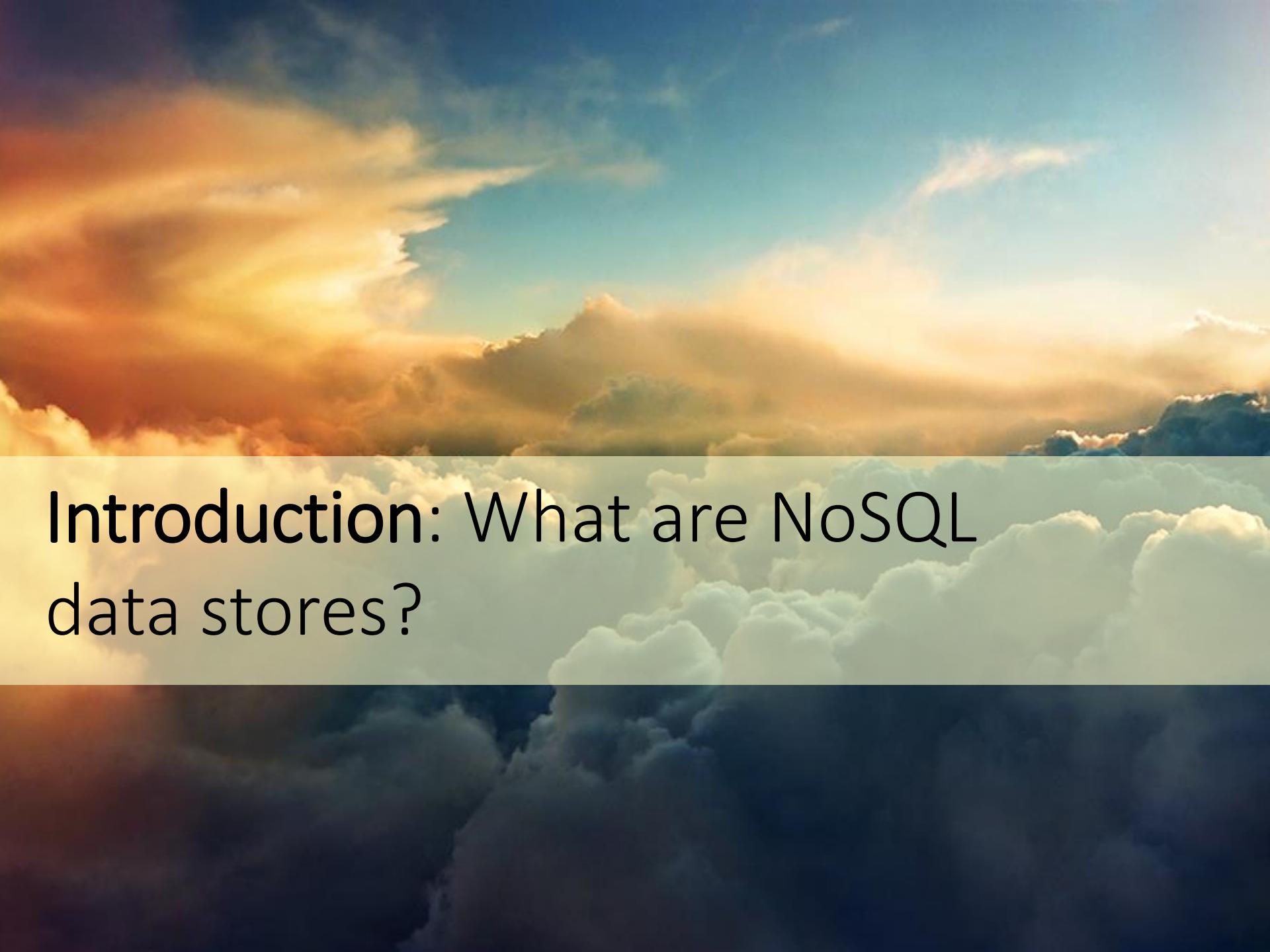
<http://blog.baqend.com/>



<http://www.dzone.com/mz/nosql>



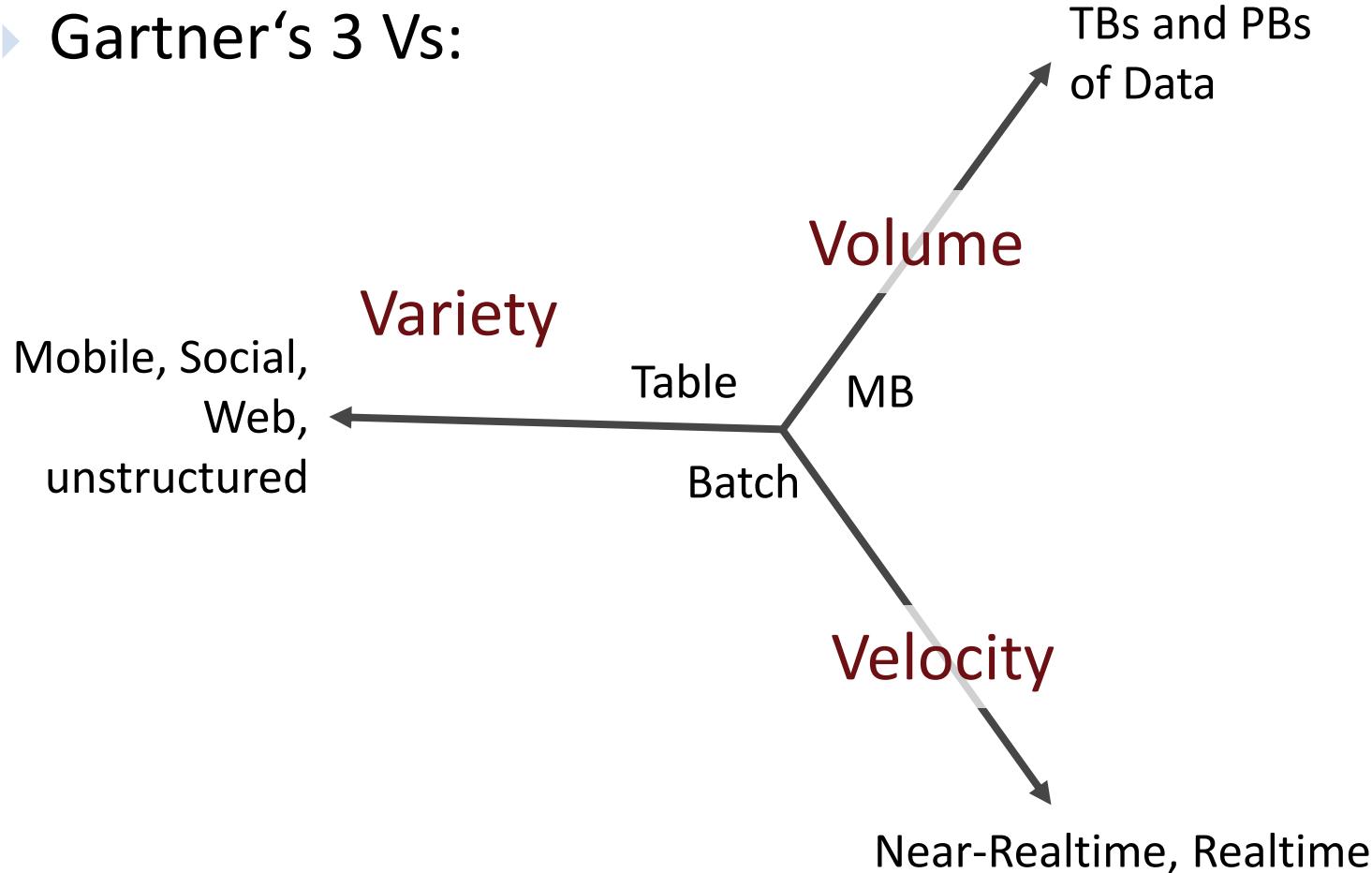
<http://highscalability.com/>

The background of the slide features a wide-angle photograph of a sky at either dawn or dusk. The upper portion of the image is dominated by a vibrant orange and yellow glow from the sun, which is partially visible on the left. This light filters through and illuminates various layers of clouds. Below this, the sky transitions into a deep, rich blue. In the lower half of the image, there are large, billowing white and grey clouds that appear to be moving across the frame. The overall effect is one of a dynamic, natural landscape.

Introduction: What are NoSQL data stores?

Motivation: Big Data

- ▶ Gartner's 3 Vs:

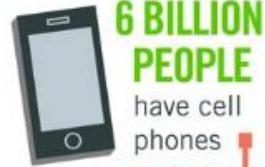


4Vs

40 ZETTABYTES

[43 TRILLION GIGABYTES]

of data will be created by 2020, an increase of 300 times from 2005.



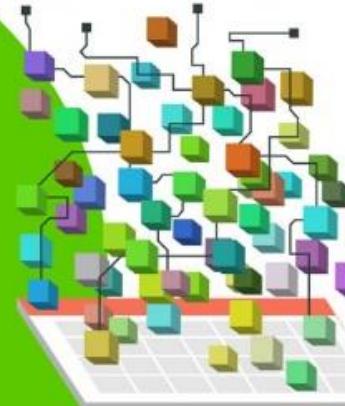
Volume SCALE OF DATA

It's estimated that

2.5 QUINTILLION BYTE

[2.3 TRILLION GIGABYTES]

of data are created each day



Most companies in the U.S. have at least

100 TERABYTES

[100,000 GIGABYTES]

of data stored



IBM Infographic (McKinsey, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS)

4Vs

The New York Stock Exchange captures
captures

1 TB OF TRADE INFORMATION

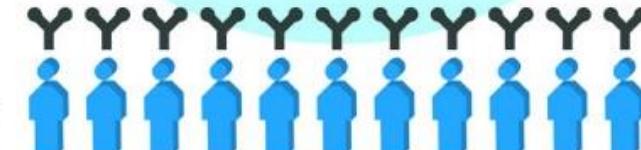
during each trading session



By 2016, it is projected there will be

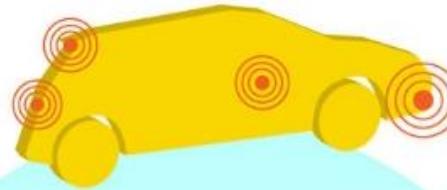
18.9 BILLION NETWORK CONNECTIONS

– almost 2.5 connections per person on earth



Velocity

ANALYSIS OF STREAMING DATA



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



IBM Infographic (McKinsey, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS)

4Vs

As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES

[161 BILLION GIGABYTES]



**30 BILLION
PIECES OF CONTENT**

are shared on Facebook every month



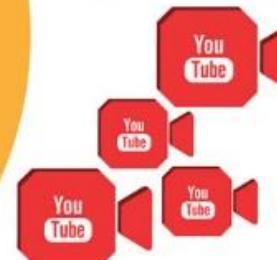
Variety
**DIFFERENT
FORMS OF DATA**

By 2014, it's anticipated there will be

**420 MILLION
WEARABLE, WIRELESS
HEALTH MONITORS**

**4 BILLION+
HOURS OF VIDEO**

are watched on YouTube each month



400 MILLION TWEETS

are sent per day by about 200 million monthly active users

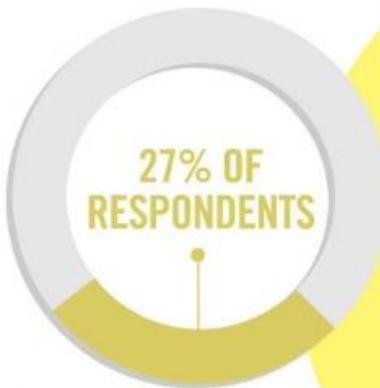


IBM Infographic (McKinsey, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS)

4Vs

1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions

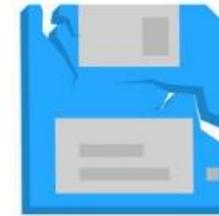


in one survey were unsure of how much of their data was inaccurate

Veracity UNCERTAINTY OF DATA

Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



IBM Infographic (McKinsey, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS)

The Database Explosion

A database system for every use case?



RDBMS

General-purpose
ACID transactions



Wide-Column Store

Long scans over
structured data



Graph Database

Graph algorithms
& queries



Parallel DWH

Aggregations/OLAP for
massive data amounts



Document Store

Deeply nested
data models



In-Memory KV-Store

Counting & statistics



NewSQL

High throughput
relational OLTP



Key-Value Store

Large-scale
session storage

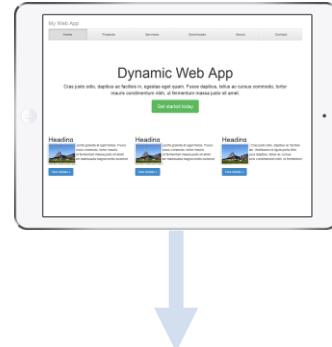


Wide-Column Store

Massive user-
generated content

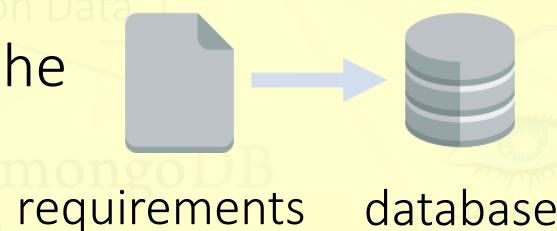
How to choose a database system?

Many Potential Candidates



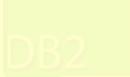
Question in this lecture:

How to approach the



decision problem?

Friend network



Cached data & metrics



redis

elasticsearch.



Apache Hadoop
Amazon Elastic MapReduce

Application Layer

Billing Data

Nested Application Data



Session data

Files



Google Cloud Storage

Recommendation Engine



Search Index

NoSQL Databases

- „NoSQL“ term coined in 2009
- Interpretation: „Not Only SQL“
- Typical properties:
 - Non-relational
 - Open-Source
 - Schema-less (*schema-free*)
 - Optimized for distribution (clusters)
 - Tunable consistency

NoSQL-Databases.org:
Current list has over 150
NoSQL systems



Wide Column Store / Column Families

Hadoop / HBase API: Java / any writer. Protocol: any write call. Query Method: MapReduce, Java / any exec. Replication: HDFS Replication, Write in: Java. Concurrency: 1, Mem: Links 3 Books [1, 2, 3].
Cassandra: massively scalable, partitioned row store. Replication: 1 to 5. Consistency: 1 to 5. Single points of failure, replication support across multiple data centers & cloud availability zones. API: Query Method: CQL and Thrift, replication: peer-to-peer, writes in: Java. Concurrency: 1 to 5. Mem: Links 3 Books [1, 2, 3]. Data compression, MapReduce support, primary/secondary index, security features. Links: Documentation, Planets, Company.

Hypertable API: Thrift[Java, PHP, Perl, Python, Ruby, C/C++], Thrift, Query Method: HQL, native Thrift.

API: Replication: HDFS Replication, Concurrency: MVCC, Consistency Model: Fully consistent. Mem: High performance implementation of Google's Bigtable.

Concurrent updates: 1 to 5.

Accumulo: Accumulo is based on BigTable and is built on top of Hadoop, Zookeeper, and Thrift. It features improvements on the BigTable design in the form of cell-based access control, improved compression, and a service-oriented programming interface that can modify replicated parts at various points in the data management process.

Amazon SimpleDB: Not open source / part of AWS. EC2 with no dependencies by DynamoDB [0].

Cloudata: Google Bigtable clone. Mem: In hac: 1. Article.

Cloudkeeper: Professional Software & Services based on Hadoop.

HPC from [ExaHello](#). Info, article.

Stratosphere: research system, massive parallel & flexible execution, MR generalization and extension (map, reduce).

[OpenComputing, Oracle, KDD]

Document Store

MongoDB API: BSON, Protocol: C, Query Method: dynamic object-based language & MapReduce, Replication: Master-Slave & Auto-Sharding, Write in: C++, Concurrency: Update in Place, Mem: Threadsafe, GridFS, Freescale, Commercial, License: MIT, Links: [MongoDB](#), [MongoDB](#).

Elasticsearch API: REST and many languages.

Protocol: REST, Query Method: via JSON, Replication: Sharding automatic and configurable, Write in: Java, Concurrency: 1 to 5, Mem: 1000, 10000, 100000, 1000000, Consistency: 1 to 5.

Couchbase Server API: Memcached API+protocol binary and ASCII, most languages Protocol: Memcached REST Interface for cluster config + memcached API, Query Method: Memcached, Replication: Peer to Peer, fully consistent, Mem: Transparency: topology changes during operation, protocols: memcached-compatible, cache: 1000000, 10000000, 100000000, supported versions available, Links: [Couchbase](#).

CouchDB API: JSON, Protocol: REST, Query Method: MapReduce, Functions, Replication: Master-Master, Write in: Erlang, Concurrency: MVCC, Mem: 100000000, Links: [3 CouchDB books](#), [Couch Lounge](#) (partitioning/ clustering), [100000000](#).

RethinkDB API: protobuf-based, Query Method: unified chainable query language (incl. JOINs, subqueries, ORDER BY, GROUP BY, LIMIT, etc.), Replication: Sync and Async Master Slave with per-table acknowledgements, Sharding: pulsed range-based, Write in: C++, Concurrency: MVCC, Mem: 100000000, Links: [RethinkDB](#).

RavenDB: Not solution. Provides HTTP/JSON access, LINQ queries & Sharing supported. [RavenDB](#).

MarkLogic Server (commercial) API: XML, PHP, Java, .NET, Protocol: HTTP, REST, native, TCP/IP. Query Method: Full Text, XML, range, and Xpath, Write in: C++, Concurrency: ACID-compliant, transactional, Shared-nothing cluster, Mem: Per-byte-scalable, document stores and full text search engine, Information ranking, Replication: Cloudable.

Cloudant Server (commercial) API: XML, PHP, Java, .NET, Protocol: HTTP, REST, native, TCP/IP. Query Method: Full Text, XML, range, and Xpath, Write in: C++, Concurrency: ACID-compliant, transactional, Shared-nothing cluster, Mem: Per-byte-scalable document stores and full text search engine, Information ranking, Replication: Cloudable.

ThruDB: please note: provides more facets! Uses Apache Thrift to map multiple database databases to Redis, MySQL, SS.

Torrastore API: Java & http, Protocol: http, Language: Java, Querying: Range queries, Predicates: Application-specific, Replication: Sync with consistent hashing, Consistency: Partition strict consistency, Mem: Based on TorraStore.

JavaDB: Lightweight open source document database written in Java for high performance, runs in memory, supports Android, Java, JSON, Java, XML, CSV, and ODBC.

QueryDB API: Java, Protocol: JDBC API.

Concurrency: Atomic document writes (locks), eventually consistent indexes.

RaptorDB: SQL based Document store database with columnar storage, distributed, replicated, automatic hybrid bloom indexing and LRU query filters.

SizedB: Document Store on top of SQL Server.

SDB: For small online databases. PHP / JSON interface, implemented in PHP.

cloudb **spanDB API**: BSON, Protocol: C++, Query Method:

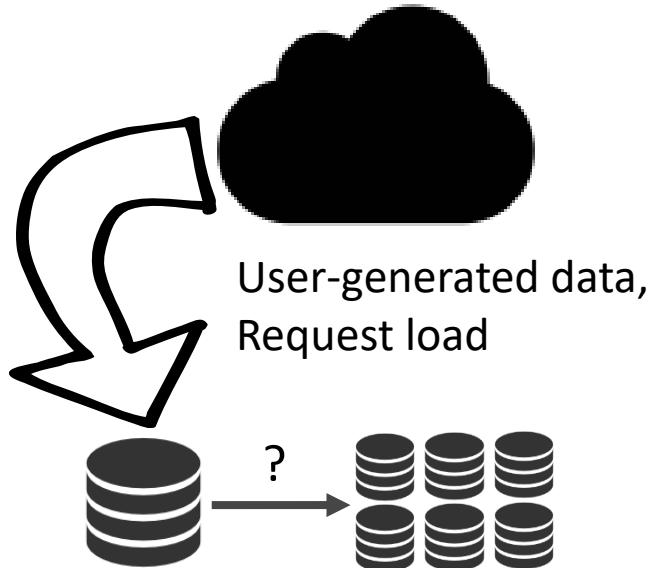
dynamic & predicates and MapReduce, Drivers: Java, C++, Python, Go, Node.js, .NET, C, C++, C#.

Protocol: BSON, API: REST, API: XMLHttpRequest.

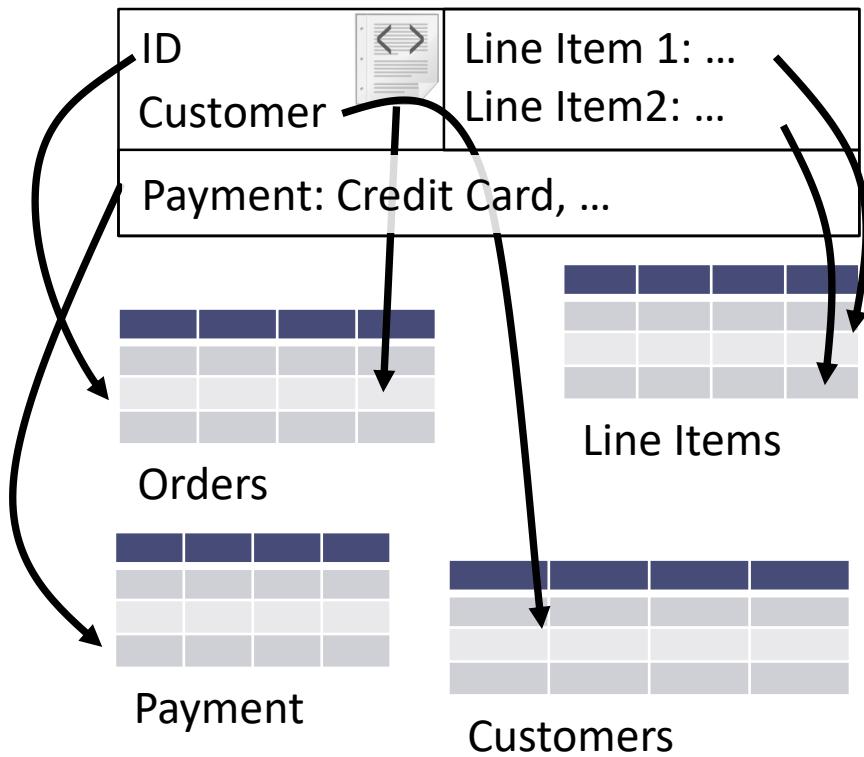
NoSQL Databases

- ▶ Two main motivations:

Scalability

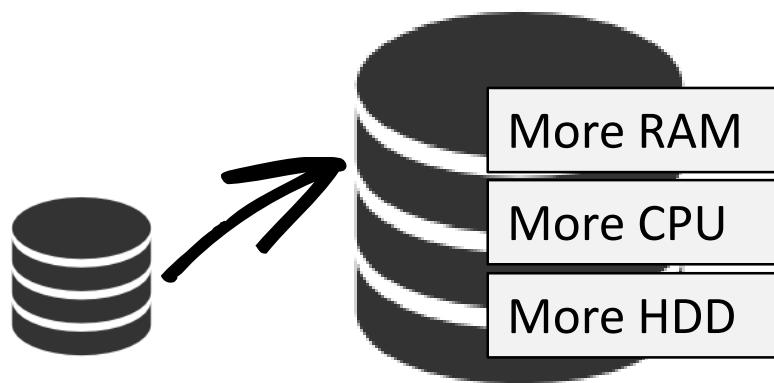


Impedance Mismatch

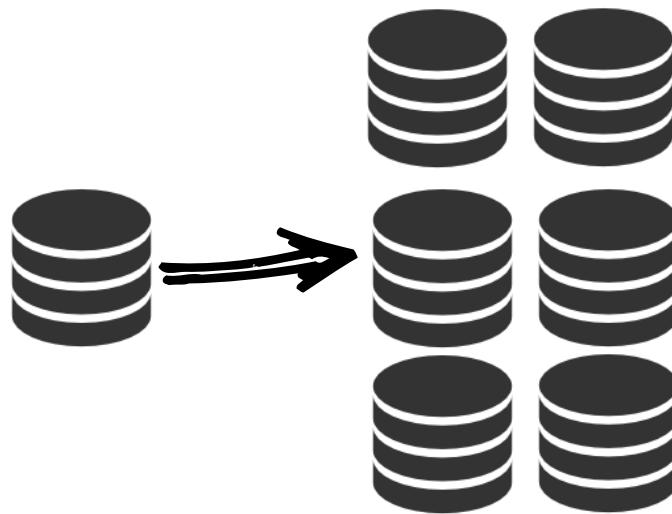


Scale-up vs Scale-out

Scale-Up (*vertical scaling*):

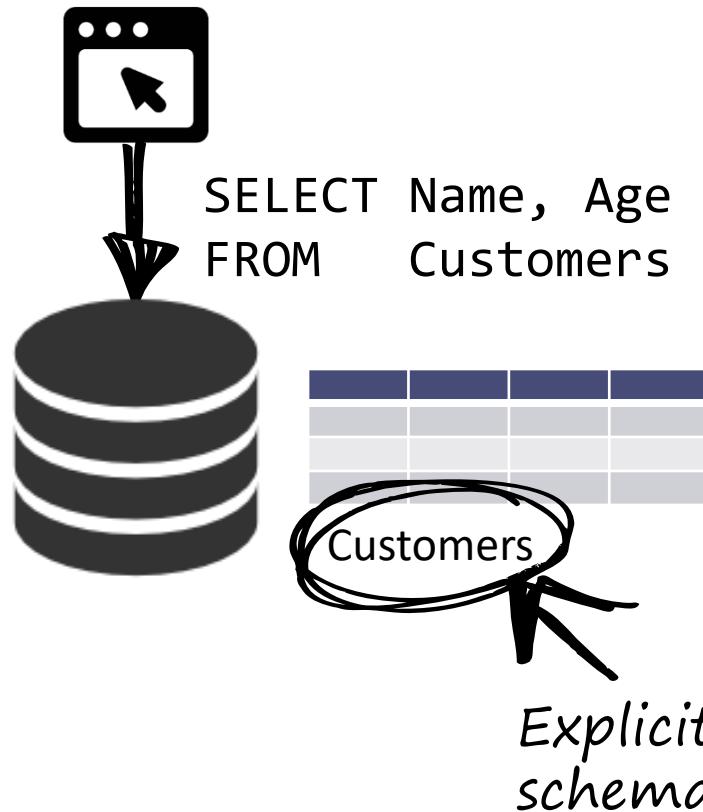


Scale-Out (*horizontal scaling*):

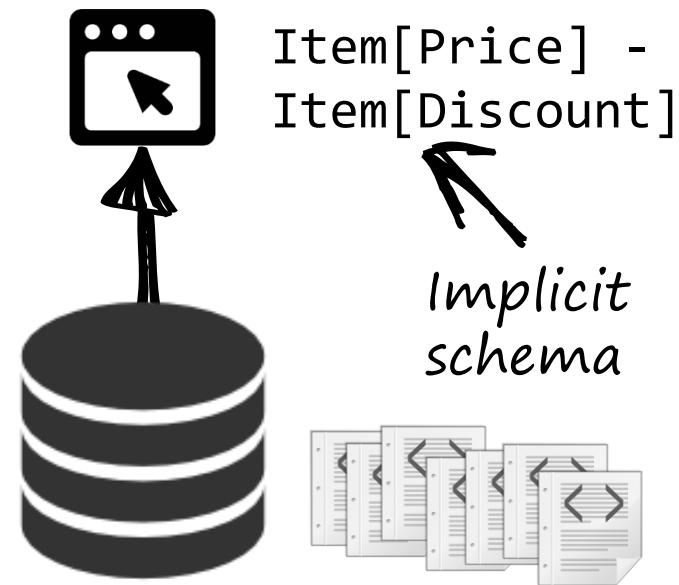


Schema-free Data Modeling

RDBMS:



NoSQL DB:



NoSQL Paradigm Shift

Open Source & Commodity Hardware



Commercial DBMS



Open-Source DBMS

Specialized DB hardware
(Oracle Exadata, etc.)



Commodity hardware

Highly available network
(Infiniband, Fabric Path, etc.)



Commodity network
(Ethernet, etc.)

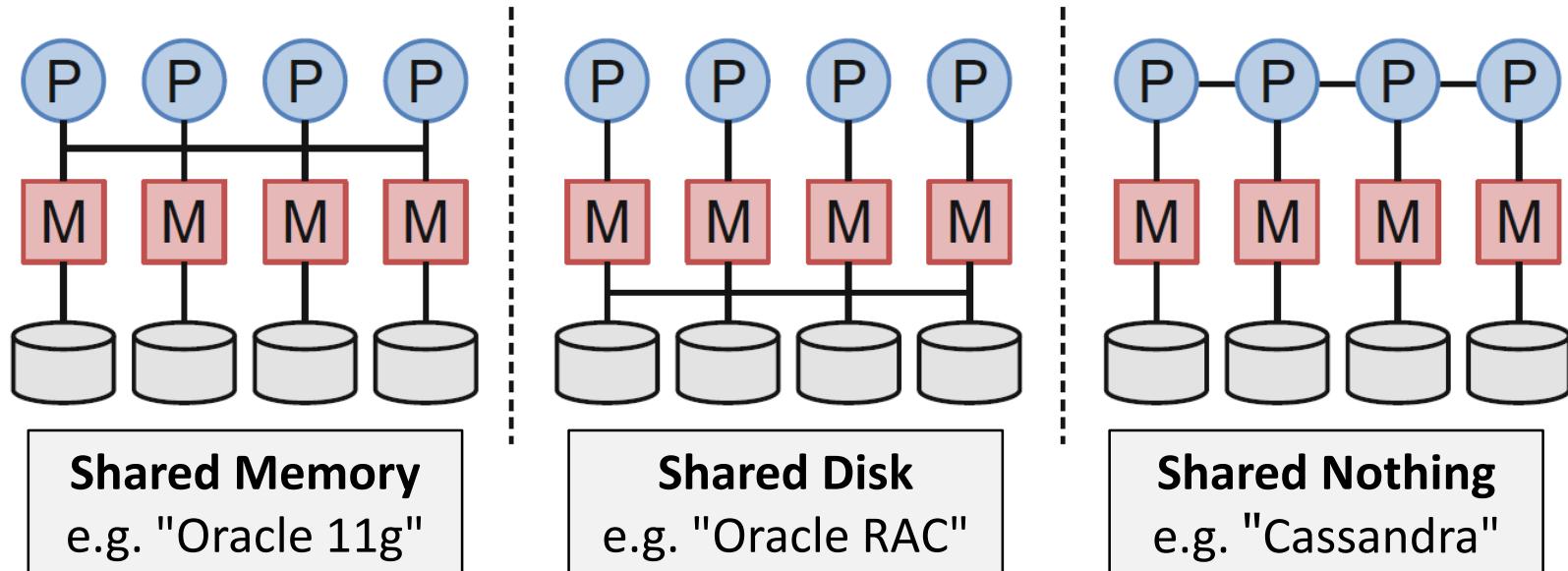
Highly Available Storage (SAN,
RAID, etc.)



Commodity drives (standard
HDDs, JBOD)

Paradigm Shift

- Shift towards **distributed computing architectures**



NoSQL and Big Data Architecture

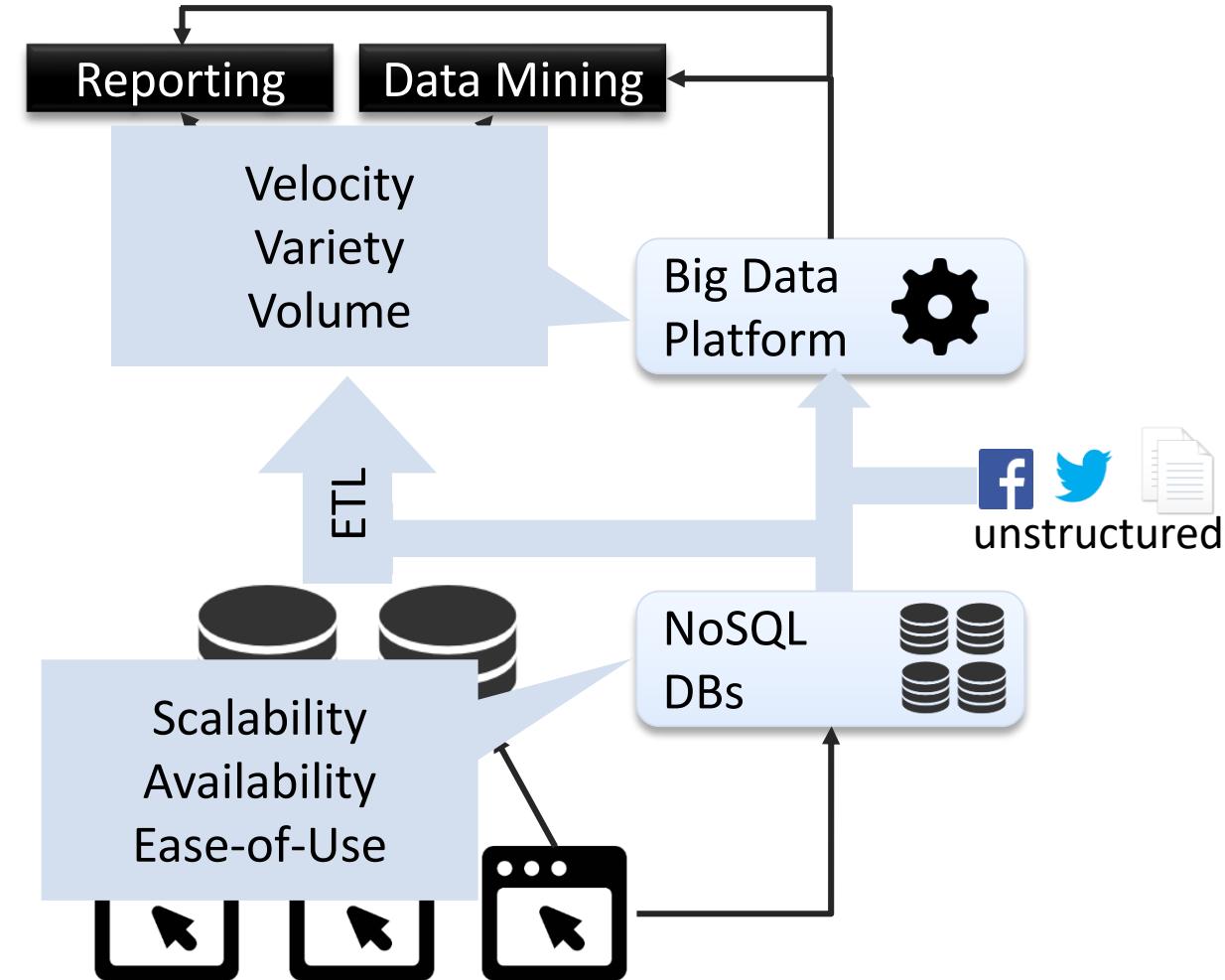
Data Management

Analytics

Data Warehouse

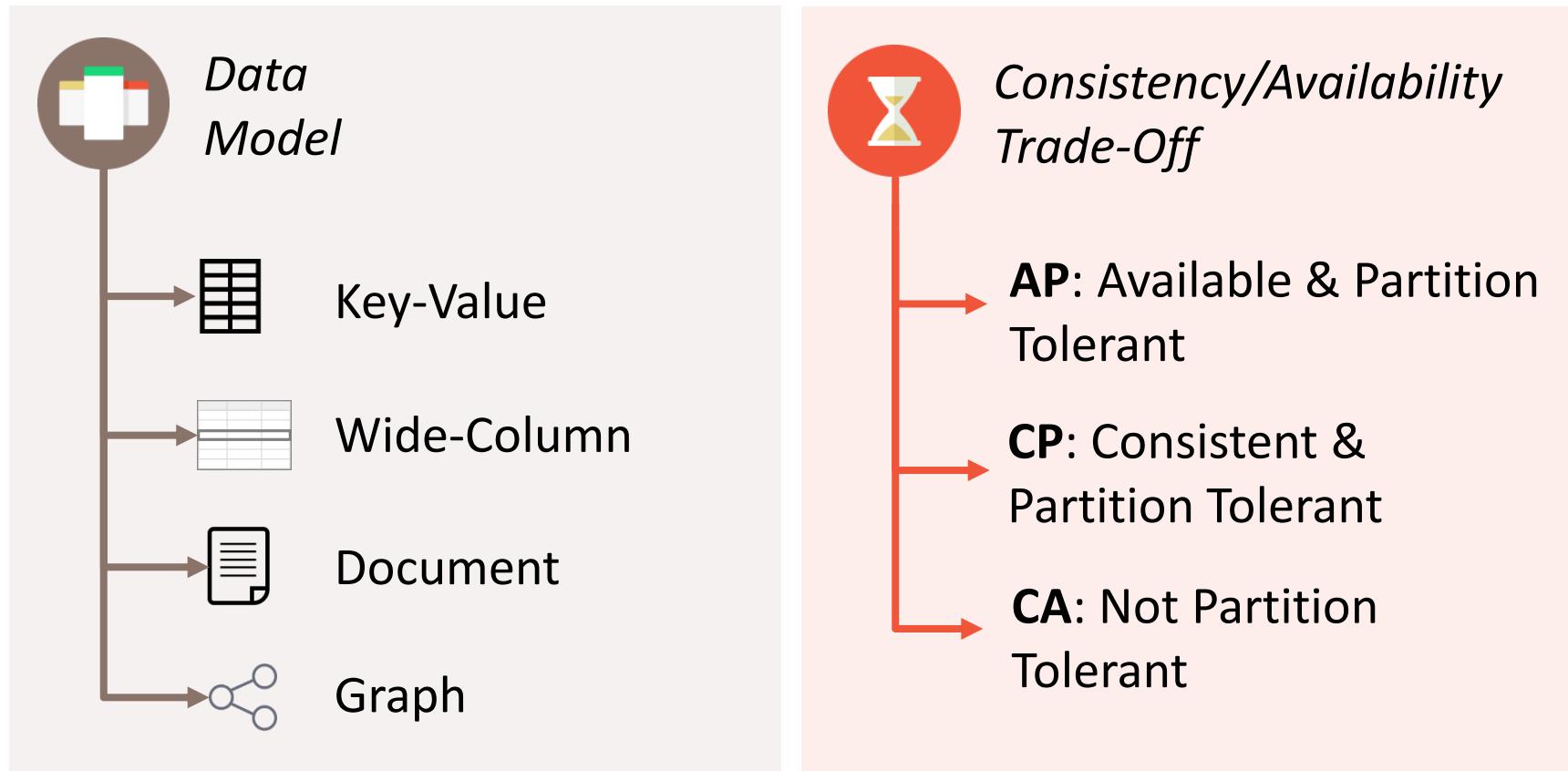
Operational Databases

Application



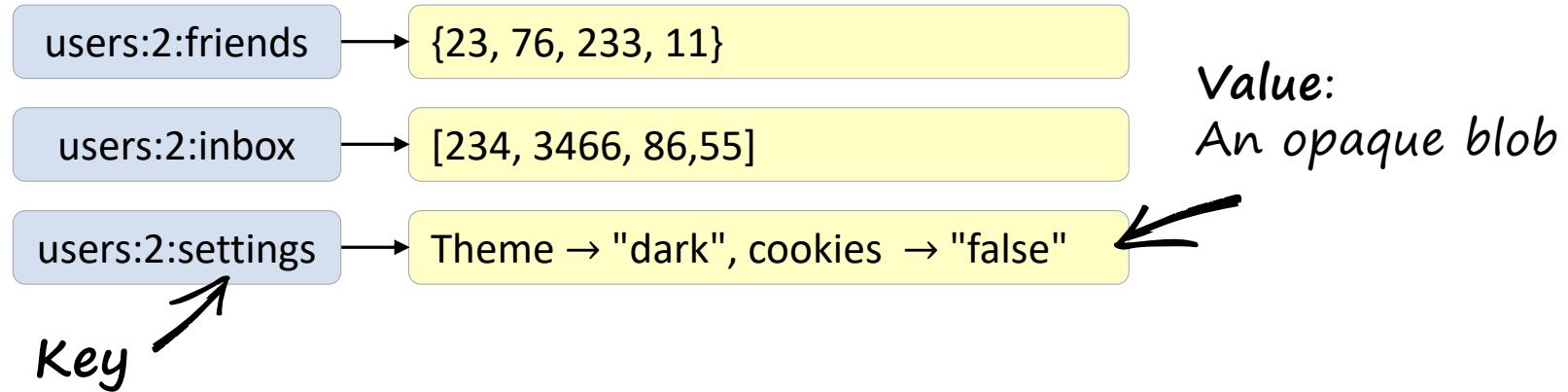
NoSQL System Classification

- ▶ Two common criteria:



Key-Value Stores

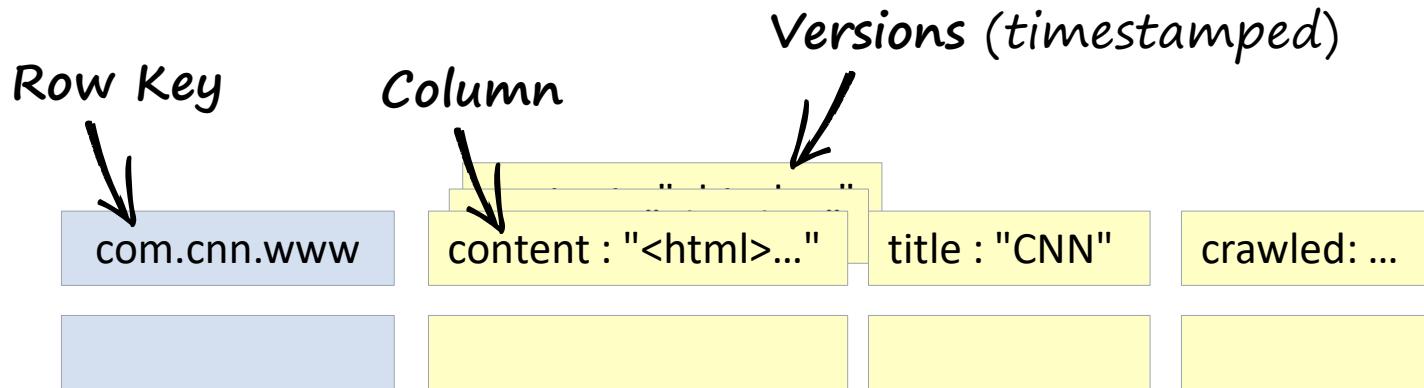
- ▶ Data model: (key) -> value
- ▶ Interface: CRUD (Create, Read, Update, Delete)



- ▶ Examples: Amazon Dynamo (AP), Riak (AP), Redis (CP)

Wide-Column Stores

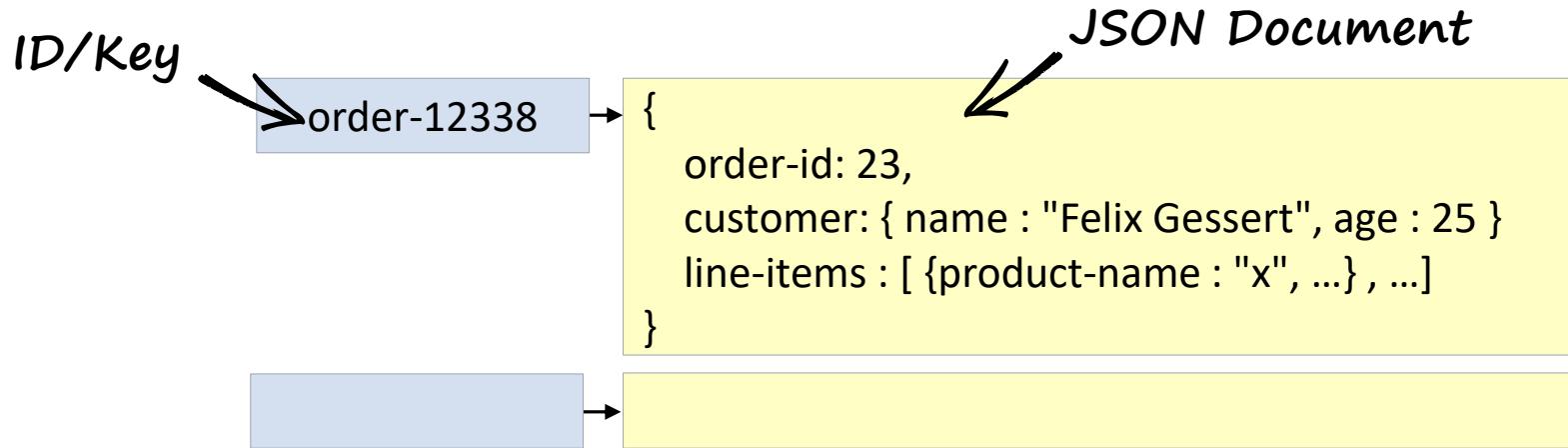
- ▶ Data model: (rowkey, column, timestamp) -> value
- ▶ Interface: CRUD, Scan



- ▶ Examples: Cassandra (AP), Google BigTable (CP), HBase (CP)

Document Stores

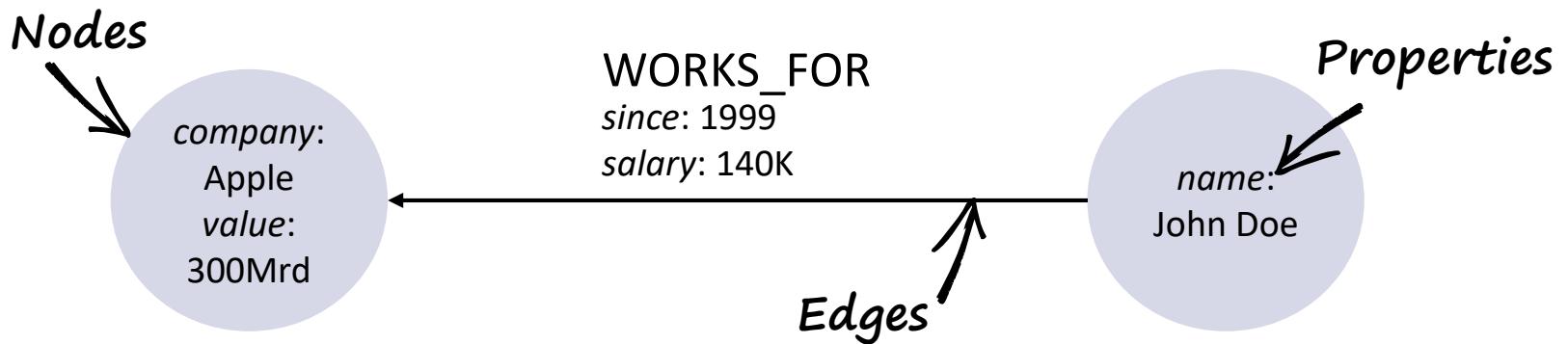
- ▶ Data model: (collection, key) -> document
- ▶ Interface: CRUD, Querys, Map-Reduce



- ▶ Examples: CouchDB (AP), RethinkDB (CP), MongoDB (CP)

Graph Databases

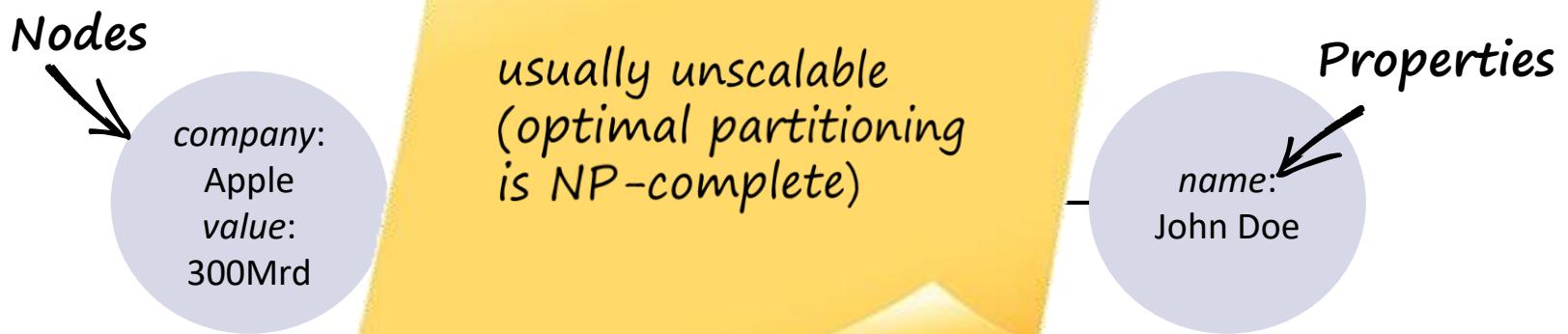
- ▶ Data model: $G = (V, E)$: Graph-Property Modell
- ▶ Interface: Traversal algorithms, querys, transactions



- ▶ Examples: Neo4j (CA), InfiniteGraph (CA), OrientDB (CA)

Graph Databases

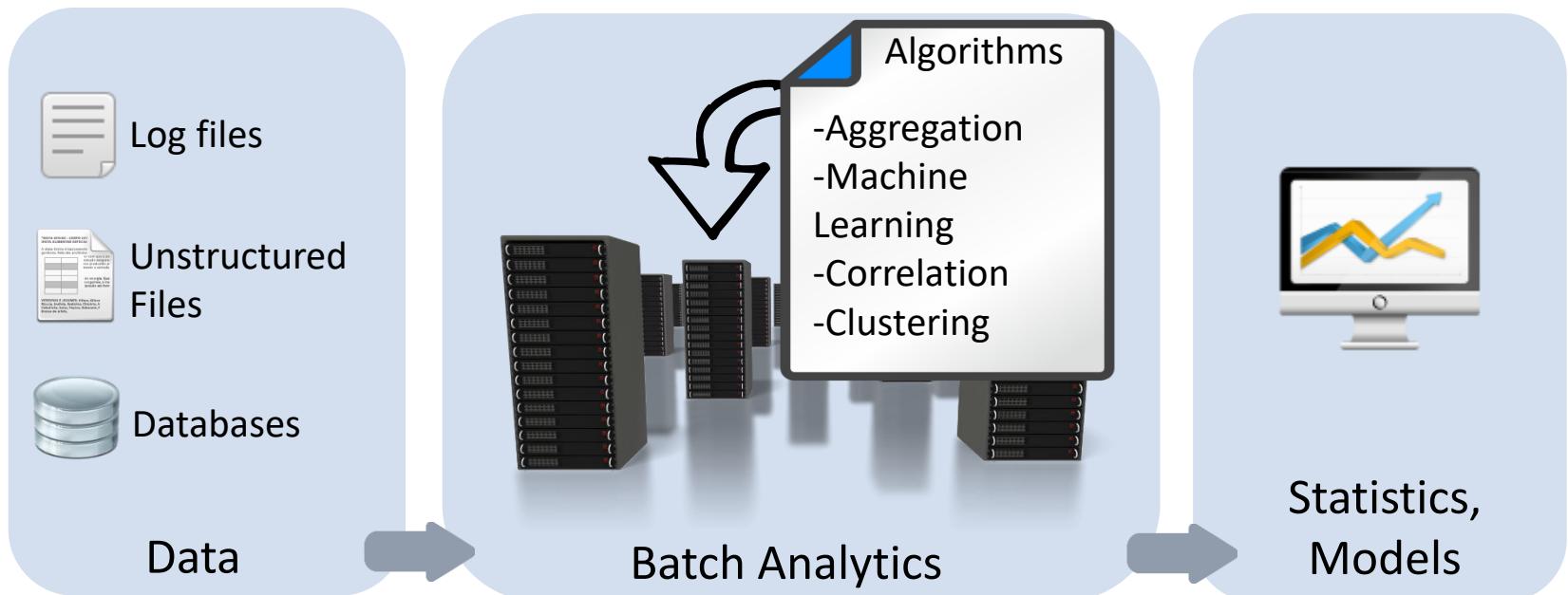
- ▶ Data model: $G = (V, E)$: Graph-Property Modell
- ▶ Interface: Traversal, Cypher, Gremlin, transactions



- ▶ Examples: Neo4j (CA), OrientGraph (CA), OrientDB (CA)

Big Data Batch Processing

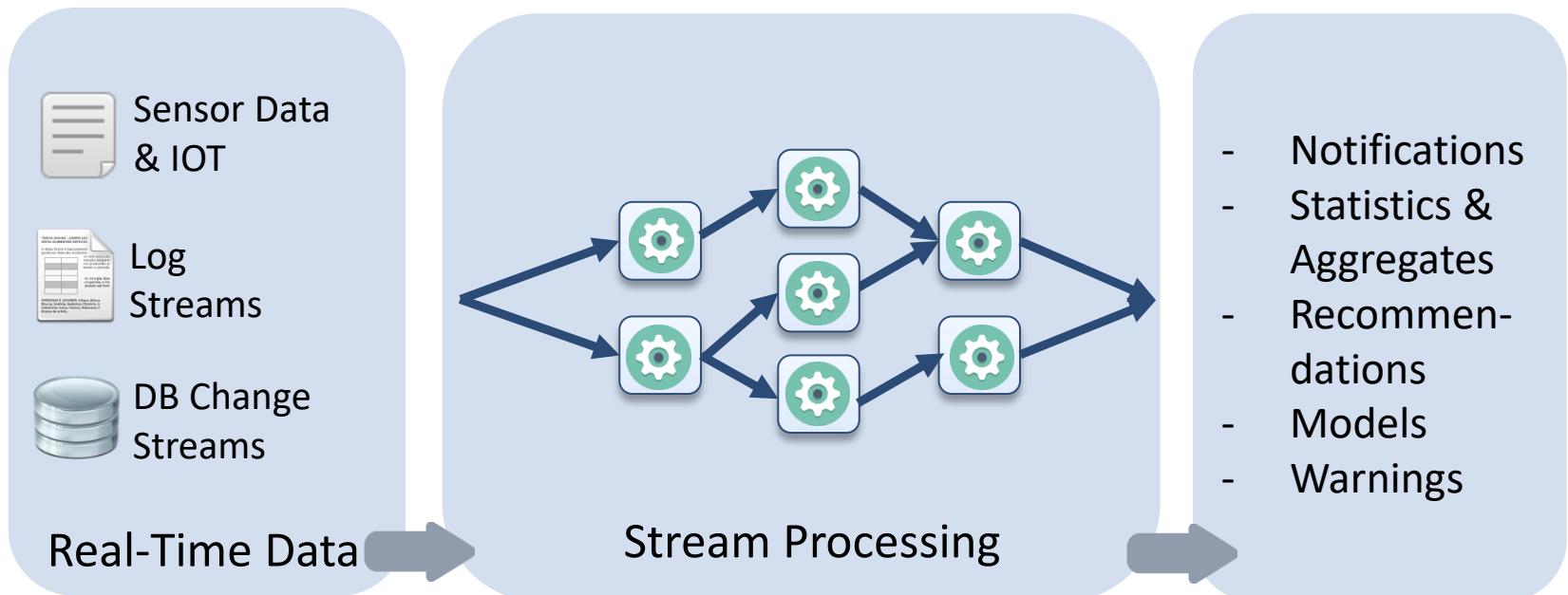
- ▶ Data model: arbitrary (frequently unstructured)
- ▶ Examples: Hadoop, Spark, Flink, DryadLink, Pregel



Big Data Stream Processing

Covered in Depth in the Last Part

- ▶ Data model: arbitrary
- ▶ Examples: Storm, Samza, Flink, Spark Streaming



Soft NoSQL Systems

Not Covered Here



Search Platforms (Full Text Search):

- No persistence and consistency guarantees for OLTP
- *Examples:* ElasticSearch (AP), Solr (AP)



Object-Oriented Databases:

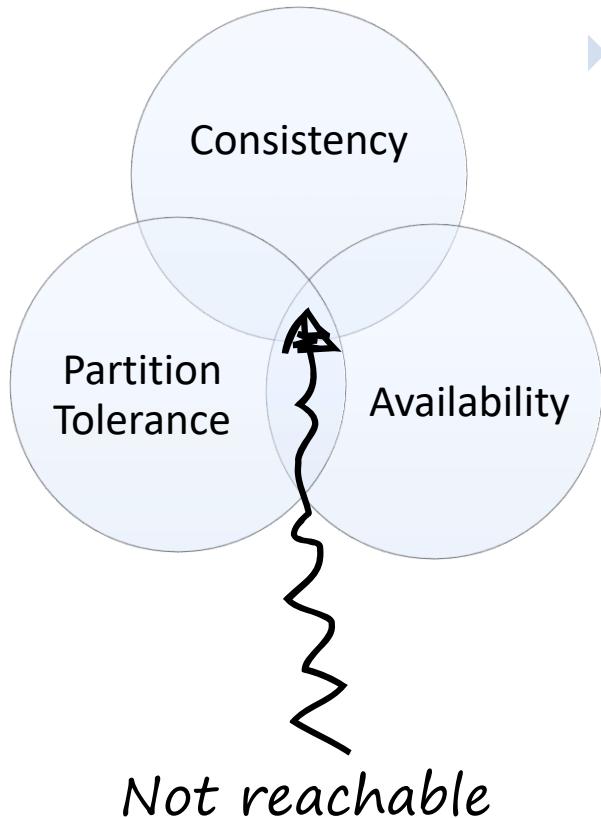
- Strong coupling of programming language and DB
- *Examples:* Versant (CA), db4o (CA), Objectivity (CA)



XML-Databases, RDF-Stores:

- Not scalable, data models not widely used in industry
- *Examples:* MarkLogic (CA), AllegroGraph (CA)

CAP-Theorem



- ▶ Classifies distributed databases
- ▶ Only 2 out of 3 properties are achievable at a time:
 - **Consistency:** all clients have the same view on the data
 - **Availability:** every request to a non-failed node must result in correct response
 - **Partition tolerance:** the system has to continue working, even under arbitrary network partitions



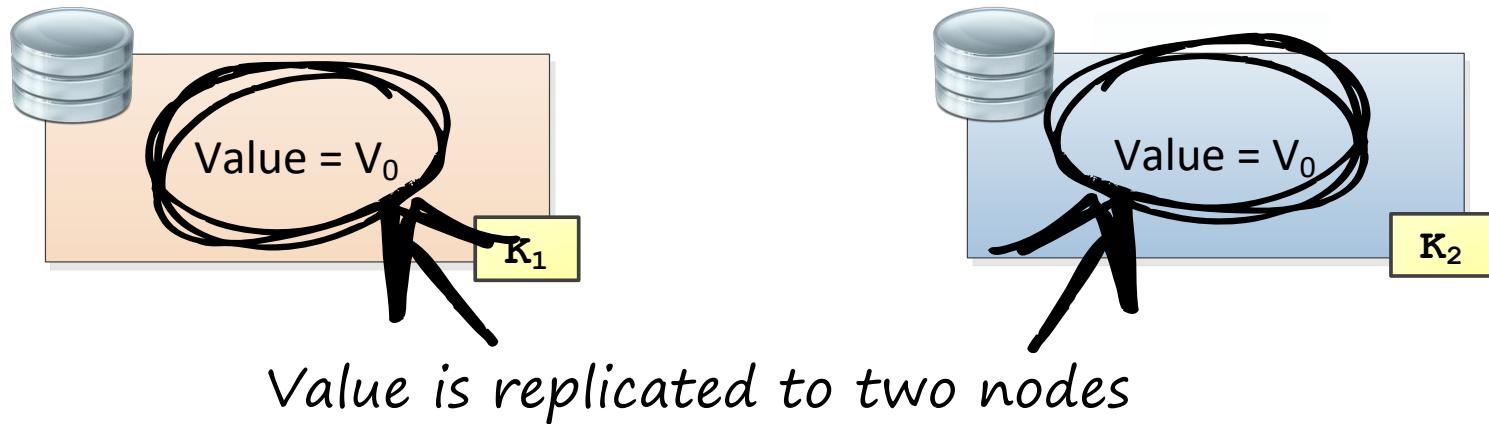
Eric Brewer, ACM-PODC Keynote, Juli 2000



Gilbert, Lynch: Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services, SigAct News 2002

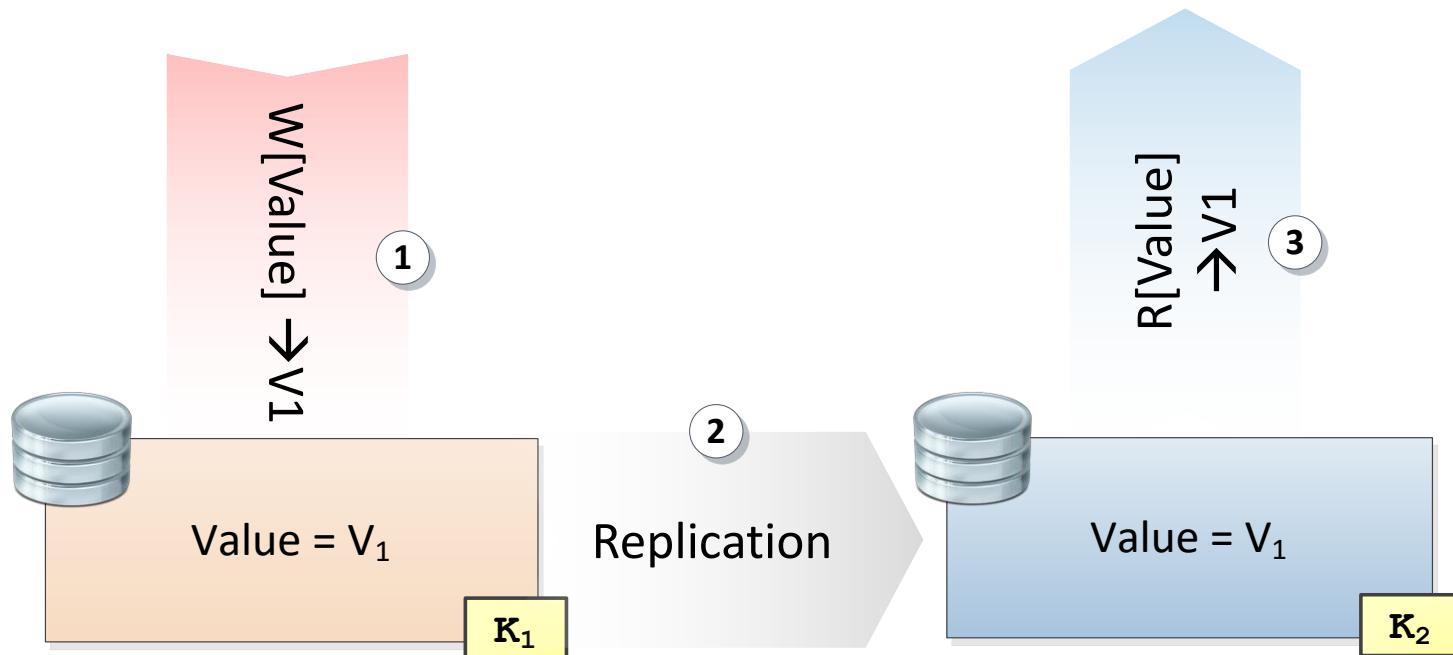
CAP-Theorem: simplified proof

- ▶ Intuition for the impossibility of simultaneously achieving C, A and P at once:



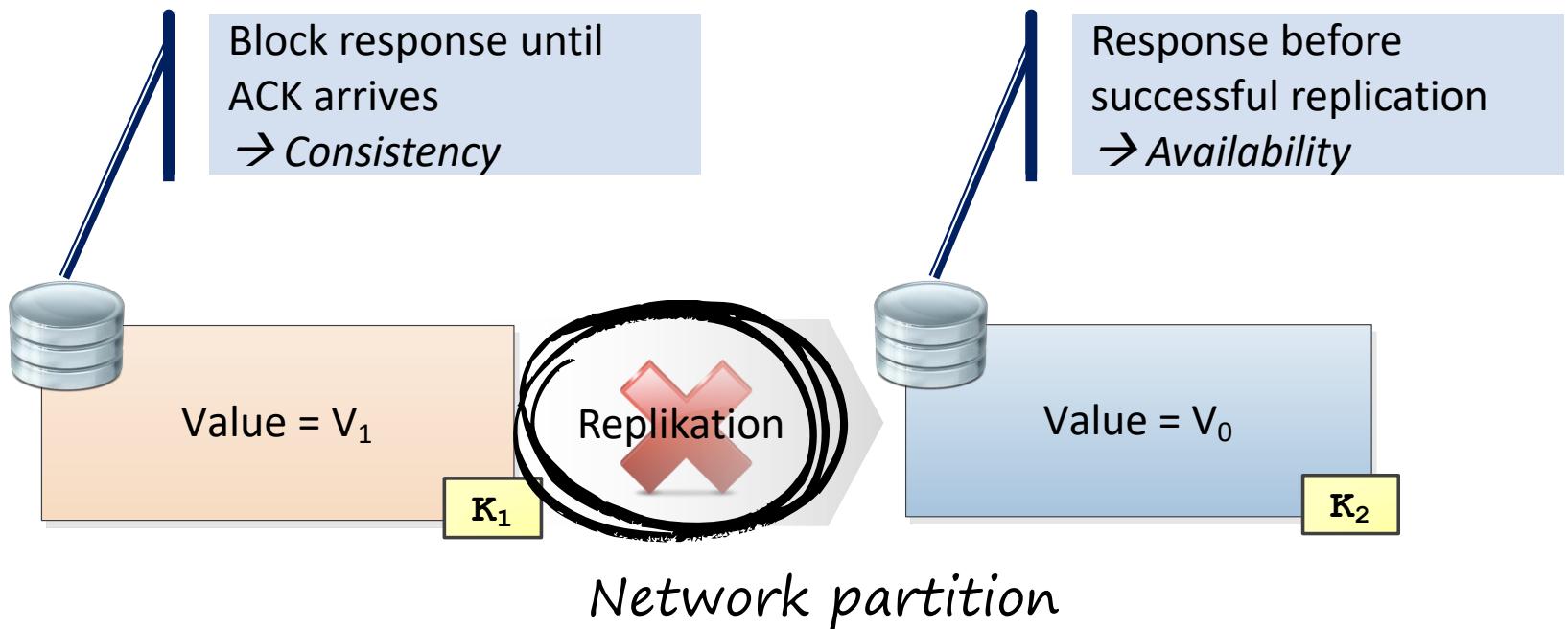
CAP-Theorem: simplified proof

- ▶ Failure-free reading and writing:

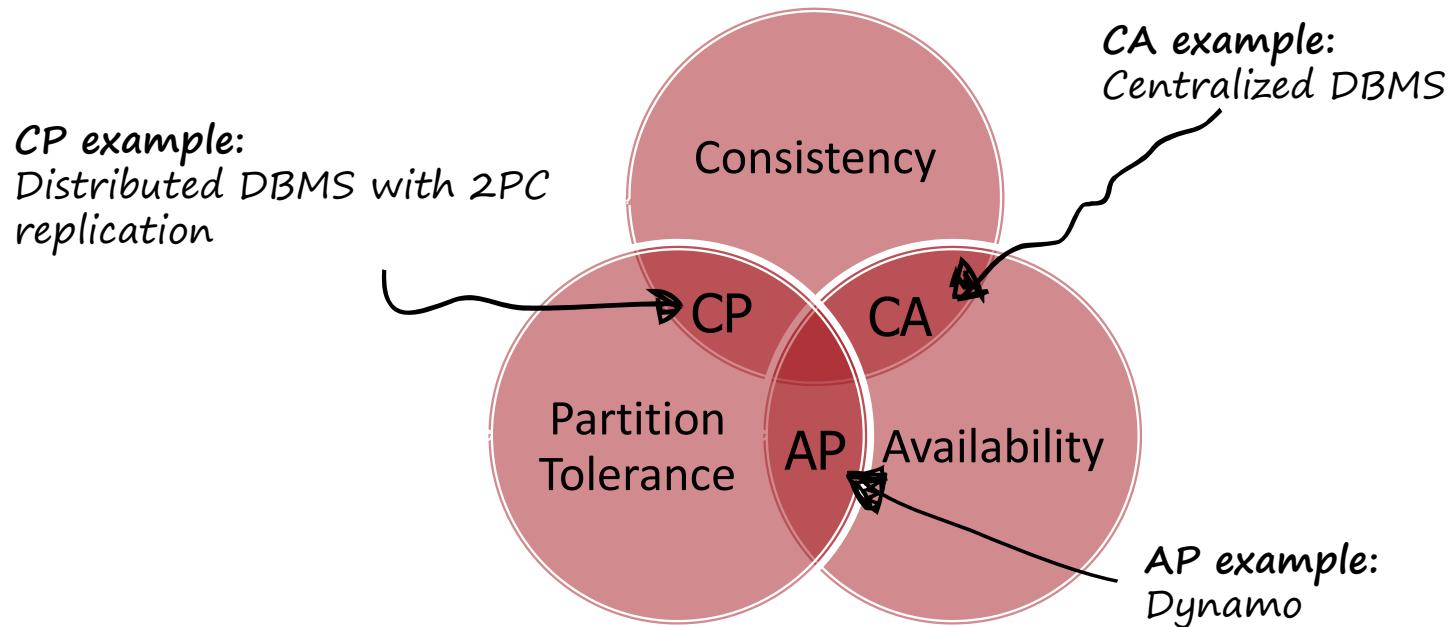


CAP-Theorem: simplified proof

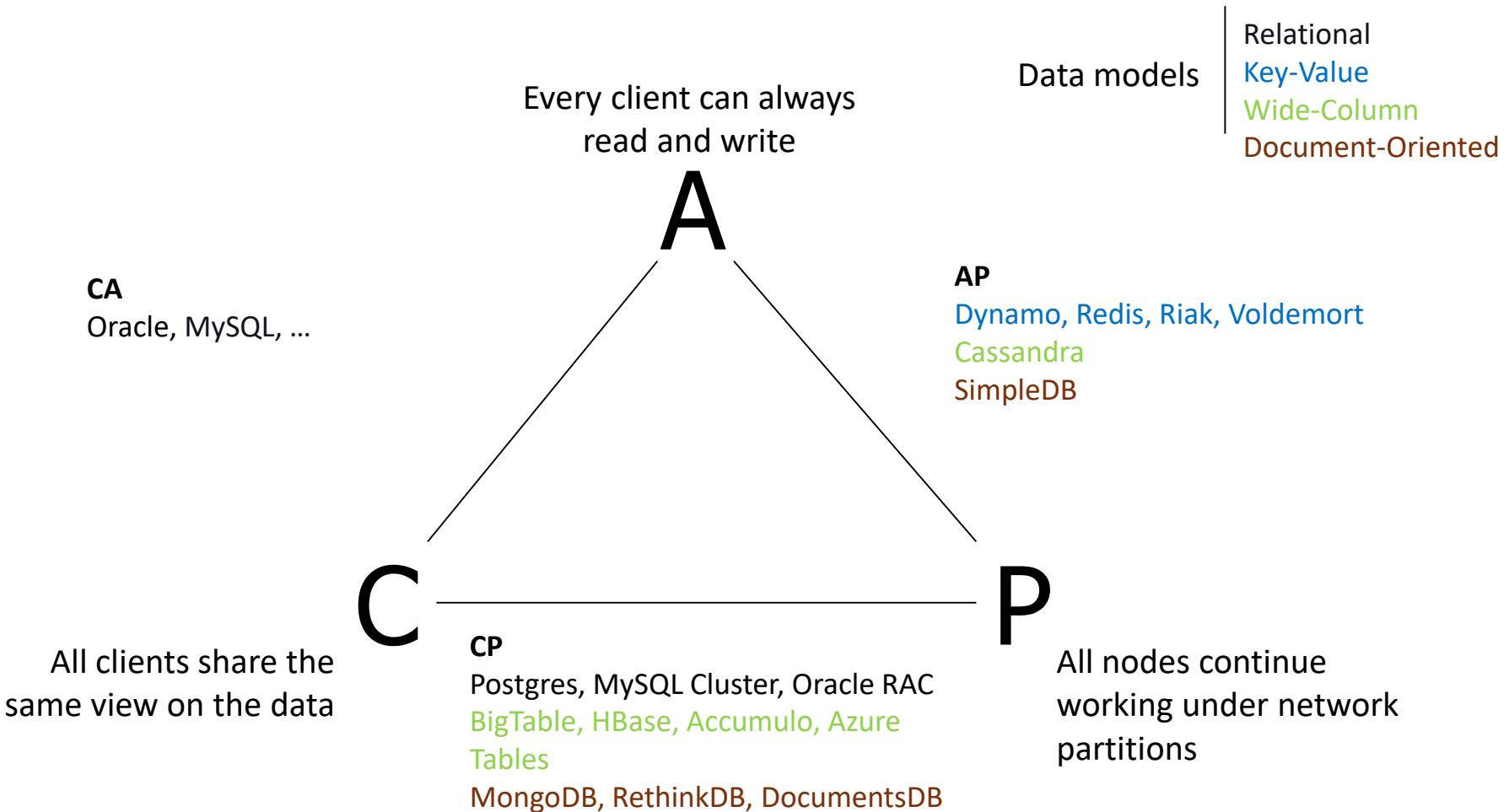
- ▶ Problem: when a network partition occurs, either consistency or availability have to be given up



CAP-Theorem: Wrap-up



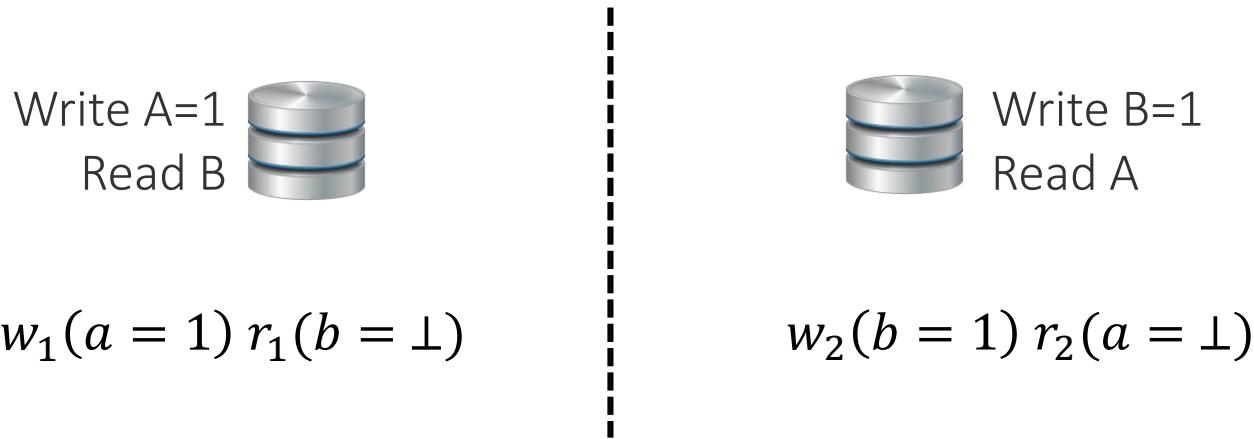
NoSQL Triangle



Serializability

Not Highly Available Either

Global serializability and availability are incompatible:



- ▶ Some weaker isolation levels allow high availability:
 - RAMP Transactions (P. Bailis, A. Fekete, A. Ghodsi, J. M. Hellerstein, und I. Stoica, „Scalable Atomic Visibility with RAMP Transactions“, SIGMOD 2014)



S. Davidson, H. Garcia-Molina, and D. Skeen. Consistency in partitioned networks. ACM CSUR, 17(3):341–370, 1985.

Where CAP fits in

Negative Results in Distributed Computing

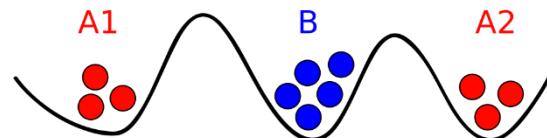
***Asynchronous Network,
Unreliable Channel***

Atomic Storage

Impossible:
CAP Theorem

Consensus

Impossible:
2 Generals Problem



***Asynchronous Network,
Reliable Channel***

Atomic Storage

Possible:
Attiya, Bar-Noy, Dolev (ABD) Algorithm

Consensus

Impossible:
Fisher Lynch Patterson (FLP) Theorem



Lynch, Nancy A. *Distributed algorithms*.
Morgan Kaufmann, 1996.

Consensus Algorithms

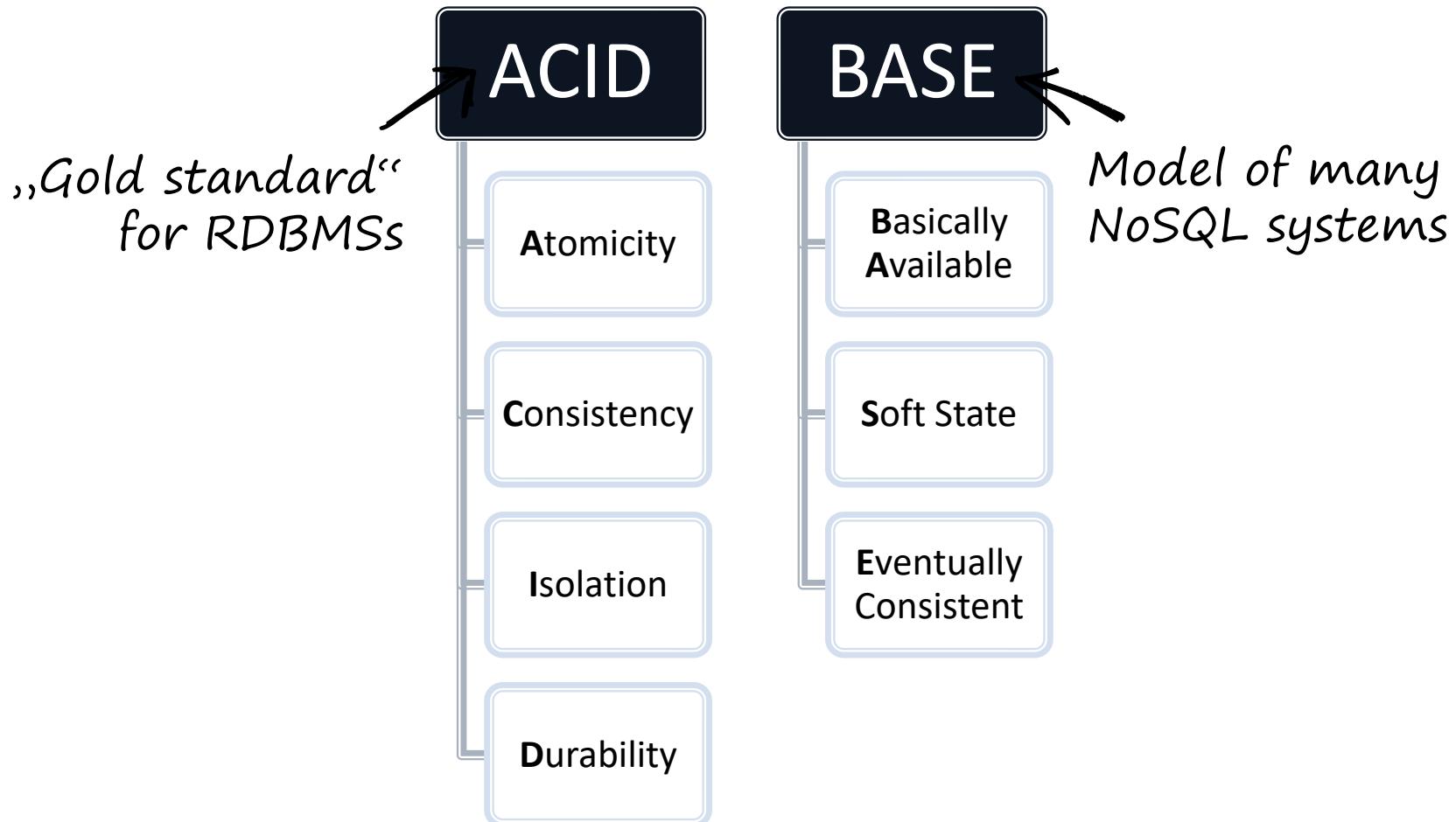
Even Harder Than CAP

Safety
Properties

- ▶ Consensus:
 - *Agreement*: No two processes can commit different decisions
 - *Validity (Non-triviality)*: If all initial values are same, nodes must commit that value
 - *Termination*: Nodes commit eventually
- ▶ No algorithm *guarantees* termination (FLP)
- ▶ Algorithms:
 - **Paxos** (e.g. Google Chubby, Spanner, Megastore, Aerospike, Cassandra Lightweight Transactions)
 - **Raft** (e.g. RethinkDB, etcd service)
 - Zookeeper Atomic Broadcast (**ZAB**)



ACID vs BASE



Weaker guarantees in a database?!

Default Isolation Levels in RDBMSs

Database	Default Isolation	Maximum Isolation
Actian Ingres 10.0/10S	S	S
Aerospike		RC
Clustrix CLX 4100		?
Greenplum 4.1		S
IBM DB2 10 for z/OS		S
IBM Informix 11.50		RR
MySQL 5.6		S
MemSQL 1b	RC	RC
MS SQL Server 2012	RC	CR
NuoDB	CR	SI
Oracle 11g	RC	SI
Oracle Berkeley DB	S	S
Postgres 9.2.2	RC	S
SAP HANA	RC	SI
ScaleDB 1.02	RC	RC
VoltDB	S	S

Depends
Theorem:

Trade-offs are central to the choice
of a database system.

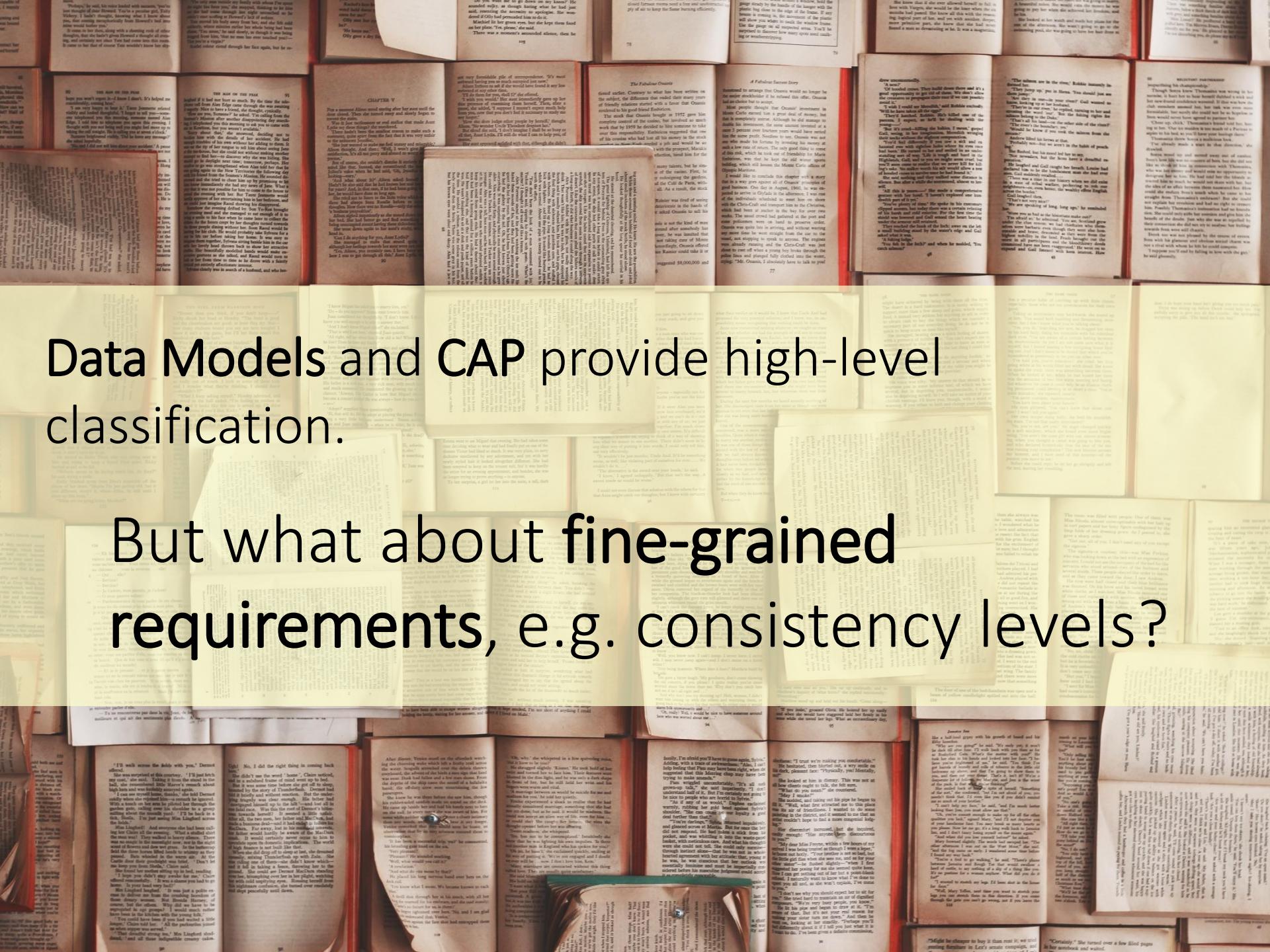
RC: read committed, RR: repeatable read, S: serializability,
SI: snapshot isolation, CS: cursor stability, CR: consistent read



Bailis, Peter, et al. "Highly available transactions: Virtues and limitations." *Proceedings of the VLDB Endowment* 7.3 (2013): 181-192.

Data Models and CAP provide high-level classification.

But what about fine-grained requirements, e.g. consistency levels?



Outline



NoSQL Foundations and Motivation



The NoSQL Toolbox:
Common Techniques

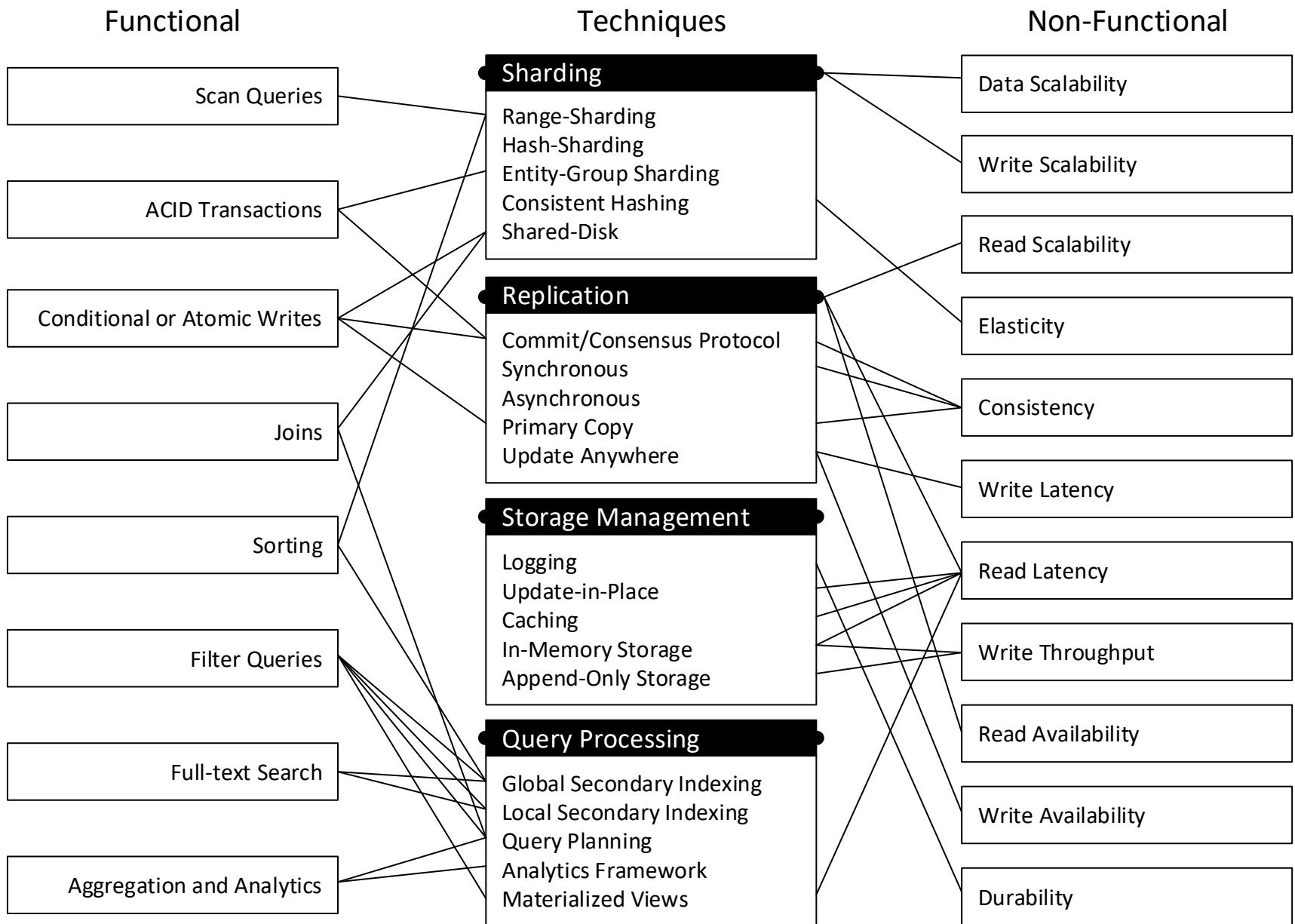


NoSQL Systems &
Decision Guidance



Scalable Real-Time
Databases and Processing

- Techniques for Functional and Non-functional Requirements
 - Sharding
 - Replication
 - Storage Management
 - Query Processing



Functional Requirements from the application



Aggregation and Analytics

Scan Queries

ACID Transactions

Conditional Writes

Joins

String

Series

Search

enable

Techniques

Sharding

Range-Sharding

Hash-Sharding

Entity-Group Sharding

Consistent Hashing

Shared-Disk

Replication

Committed Consistency

Synchronous

Asynchronous

Primary Copy

Updatability

Storage Management

Logging

Update-in-Place

Caching

In-Mem

Append

Query Processing

Global Secondary Indexing

Local Secondary Indexing

Query Planning

Analytics Framework

Materialized Views

enable

Non-Functional

Data Scalability

Write Scalability

Read Scalability

Operational Requirements

Consistency

Write Latency

Read Latency

Write

Read

Write Availability

Durability



NoSQL Database Systems: A Survey and Decision Guidance

Felix Gessert, Wolfram Wingerath, Steffen Friedrich, and Norbert Ritter

Universität Hamburg, Germany

{gessert, wingerath, friedrich, ritter}@informatik.uni-hamburg.de

Abstract. Today, data is generated and consumed at unprecedented scale. This has lead to novel approaches for scalable data management subsumed under the term “NoSQL” database systems to handle the ever-increasing data volume and request loads. However, the heterogeneity and diversity of the numerous existing systems impede the well-informed selection of a data store appropriate for a given application context. Therefore, this article gives a top-down overview of the field: Instead of contrasting the implementation specifics of individual representatives, we propose a comparative classification model that relates functional and non-functional requirements to techniques and algorithms employed in NoSQL databases. This NoSQL Toolbox allows us to derive a simple decision tree to help practitioners and researchers filter potential system candidates based on central application requirements.

1 Introduction

Traditional relational database management systems (RDBMSs) provide powerful mechanisms to store and query structured data under strong consistency and transaction guarantees and have reached an unmatched level of reliability, stability and support through decades of development. In recent years, however, the amount of useful data in some application areas has become so vast that it cannot be stored or processed by traditional database solutions. User-generated content in social networks or data retrieved from large sensor networks are only two examples of this phenomenon commonly referred to as **Big Data** [35]. A class of novel data storage systems able to cope with Big Data are subsumed under the term **NoSQL databases**, many of which offer horizontal scalability and higher availability than relational databases by sacrificing querying capabilities and consistency guarantees. These trade-offs are pivotal for service-oriented computing and as-a-service models, since any stateful service can only be as scalable and fault-tolerant as its underlying data store.

There are dozens of NoSQL database systems and it is hard to keep track of where they excel, where they fail or even where they differ, as implementation details change quickly and feature sets evolve over time. In this article, we therefore aim to provide an overview of the NoSQL landscape by discussing employed concepts rather than system specificities and explore the requirements typically posed to NoSQL database systems, the techniques used to fulfil these requirements and the trade-offs that have to be made in the process. Our focus lies on key-value, document and wide-column stores, since these NoSQL categories

NoSQL Databases: a Survey and Decision Guidance



Felix Gessert

Aug 15, 2016 · 26 min read

Together with our colleagues at the University of Hamburg, we—that is [Felix Gessert](#), [Wolfram Wingerath](#), [Steffen Friedrich](#) and [Norbert Ritter](#)—presented an overview over the NoSQL landscape at [SummerSOC'16](#) last month. Here is the written gist. We give our best to convey the condensed NoSQL knowledge we gathered building [Baqend](#).



As a blog post: blog.baqend.com

Functional

Techniques

Non-Functional

Scan Queries

ACID Transactions

Conditional or Atomic Writes

Joins

Sorting

Sharding

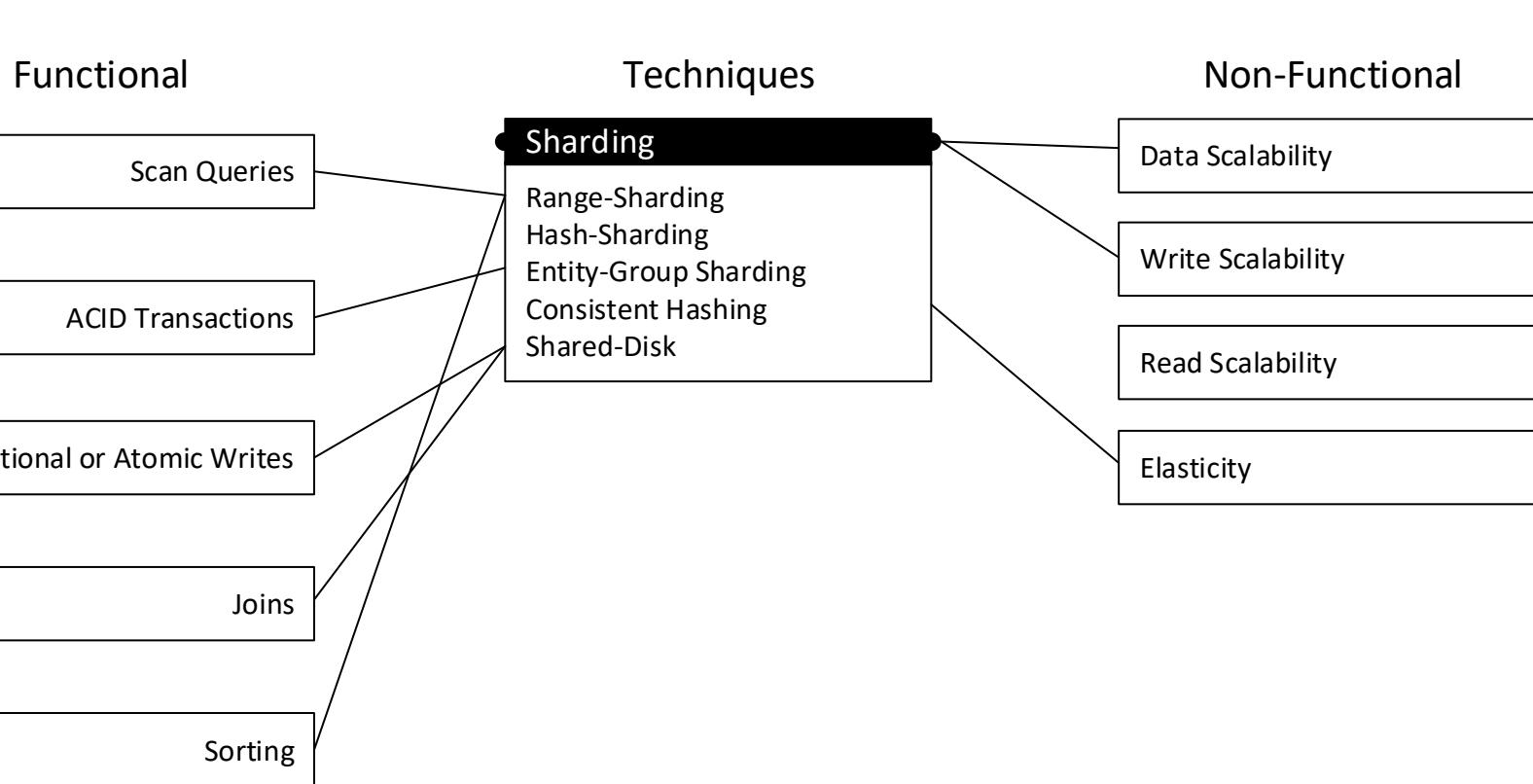
- Range-Sharding
- Hash-Sharding
- Entity-Group Sharding
- Consistent Hashing
- Shared-Disk

Data Scalability

Write Scalability

Read Scalability

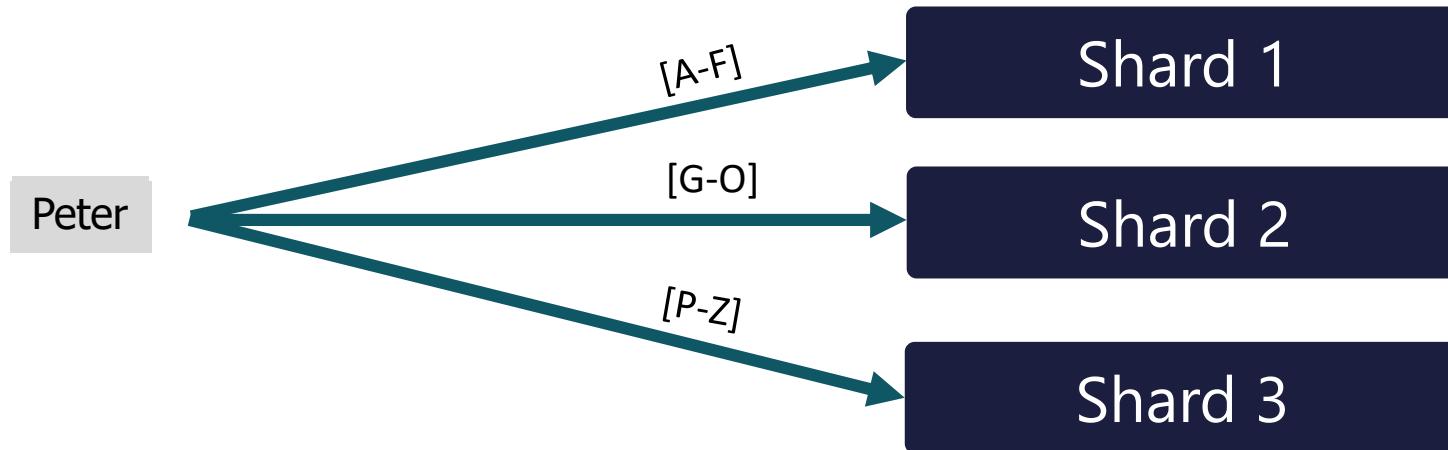
Elasticity



Sharding (aka Partitioning, Fragmentation)

Scaling Storage and Throughput

- ▶ Horizontal distribution of data over nodes



- ▶ **Partitioning strategies:** Hash-based vs. Range-based
- ▶ **Difficulty:** Multi-Shard-Operations (join, aggregation)

Sharding

Approaches

Hash-based Sharding

- Hash of data values (e.g. key) determines partition (shard)
- **Pro:** Even distribution
- **Contra:** No data locality

Implemented in

Range-based Sharding

- Assigns ranges defined over fields (shard keys) to partitions
- **Pro:** Enables *Range Scans* and *Sorting*
- **Contra:** Repartitioning/balancing required

Entity-Group Sharding

- Explicit data co-location for single-node-transactions
- **Pro:** Enables *ACID Transactions*
- **Contra:** Partitioning not easily changeable



David J DeWitt and Jim N Gray: "Parallel database systems: The future of high performance database systems," Communications of the ACM, volume 35, number 6, pages 85–98, June 1992.

Sharding

Approaches

Hash-based Sharding

- Hash of data values (e.g. key) determine shard
- Pro: Even distribution
- Contra: No data locality

Implemented in

MongoDB, Riak, Redis, Cassandra, Azure Table, Dynamo

Implemented in

BigTable, HBase, DocumentDB Hypertable, MongoDB, RethinkDB, Espresso

Implemented in

G-Store, MegaStore, Relational Cloud, Cloud SQL Server

Range-based Sharding

- Assigns ranges defined over fields
- Pro: Enables *Range Scans* and *Scalability*
- Contra: Repartitioning/balancing

Entity-Group Sharding

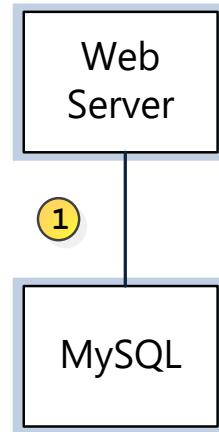
- Explicit data co-location for single entities
- Pro: Enables *ACID Transactions* and *Consistency*
- Contra: Partitioning not easily done



Problems of Application-Level Sharding

Example: Tumblr

- ▶ Caching
- ▶ Sharding from application



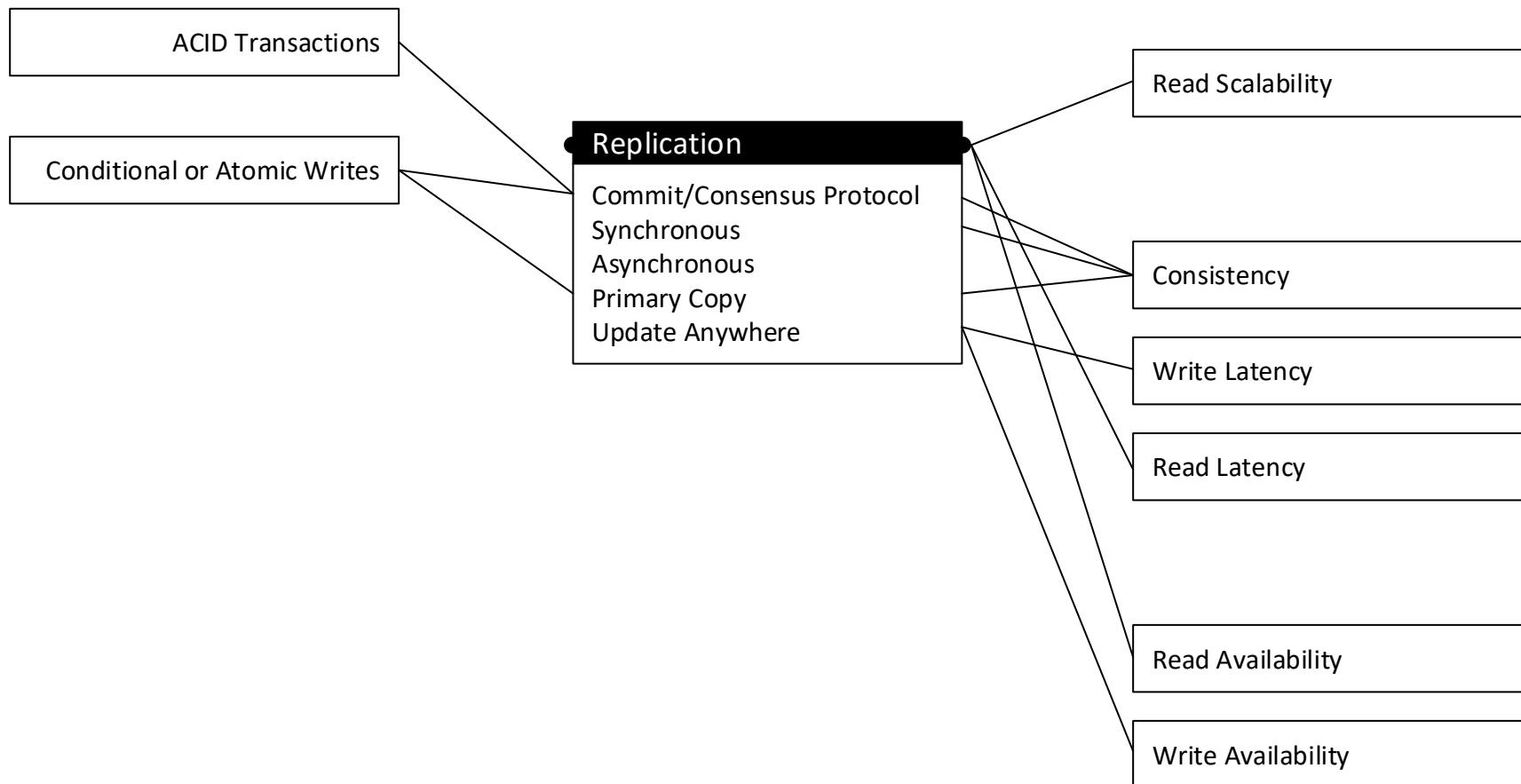
Moved towards:

- ▶ Redis
- ▶ HBase

Functional

Techniques

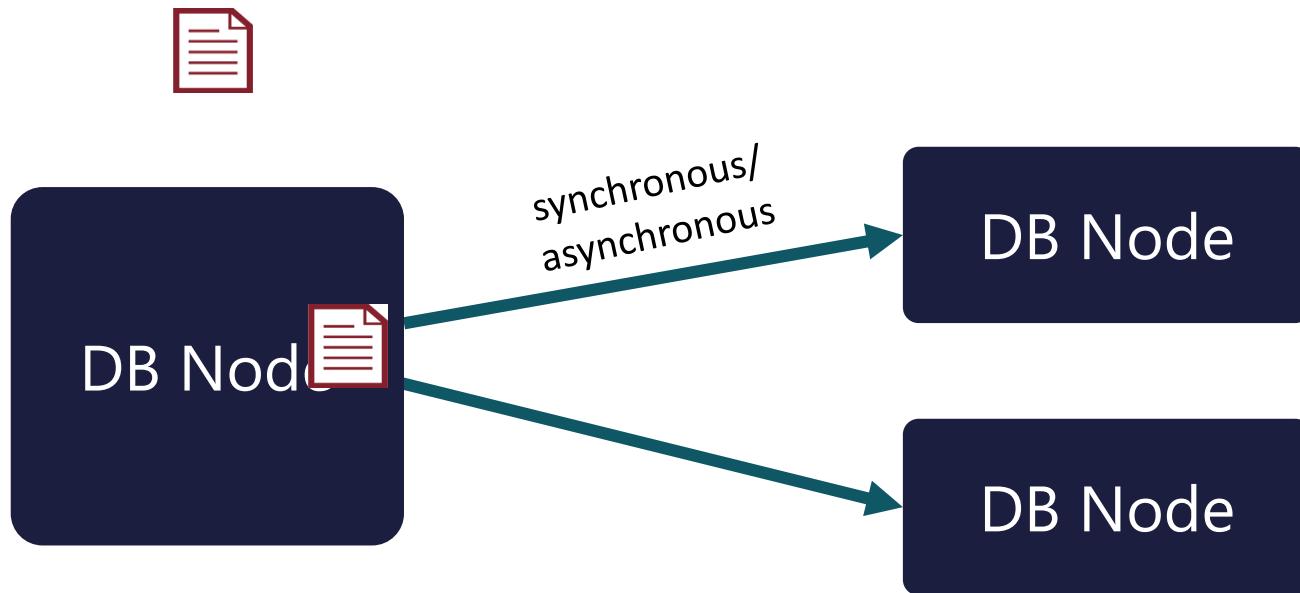
Non-Functional



Replication

Read Scalability + Failure Tolerance

- ▶ Stores N copies of each data item



- ▶ **Consistency model:** synchronous vs asynchronous
- ▶ **Coordination:** Multi-Master, Master-Slave



Replication: When

Asynchronous (lazy)

- Writes are acknowledged immediately
- Performed through *log shipping* or *update propagation*
- **Pro:** Fast writes, no coordination needed
- **Contra:** Replica data potentially stale (*inconsistent*)

Synchronous (eager)

- The node accepting writes synchronously propagates updates/transactions before acknowledging
- **Pro:** Consistent
- **Contra:** needs a commit protocol (more roundtrips), unavailable under certain network partitions



Replication: When

Asynchronous (lazy)

- Writes are acknowledged immediately
- Performed through *log shipping*
- **Pro:** Fast writes, no coordination required
- **Contra:** Replica data potentially stale

Implemented in

Dynamo, Riak, CouchDB, Redis, Cassandra, Voldemort, MongoDB, RethinkDB

Synchronous (eager)

- The node accepting writes synchronizes updates/transactions before acknowledging them
- **Pro:** Consistent
- **Contra:** needs a commit protocol, unavailable under certain network partitions

Implemented in

BigTable, HBase, Accumulo, CouchBase, MongoDB, RethinkDB



Replication: Where

Master-Slave (*Primary Copy*)

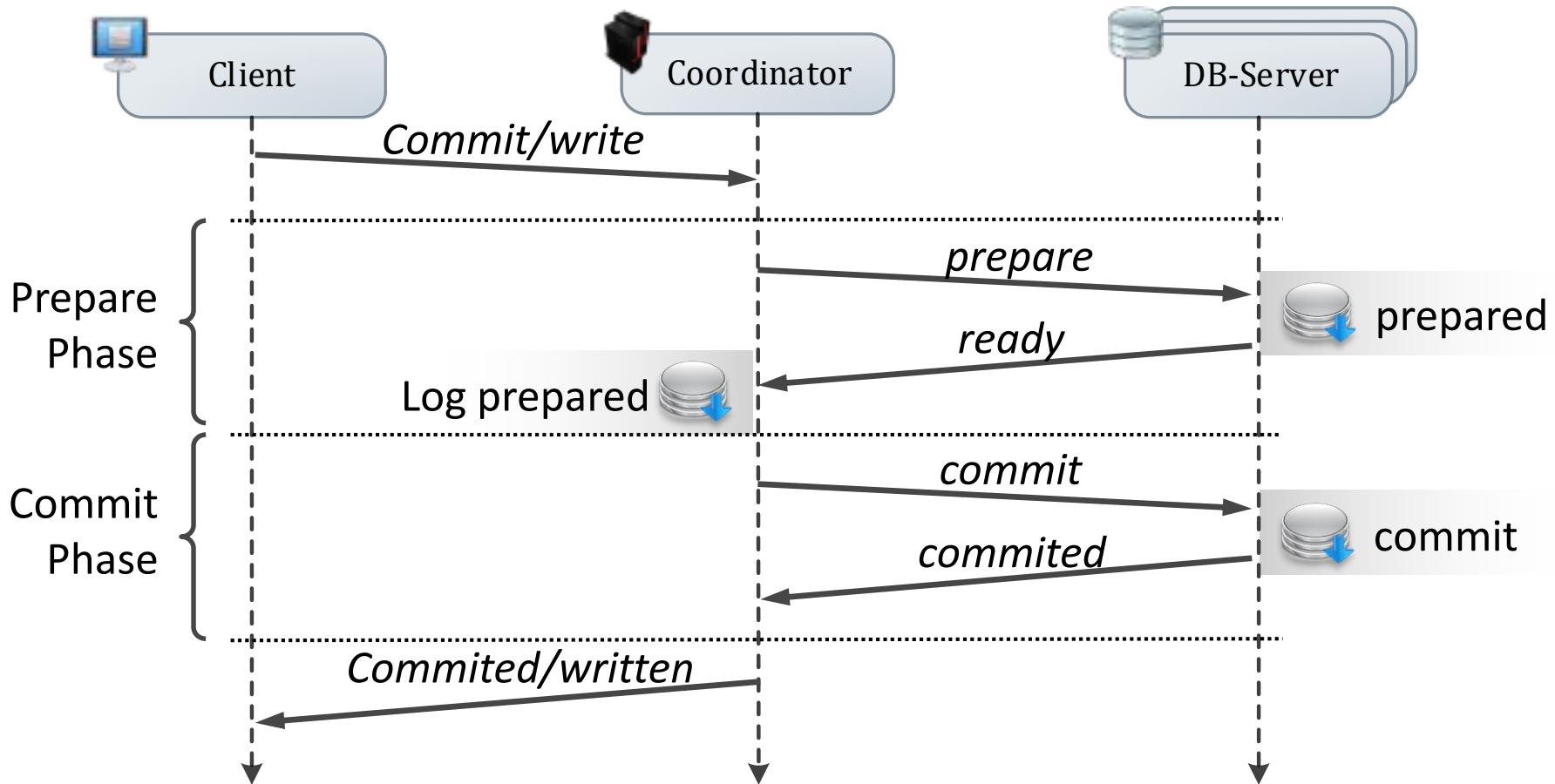
- Only a dedicated master is allowed to accept writes, slaves are read-replicas
- **Pro:** reads from the master are consistent
- **Contra:** master is a bottleneck and SPOF

Multi-Master (*Update anywhere*)

- The server node accepting the writes synchronously propagates the update or transaction before acknowledging
- **Pro:** fast and highly-available
- **Contra:** either needs coordination protocols (e.g. Paxos) or is inconsistent



Synchronous Replication: 2PC



Error scenarios (timeouts)

In INITIAL:

- No consequence

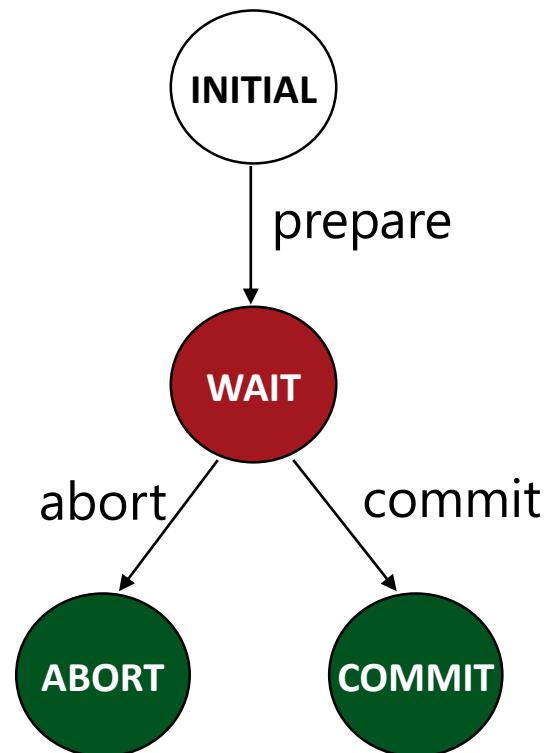
In WAIT:

- Abort

In ABORT oder COMMIT:

- Resend *commit/abort* and wait for all responses

Coordinator



Error scenarios (timeouts)

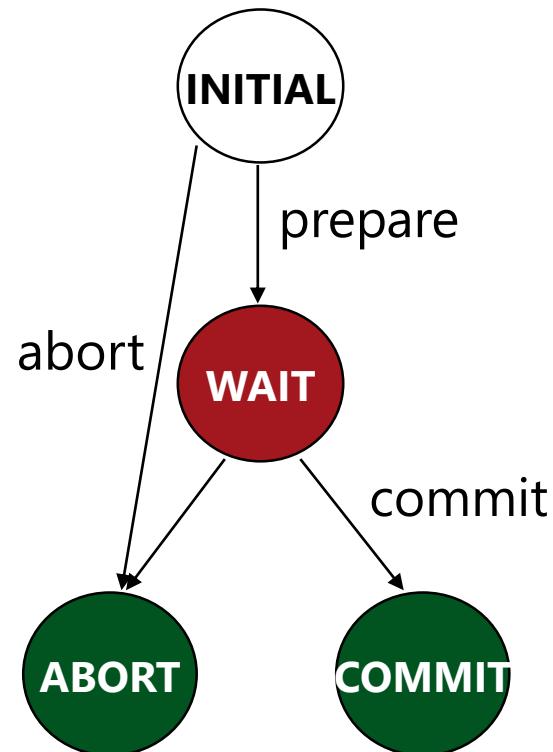
In INITIAL:

- Coordinator probably crashed → Abort („*Presumed-Abort-Protocol*“)

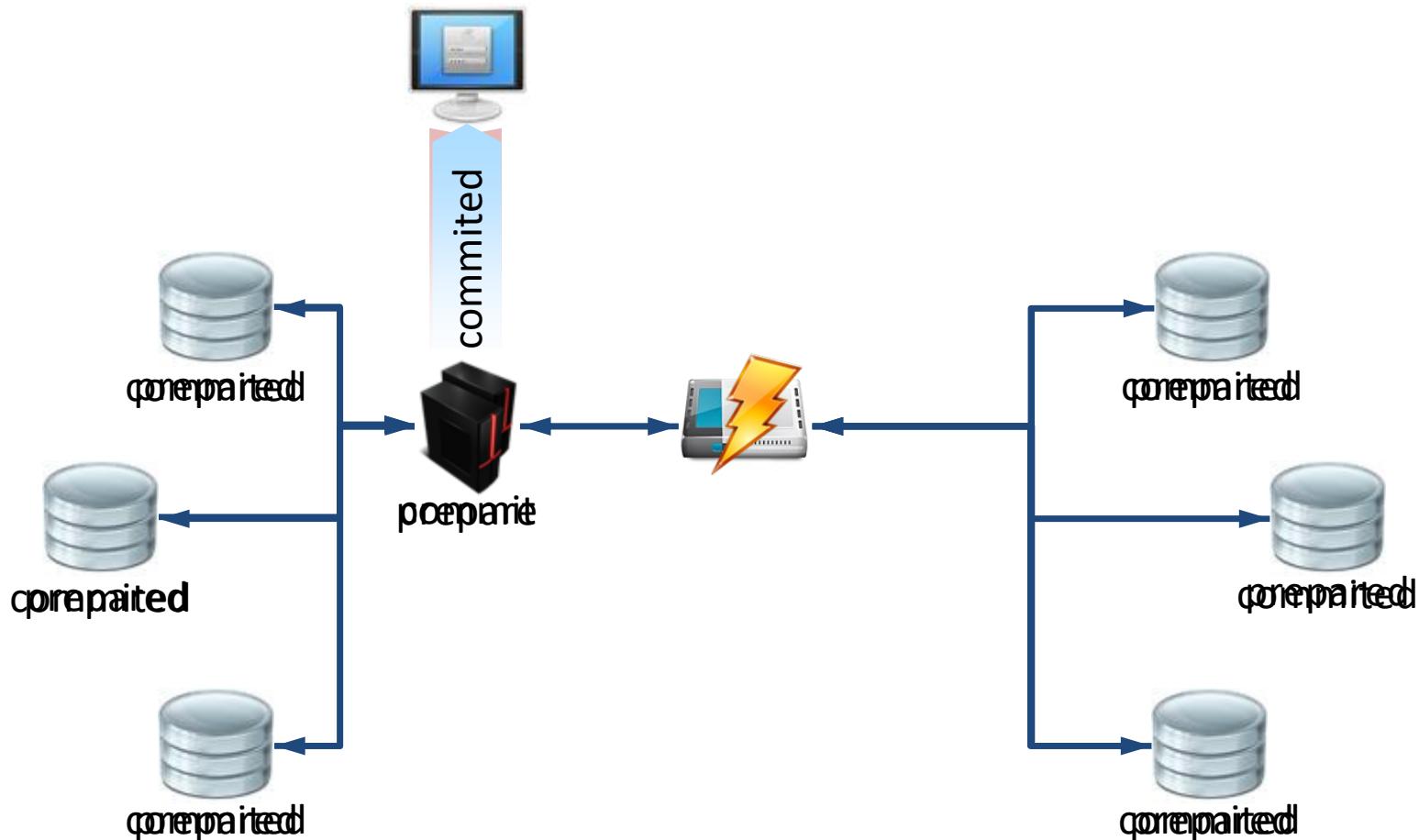
In WAIT:

- Wait for message from coordinator

Resource-Manager

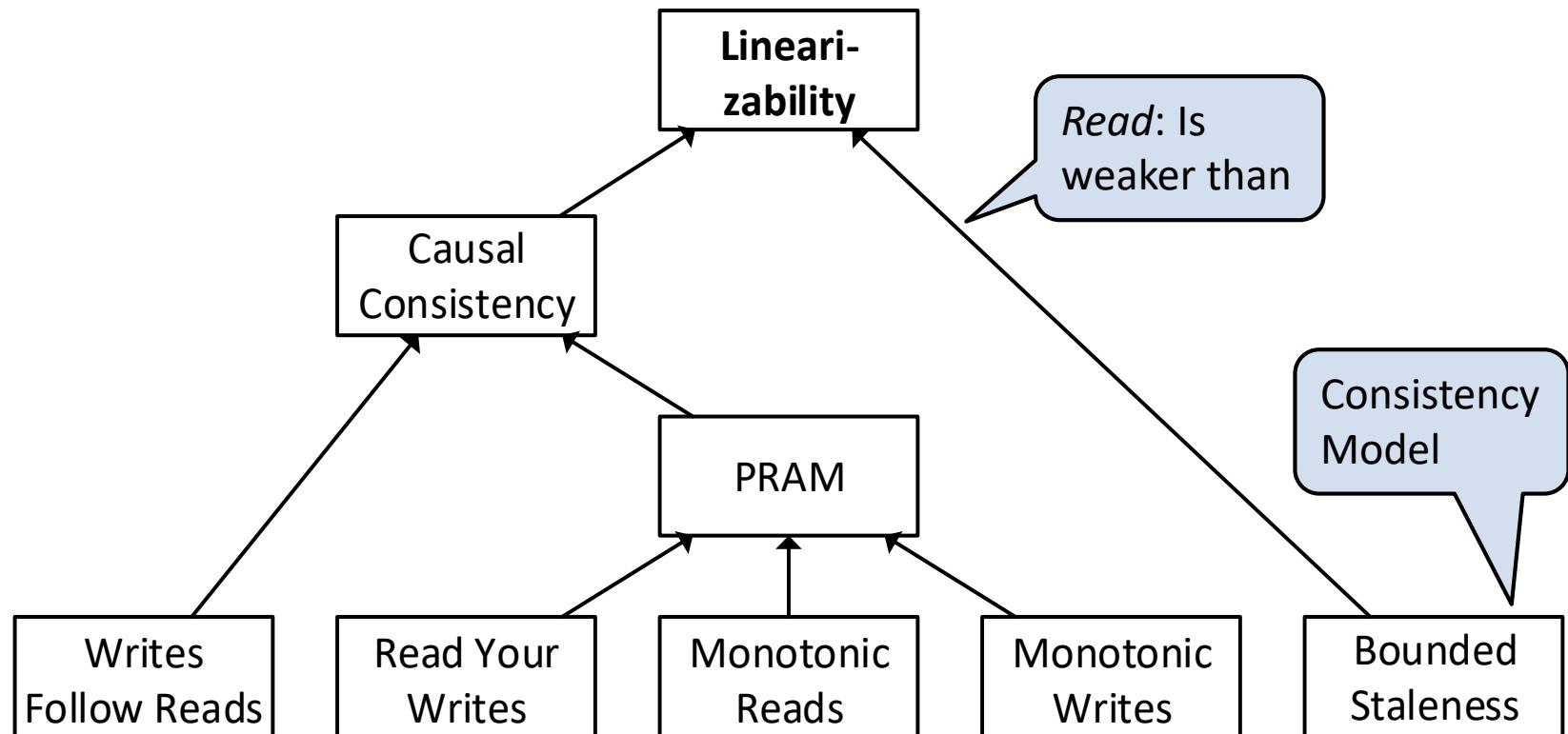


2PC is not *highly available*



NoSQL Consistency Models

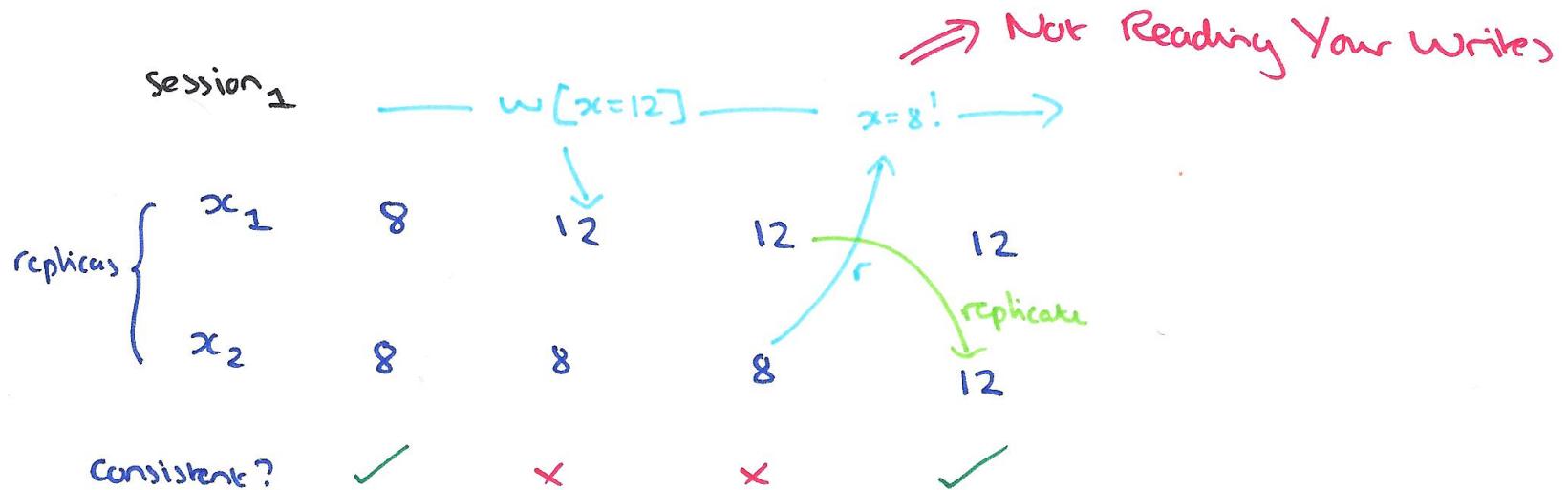
- ▶ Define models that relax strong consistency (=linearizability) in different aspects



Bailis, Peter, et al. "Highly available transactions: Virtues and limitations." Proceedings of the VLDB Endowment 7.3 (2013): 181-192.

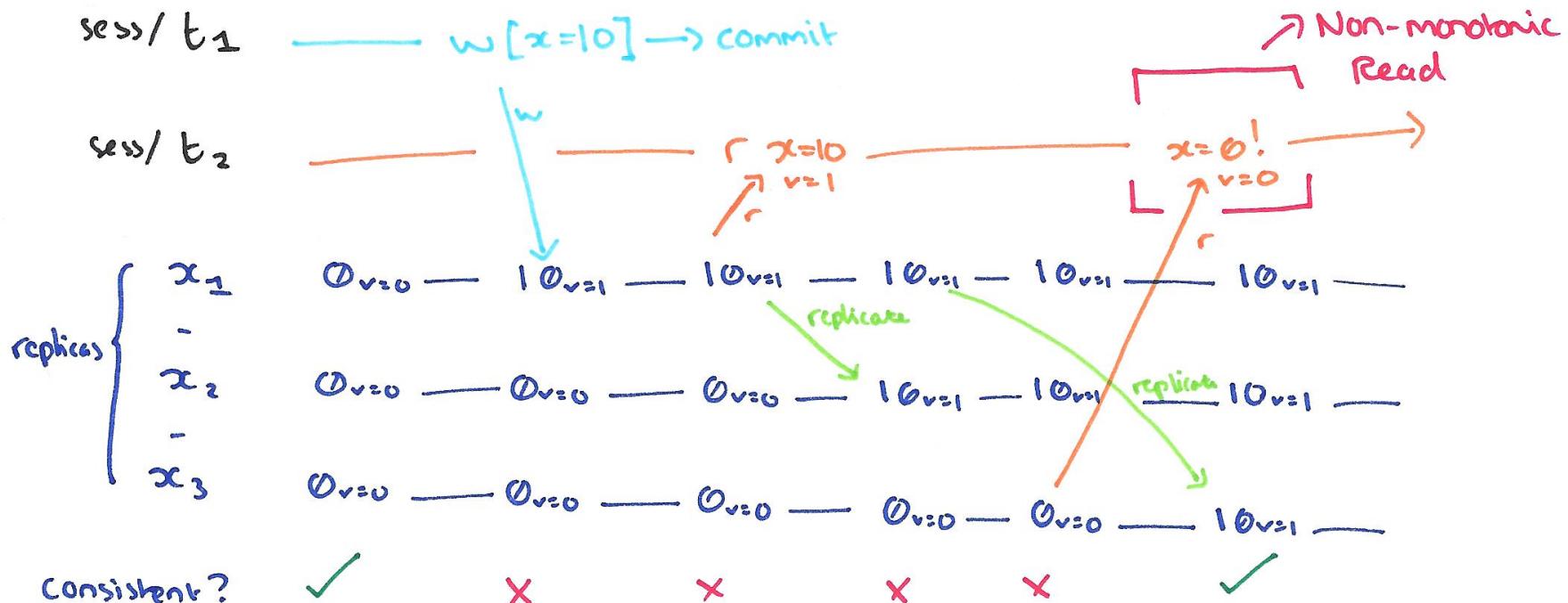
Read Your Writes (RYW)

Definition: Once the user has written a value, subsequent reads will return this value (or newer versions if other writes occurred in between); the user will never see versions older than his last write.



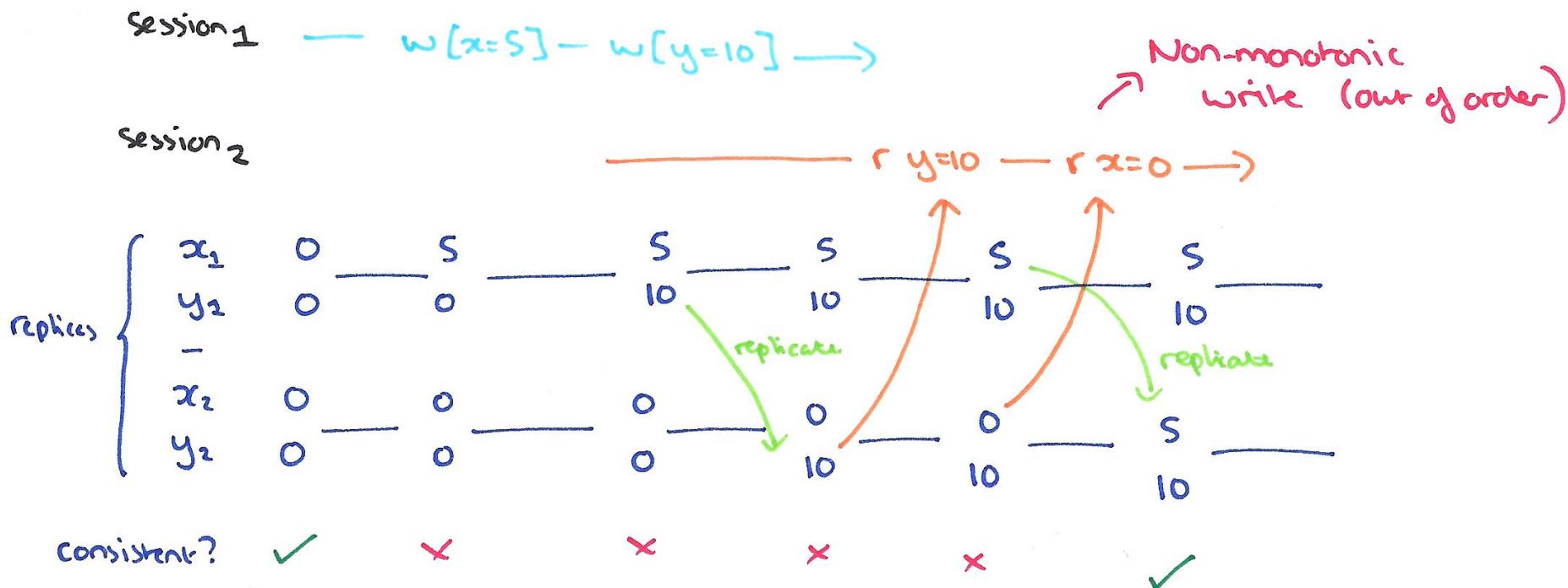
Monotonic Reads (MR)

Definition: Once a user has read a version of a data item on one replica server, it will never see an older version on any other replica server



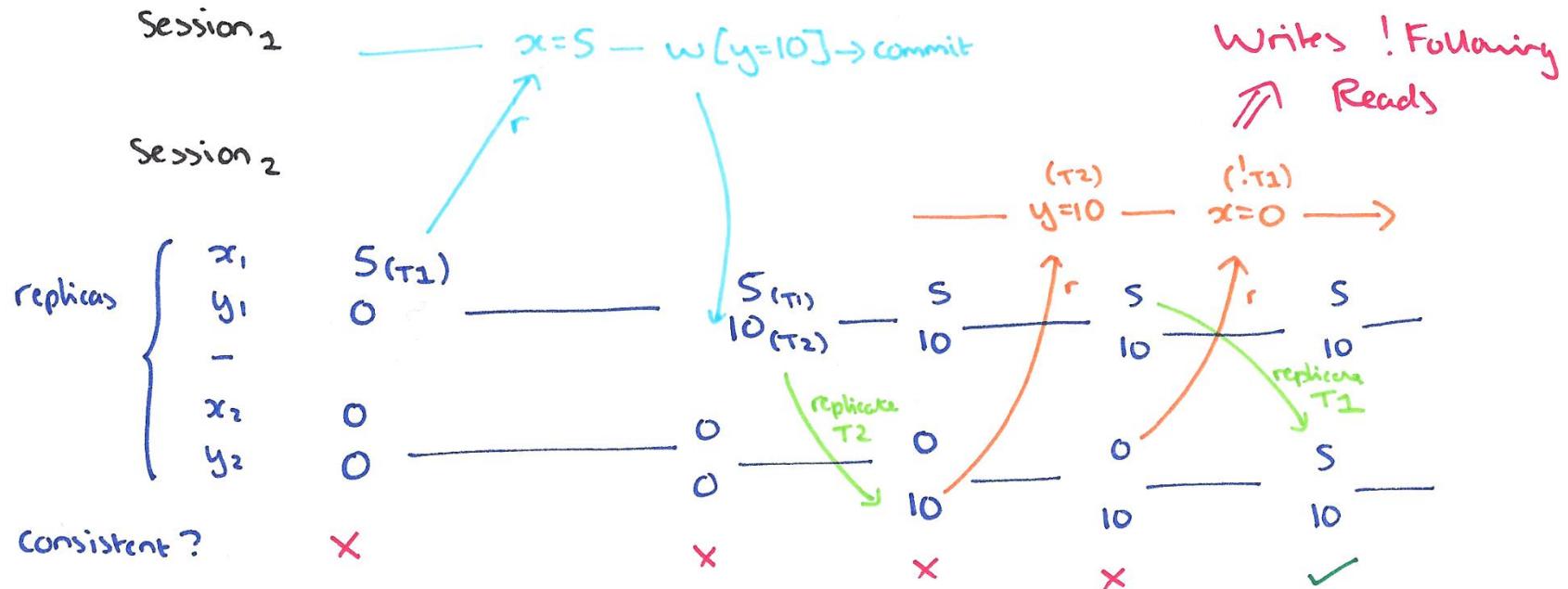
Monotonic Writes (MW)

Definition: Once a user has written a new value for a data item in a session, any previous write has to be processed before the current one. I.e., the order of writes inside the session is strictly maintained.



Writes Follow Reads (WFR)

Definition: When a user reads a value written in a session after that session already read some other items, the user must be able to see those *causally relevant* values too.



Bounded Staleness

- ▶ Either **time-based**:

t-Visibility (Δ -atomicity): the inconsistency window comprises at most t time units; that is, any value that is returned upon a read request was up to date t time units ago.

- ▶ Or **version-based**:

k-Staleness: the inconsistency window comprises at most k versions; that is, lags at most k versions behind the most recent version.

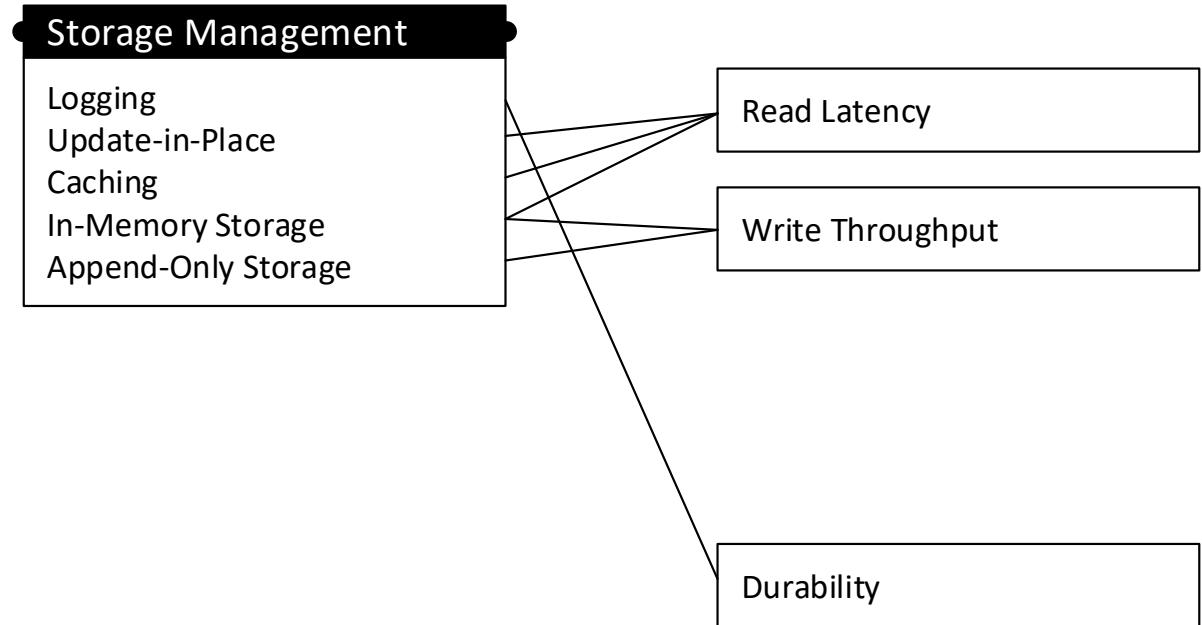
- ▶ Both are *not* achievable with high availability



Functional

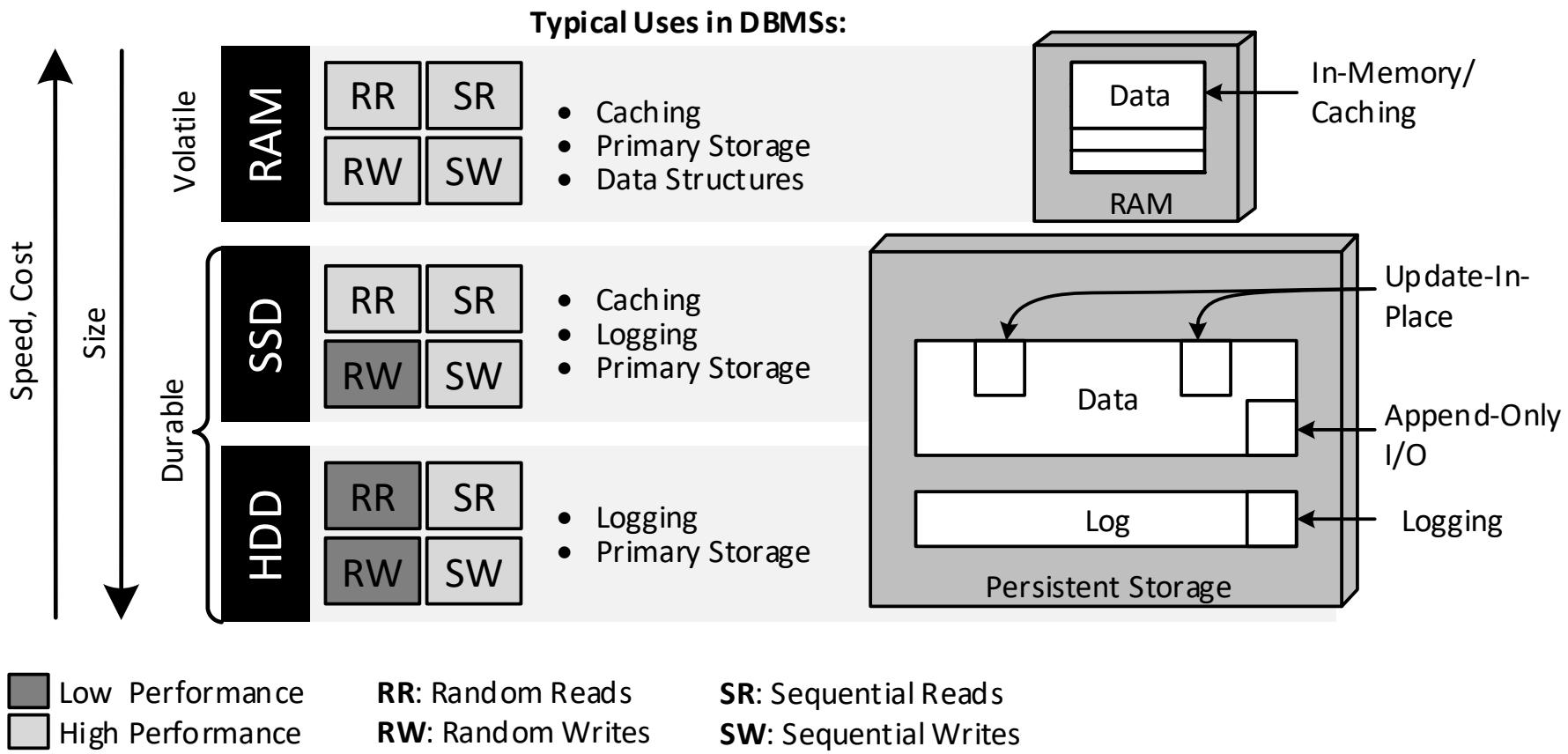
Techniques

Non-Functional



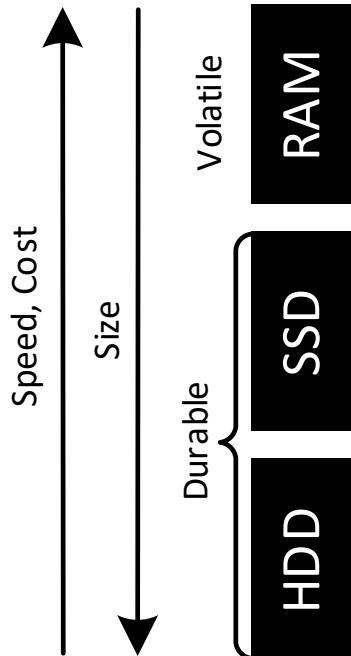
NoSQL Storage Management

In a Nutshell



NoSQL Storage Management

In a Nutshell

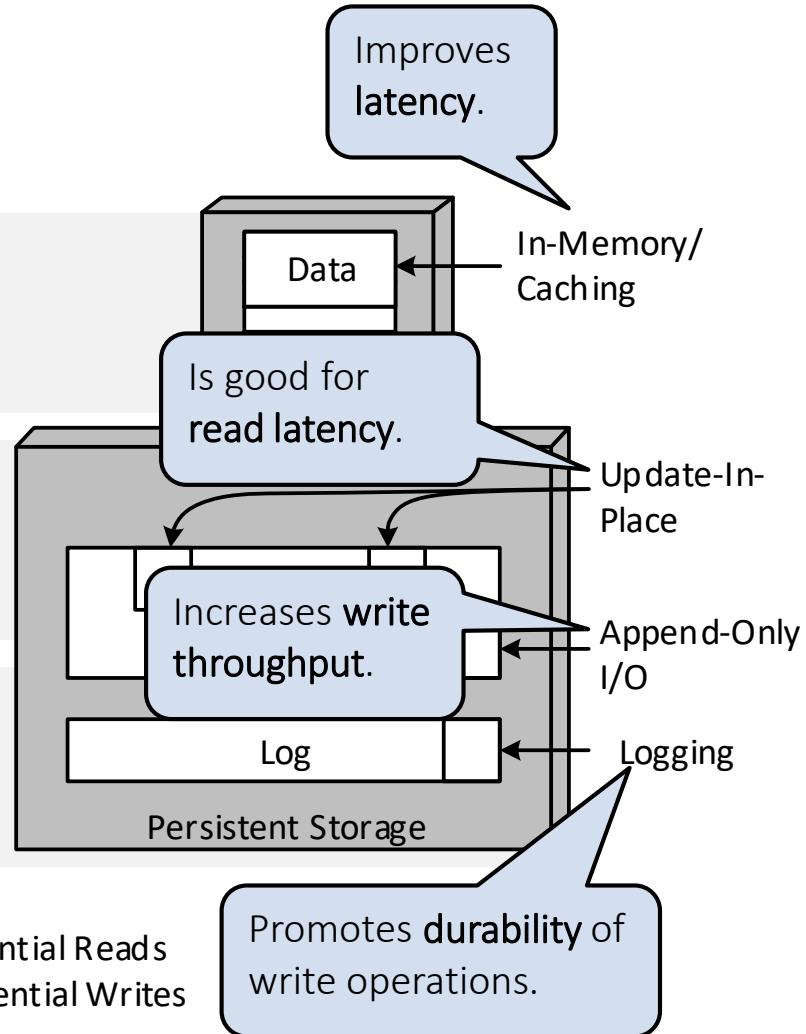


Typical Uses in DBMSs:

- Caching
- Primary Storage
- Data Structures

- Caching
- Logging
- Primary Storage

- Logging
- Primary Storage



Low Performance
High Performance

RR: Random Reads

RW: Random Writes

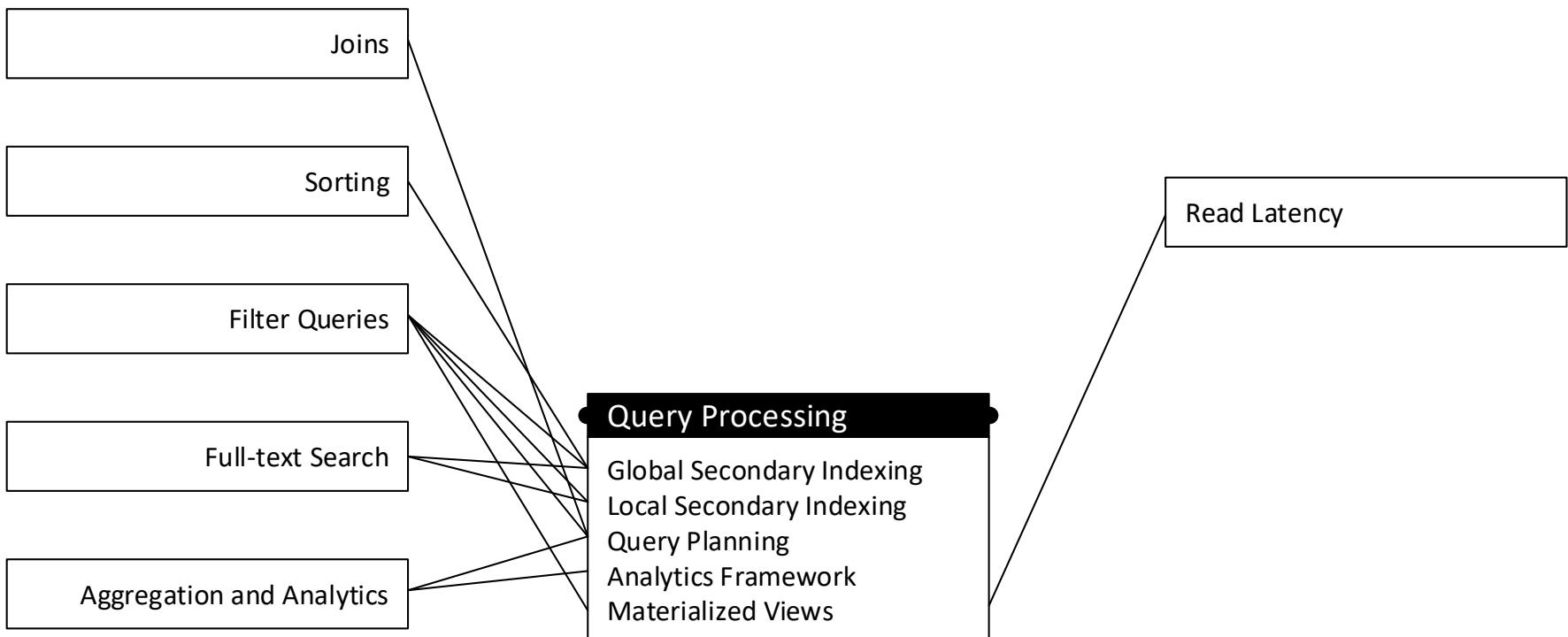
SR: Sequential Reads

SW: Sequential Writes

Functional

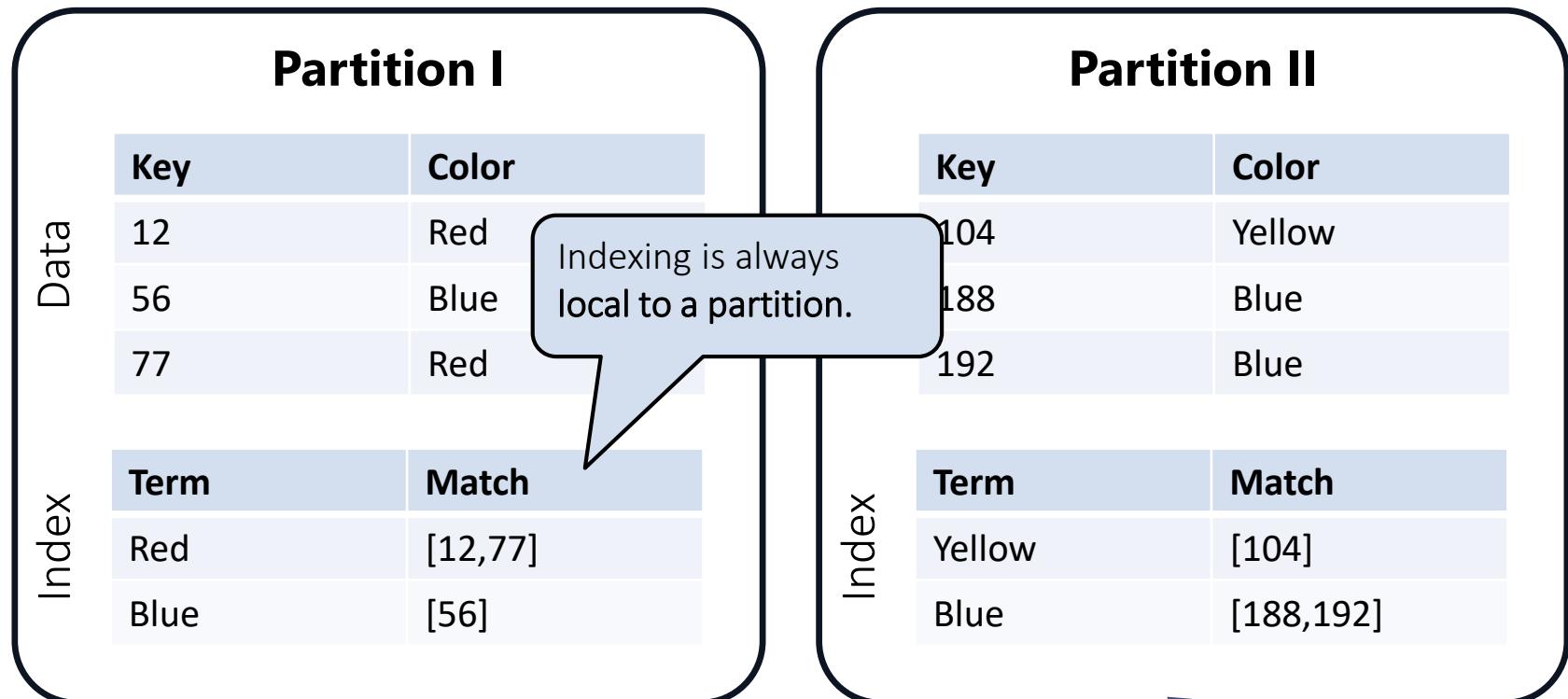
Techniques

Non-Functional



Local Secondary Indexing

Partitioning By Data Item



Scatter-gather query pattern.

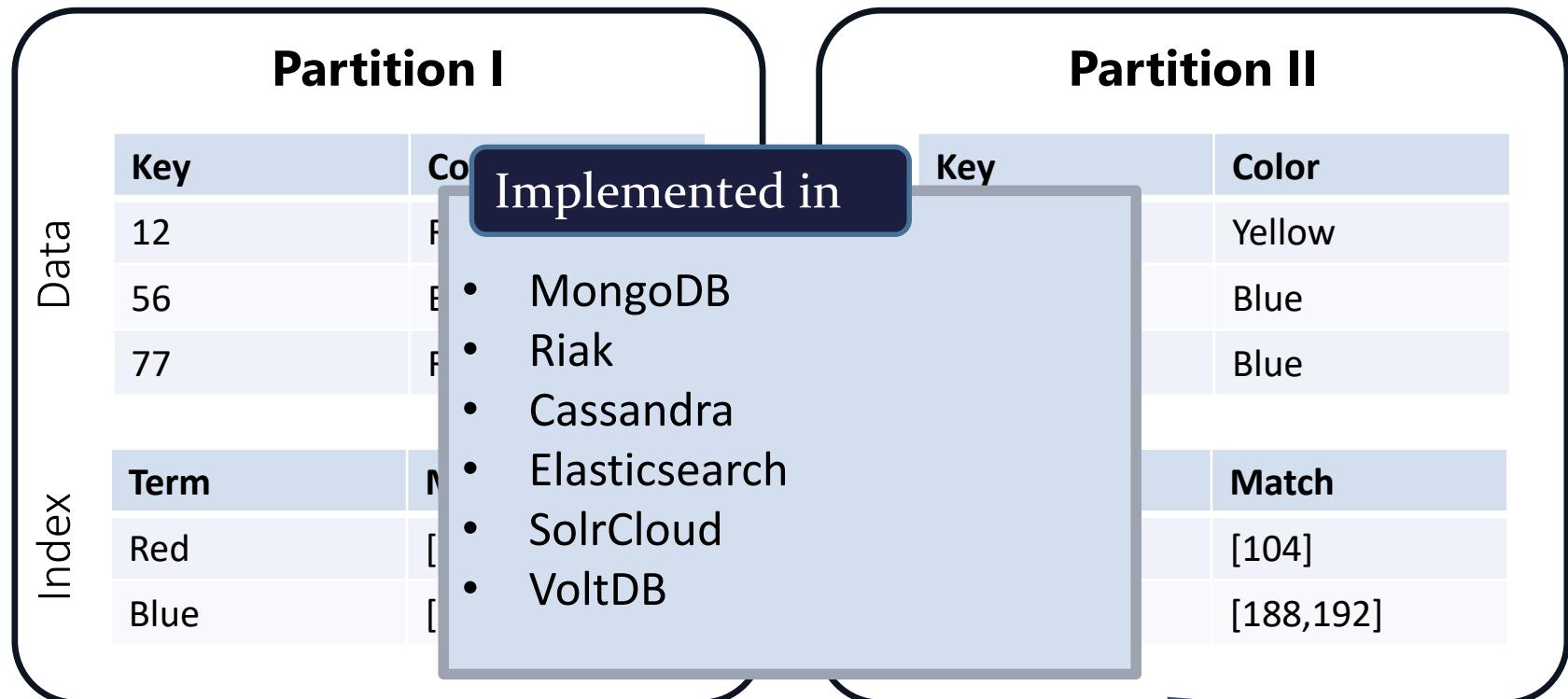
WHERE color=blue



Kleppmann, Martin. "Designing data-intensive applications." (2016).

Local Secondary Indexing

Partitioning By Data Item



Scatter-gather query pattern.

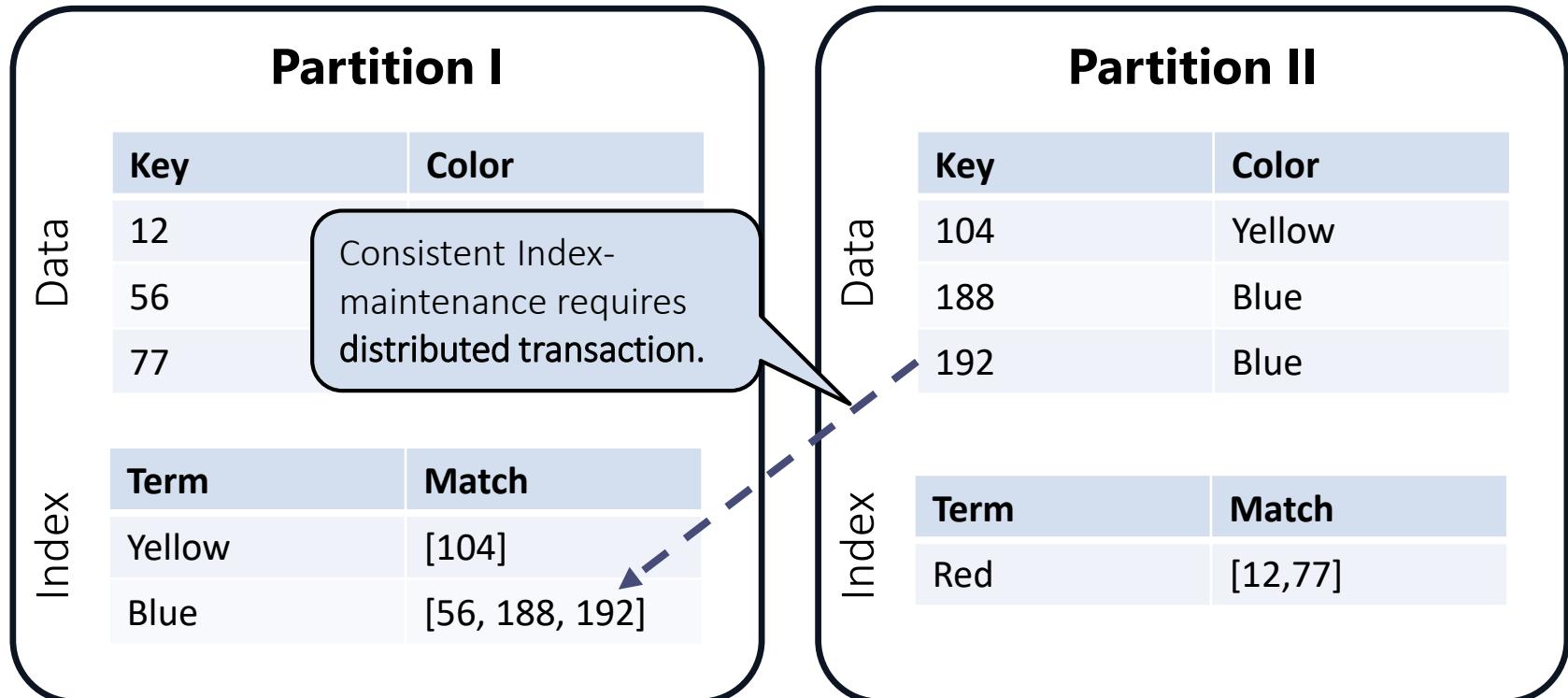
WHERE color=blue



Kleppmann, Martin. "Designing data-intensive applications." (2016).

Global Secondary Indexing

Partitioning By Query Term



Targeted Query

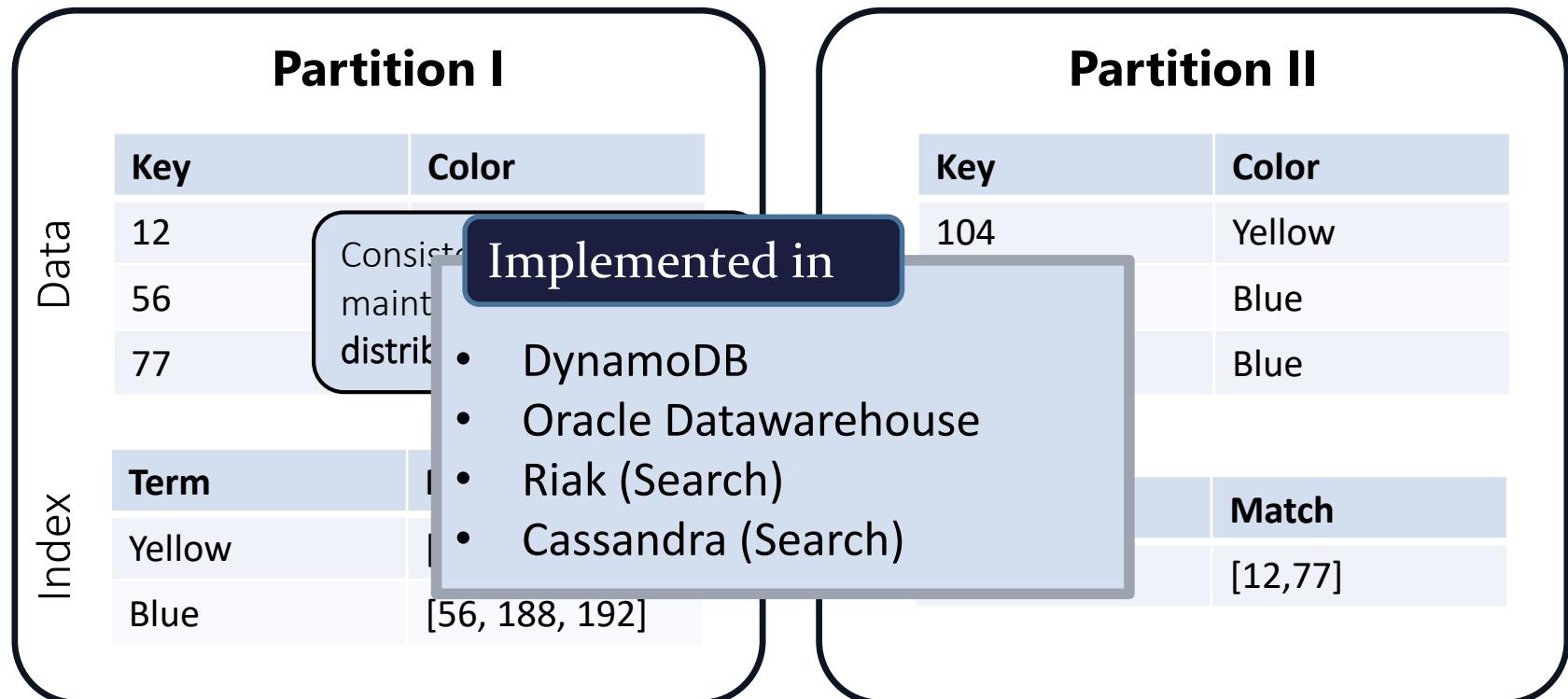
WHERE color=blue



Kleppmann, Martin. "Designing data-intensive applications." (2016).

Global Secondary Indexing

Partitioning By Query Term



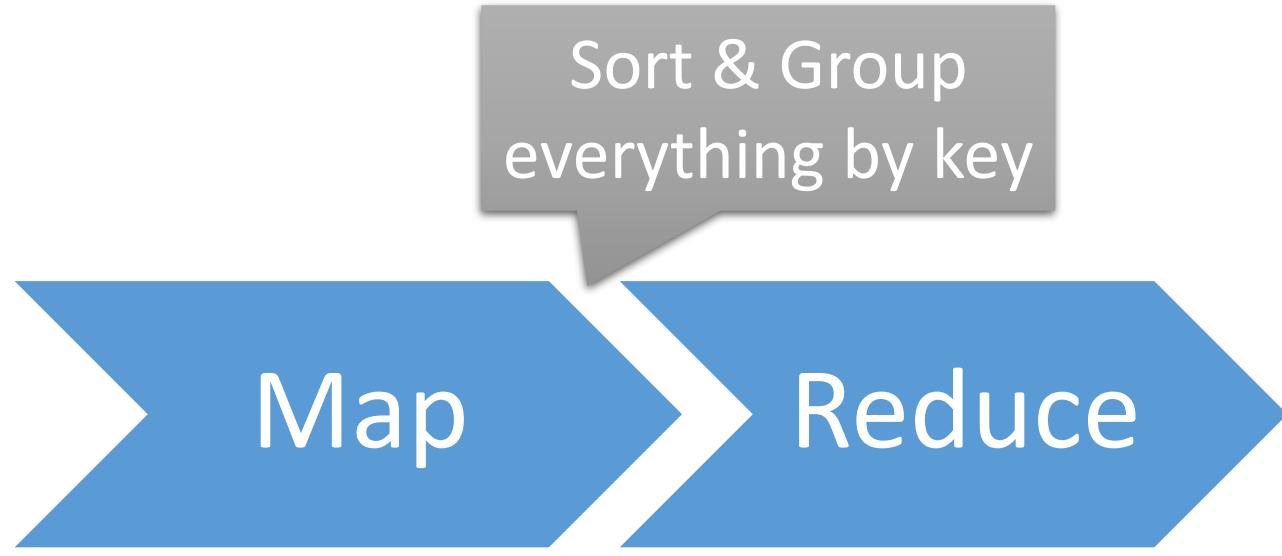
Targeted Query

WHERE color=blue



Kleppmann, Martin. "Designing data-intensive applications." (2016).

The MapReduce query concept



Convert Each Input to a key-value pair

Process the list of values for each key

MapReduce: Word Count

```
map(zeilennr, text) :
```

1: ich bin ich

```
    for each word in text:
```

```
        emit(word, 1)
```

(ich, 1)

(bin, 1)

(ich, 1)

```
reduce(word, values) :
```

(ich, [1,1])

```
    sum = 0
```

```
    for each v in values:
```

(bin, [1])

```
        sum = sum + v
```

```
    emit(word, sum)
```

(ich, 2)

(bin, 1)

Query Processing Techniques

Summary

- ▶ **Local Secondary Indexing:** Fast writes, scatter-gather queries
- ▶ **Global Secondary Indexing:** Slow or inconsistent writes, fast queries
- ▶ **(Distributed) Query Planning:** scarce in NoSQL systems but increasing (e.g. left-outer equi-joins in MongoDB and θ -joins in RethinkDB)
- ▶ **Analytics Frameworks:** fallback for missing query capabilities (e.g. through MapReduce)
- ▶ **Materialized Views:** similar to global indexing



Next Lecture: How are the techniques from the NoSQL toolbox used in actual data stores?

Summary



- ▶ High **data volumes**, unstructured sources and **new applications** triggered NoSQL technologies
- ▶ Shared Nothing architectures for **horizontal scalability**
 - **Replication** enables read scalability and fault tolerance
 - **Sharding** enables write scalability and data volume scalability
- ▶ **CAP Theorem**: Consistency, Availability and Partition Tolerance cannot be achieved at the same time
- ▶ **4 NoSQL Data Model**: Key-Value Stores, Document Stores, Wide-Column Stores, Graph Databases
- ▶ **2-Phase-Commit (2PC)**: popular protocol for atomic commitment; without availability guarantee
- ▶ **Consistency** can be relaxed in various ways

Outline



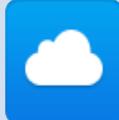
NoSQL Foundations and Motivation



The NoSQL Toolbox:
Common Techniques



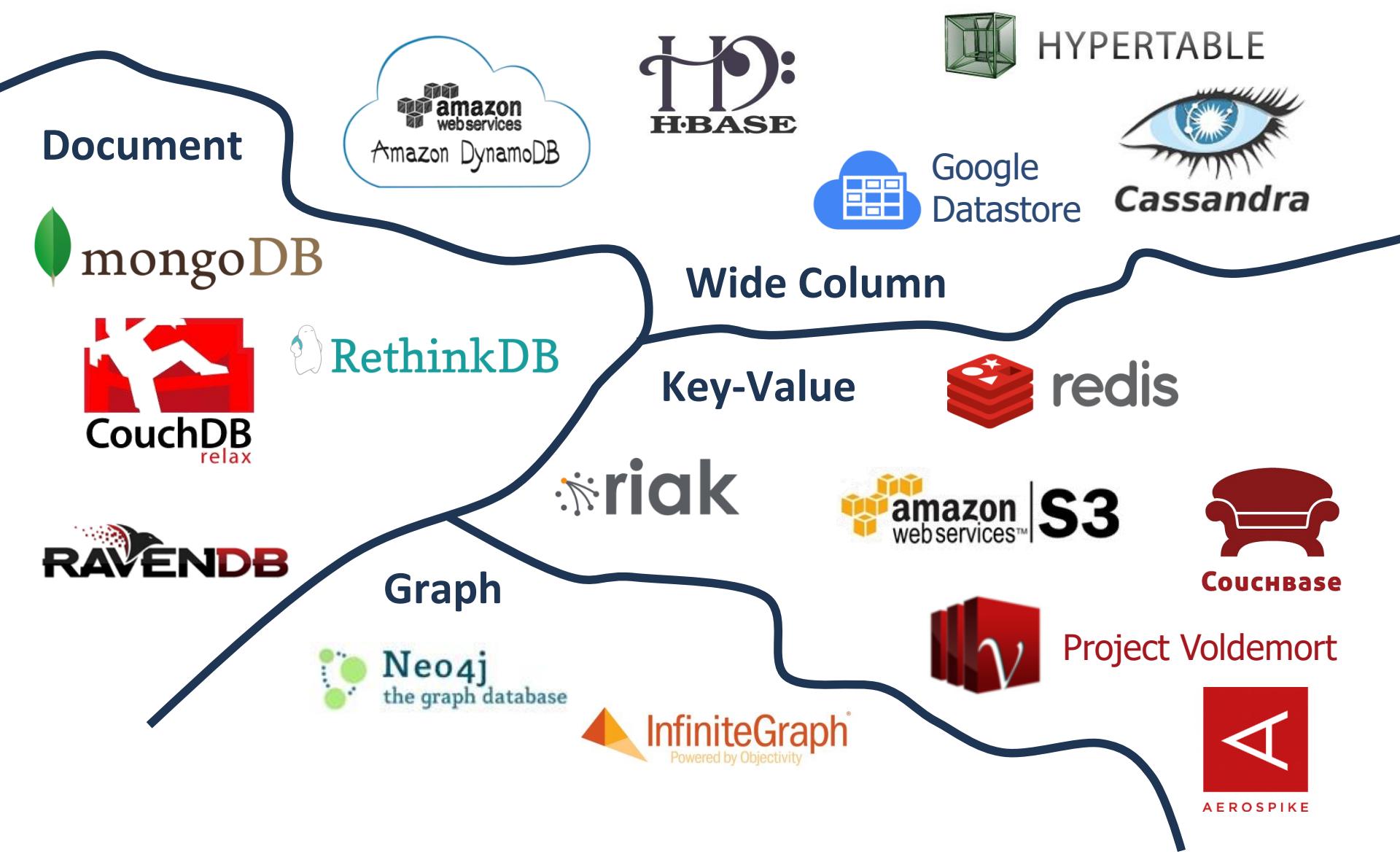
NoSQL Systems &
Decision Guidance



Scalable Real-Time
Databases and Processing

- Overview & Popularity
- Core Systems:
 - Dynamo
 - BigTable
- Riak
- HBase
- Cassandra
- Redis
- MongoDB

NoSQL Landscape



Popularity (Feb 2019)

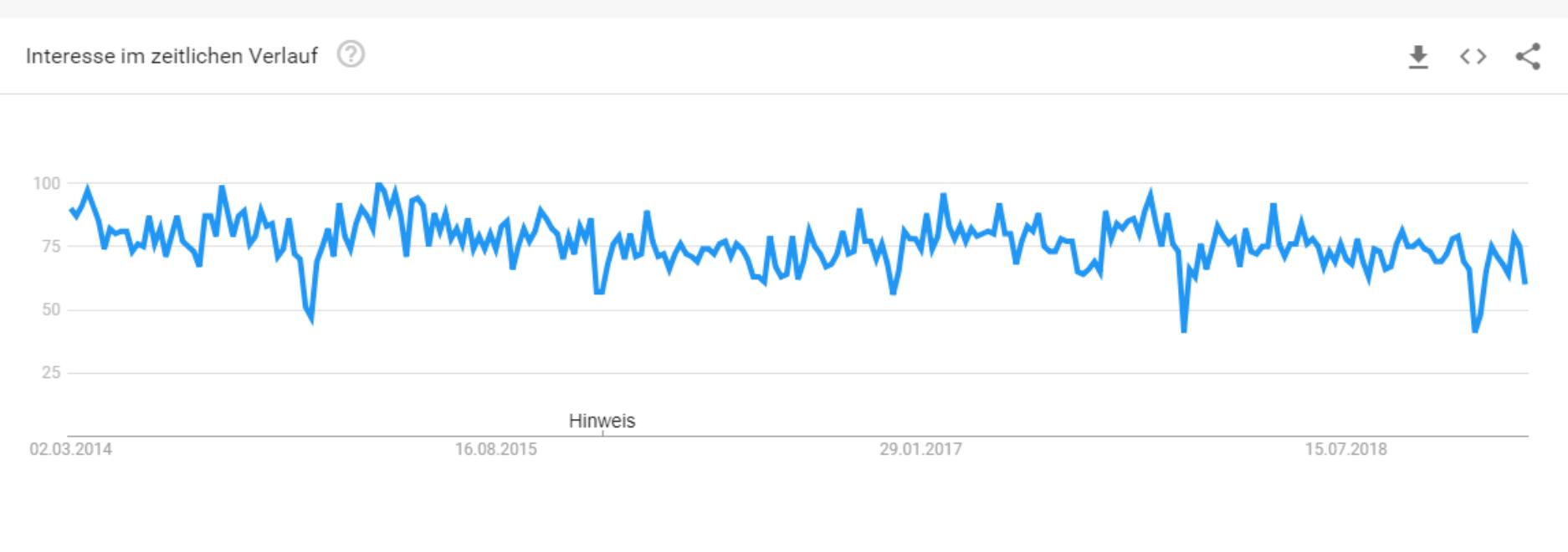
#	System	Model
1.	Oracle	Relational DBMS
2.	MySQL	Relational DBMS
3.	MS SQL Server	Relational DBMS
4.	PostgreSQL	Relational DBMS
5.	MongoDB	Document store
6.	DB2	Relational DBMS
7.	Microsoft Access	Relational
8.	Redis	Key-value store
9.	ElasticSearch	Search engine
10.	SQLite	Relational DBMS

11.	Cassandra	Wide column store
12.	MariaDB	Relational DBMS
13.	Splunk	Search engine
14.	Teradata	Search engine
15.	Hive	Relational
16.	Solr	Relational DBMS
17.	HBase	Relational DBMS
18.	FileMaker	Relational
19.	SAP Adaptive Server	Relational DBMS
20.	SAP HANA	Relational DBMS
21.	Amazon DynamoDB	Multi-model
22.	Neo4j	Graph DB
23.	Couchbase	Document store
24.	Memcached	Key-value store
25.	SQL Azure	Multi-model

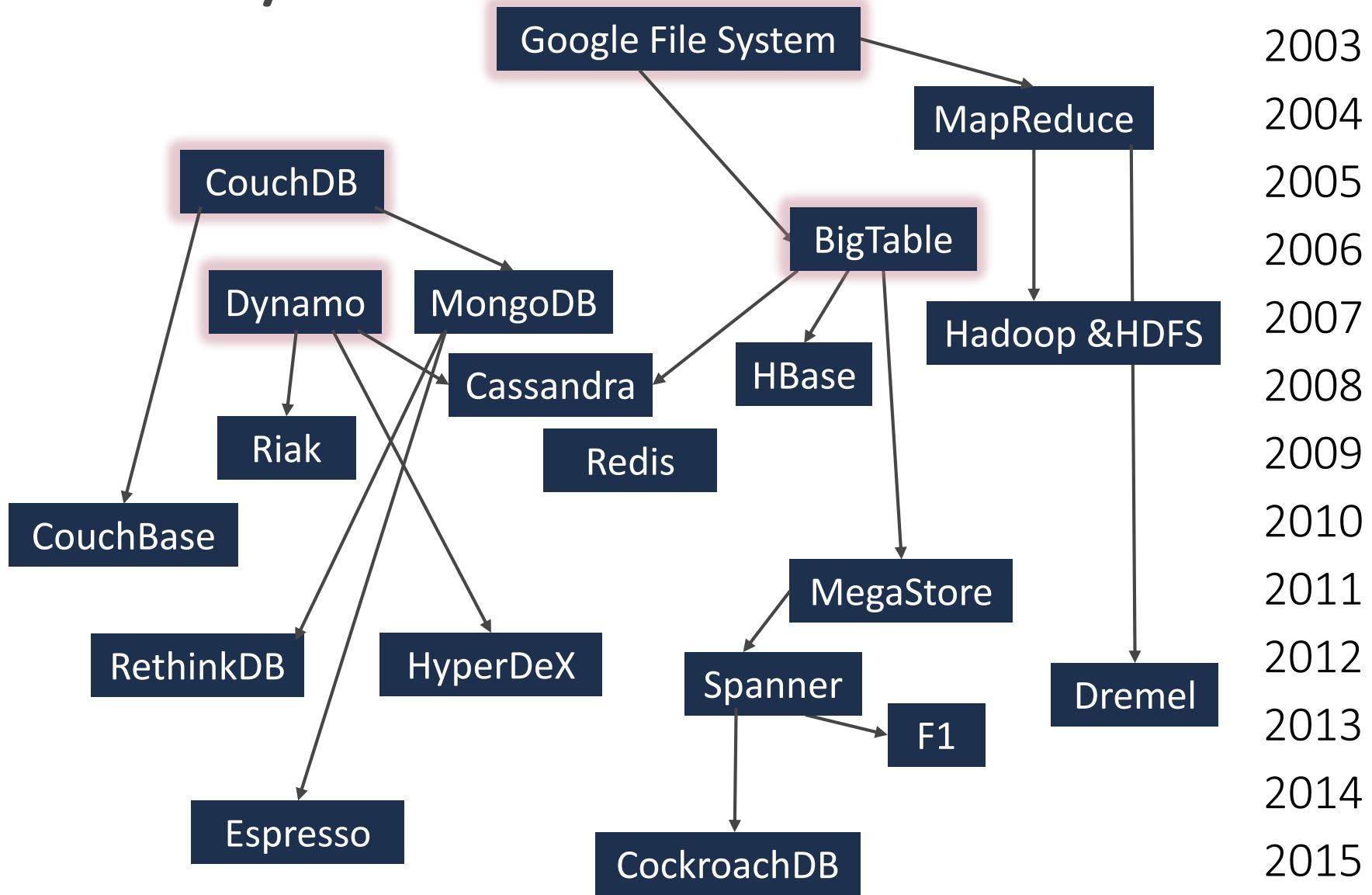
Scoring: Google/Bing results, Google Trends, Stackoverflow, job offers, LinkedIn



NoSQL: Still a Thing in 2019



History



NoSQL foundations

- ▶ **BigTable** (2006, Google)
 - Consistent, Partition Tolerant
 - **Wide-Column** data model
 - Master-based, fault-tolerant, large clusters (1.000+ Nodes),
HBase, Cassandra, HyperTable, Accumulo
- ▶ **Dynamo** (2007, Amazon)
 - Available, Partition tolerant
 - **Key-Value** interface
 - Eventually Consistent, always writable, fault-tolerant
 - **Riak, Cassandra, Voldemort, DynamoDB**



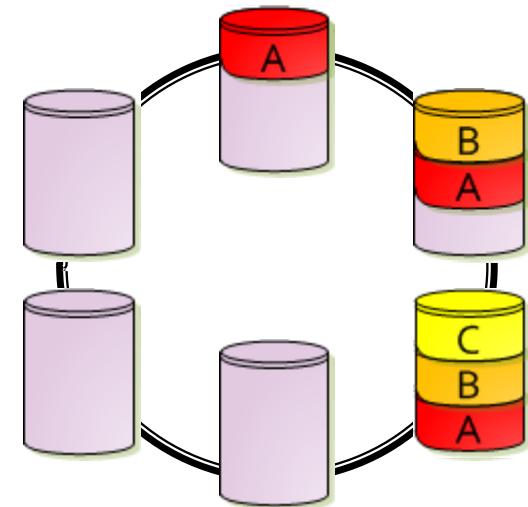
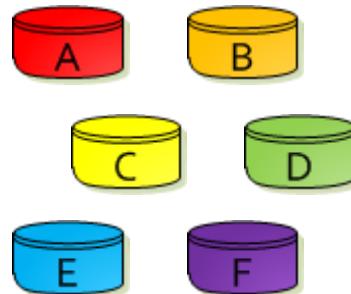
Chang, Fay, et al. "Bigtable: A distributed storage system for structured data."



DeCandia, Giuseppe, et al. "Dynamo: Amazon's highly available key-value store."

Dynamo (AP)

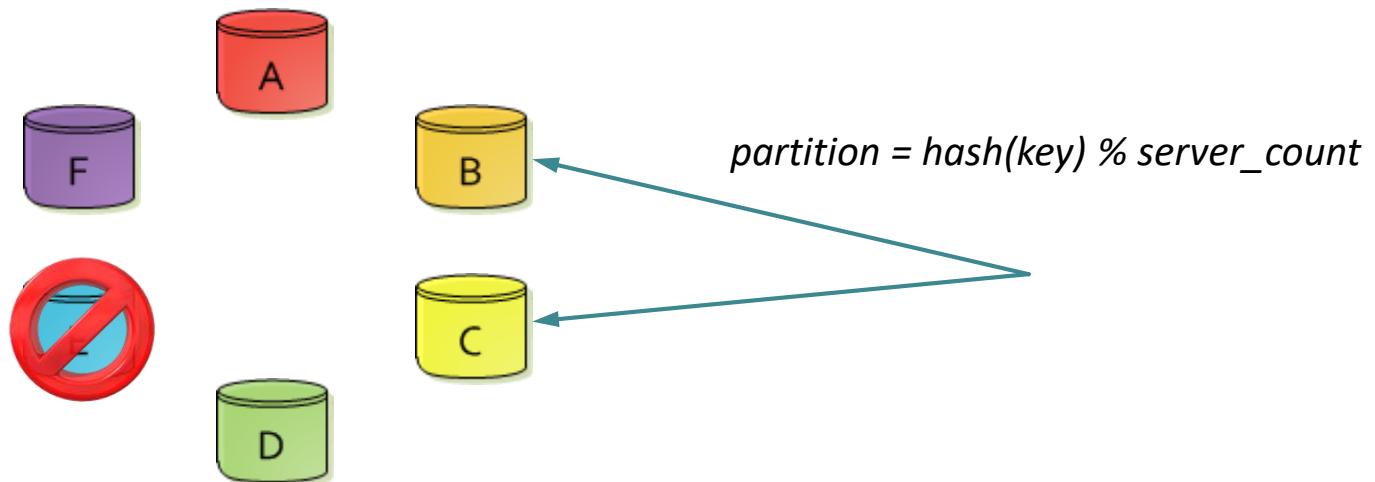
- ▶ Developed at Amazon (2007)
- ▶ Sharding of data over a ring of nodes
- ▶ Each node holds multiple partitions
- ▶ Each partition replicated **N** times



DeCandia, Giuseppe, et al. "Dynamo: Amazon's highly available key-value store."

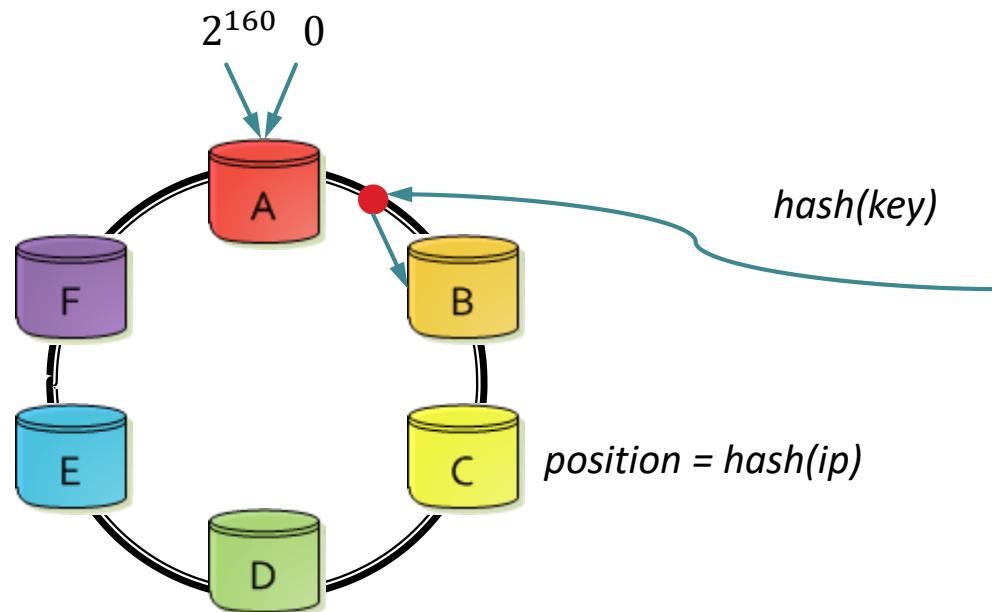
Consistent Hashing

- ▶ Naive approach: Hash-partitioning (e.g. in Memcache, Redis Cluster)



Consistent Hashing

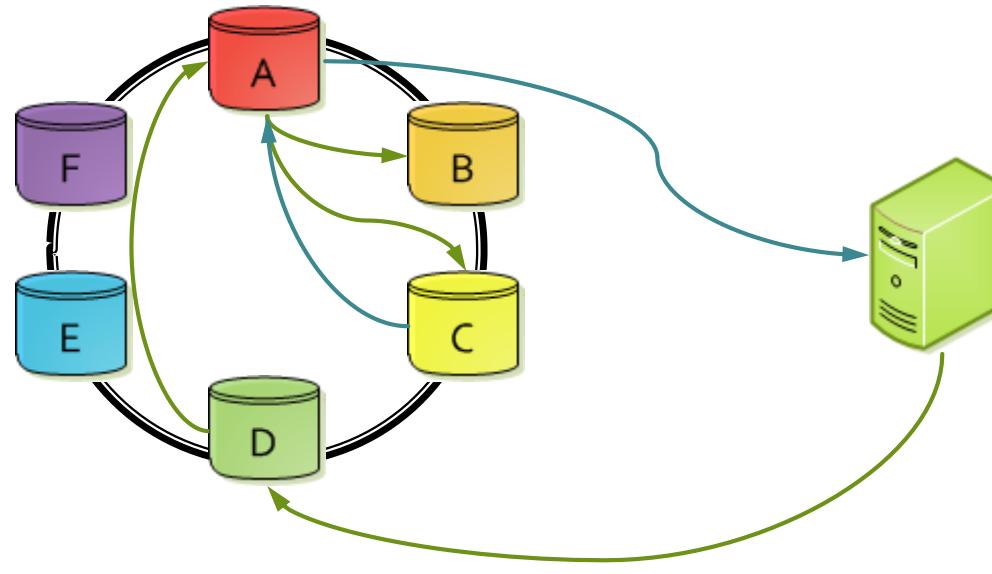
- ▶ Solution: **Consistent Hashing** – mapping of data to nodes is stable under topology changes



Reading

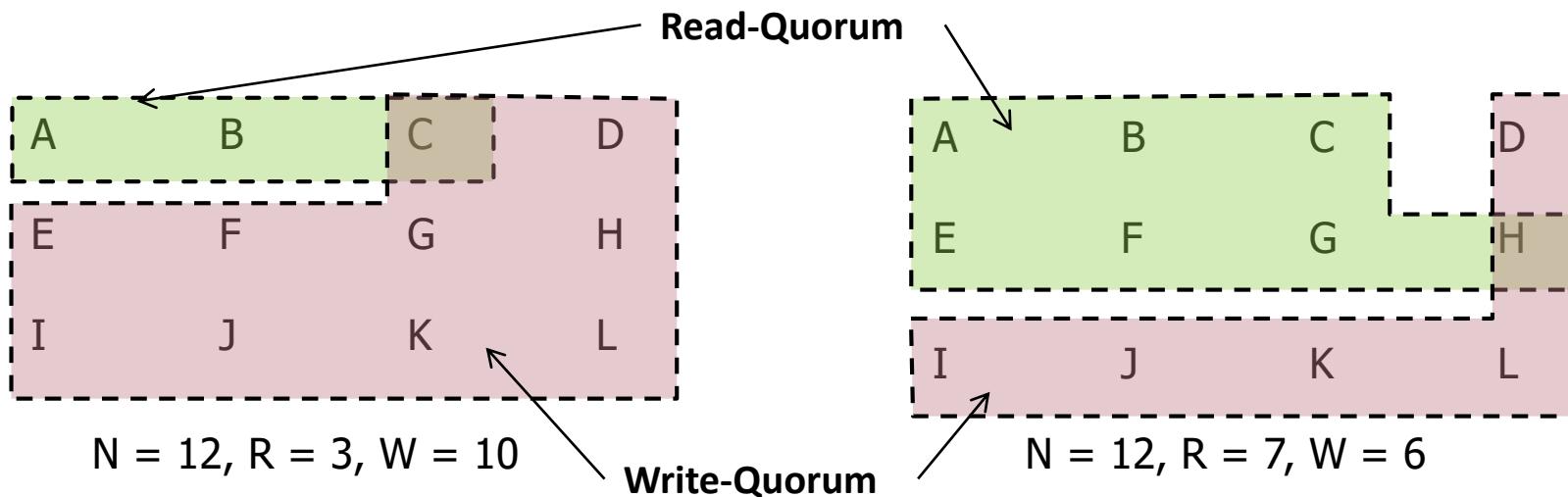
Parameters R, W, N

- ▶ An arbitrary node acts as a coordinator
- ▶ N: number of replicas
- ▶ R: number of nodes that need to confirm a read
- ▶ W: number of nodes that need to confirm a write



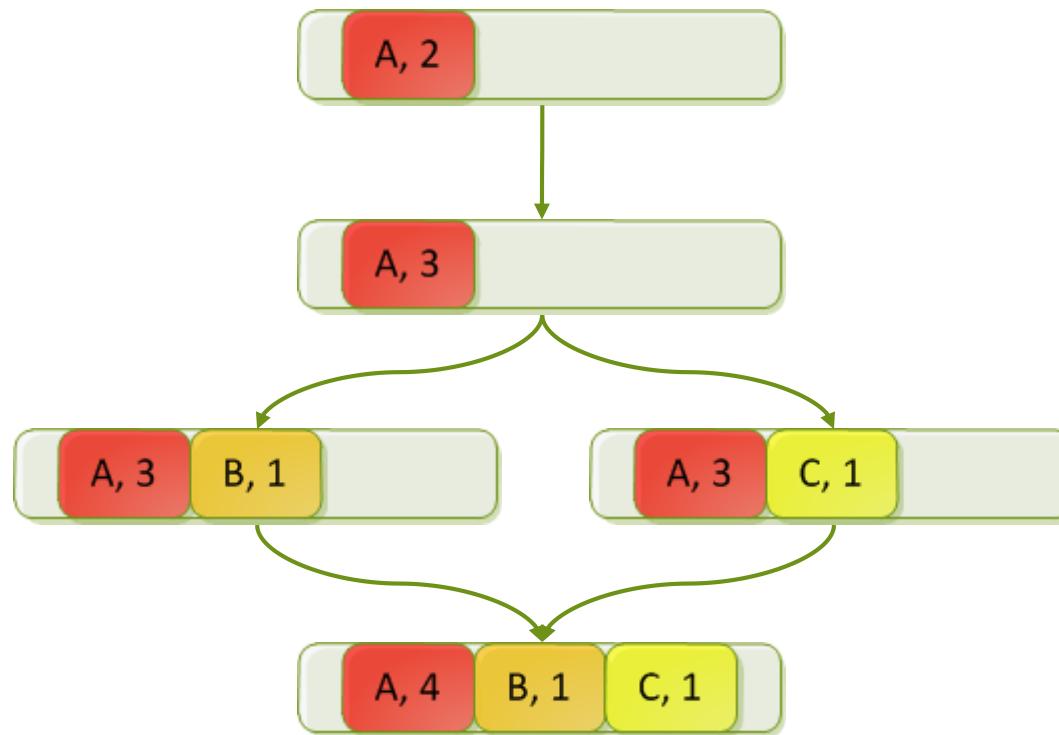
Quorums

- ▶ N (Replicas), W (Write Acks), R (Read Acks)
 - $R + W \leq N \Rightarrow$ No guarantee
 - $R + W > N \Rightarrow$ newest version included



Vector clocks

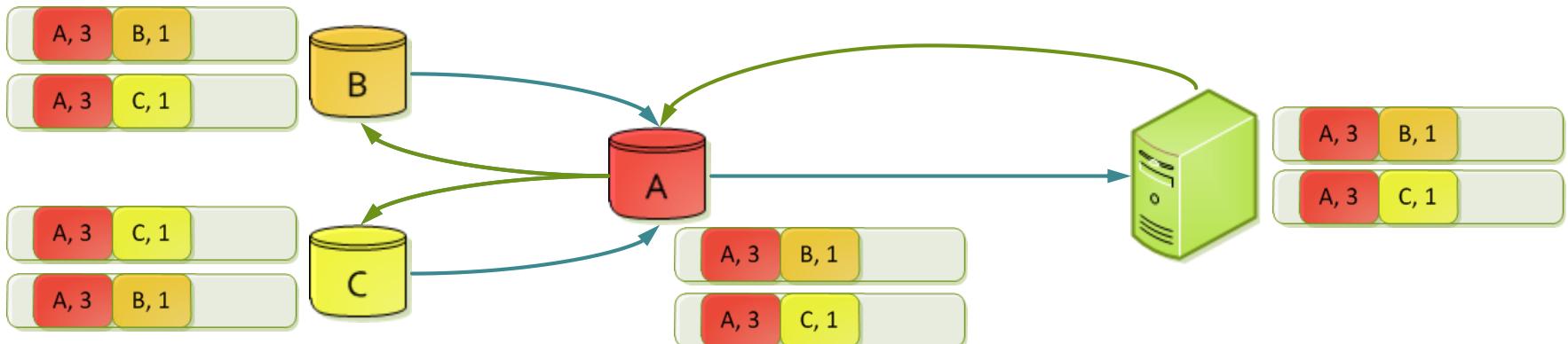
- ▶ Dynamo uses **Vector Clocks** for versioning



C. J. Fidge, Timestamps in message-passing systems
that preserve the partial ordering (1988)

Versioning and Consistency

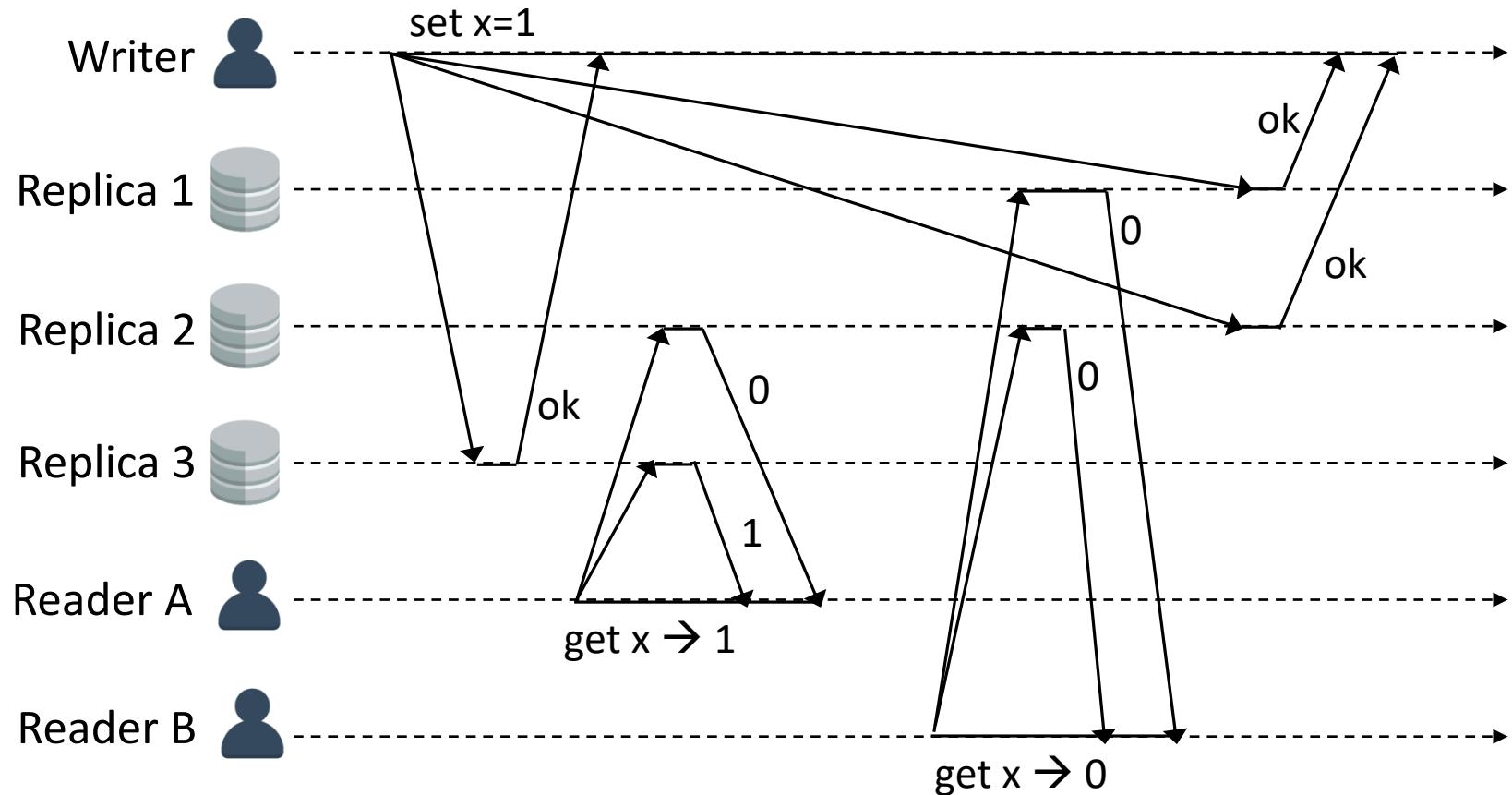
- ▶ $R + W \leq N \Rightarrow$ no consistency guarantee
- ▶ $R + W > N \Rightarrow$ newest acked value included in reads
- ▶ **Vector Clocks** used for versioning



Read Repair

$R + W > N$ does not imply linearizability

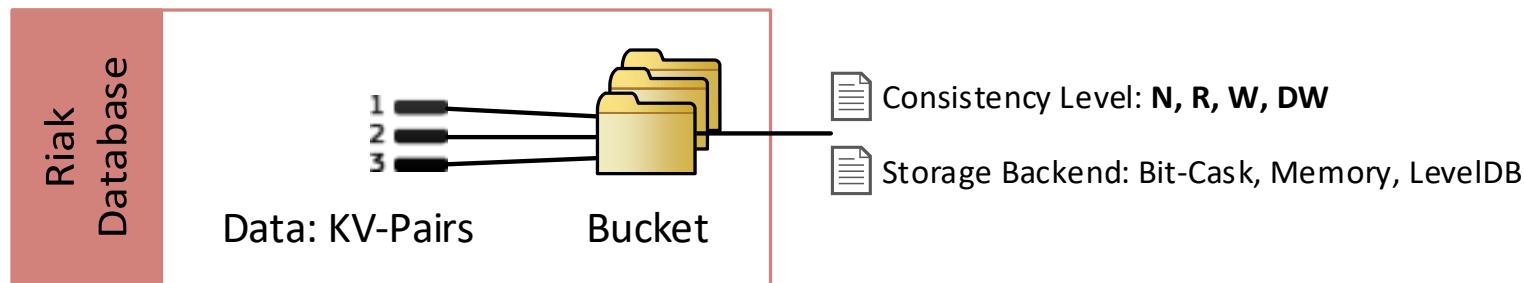
- Consider the following execution:



Riak (AP)

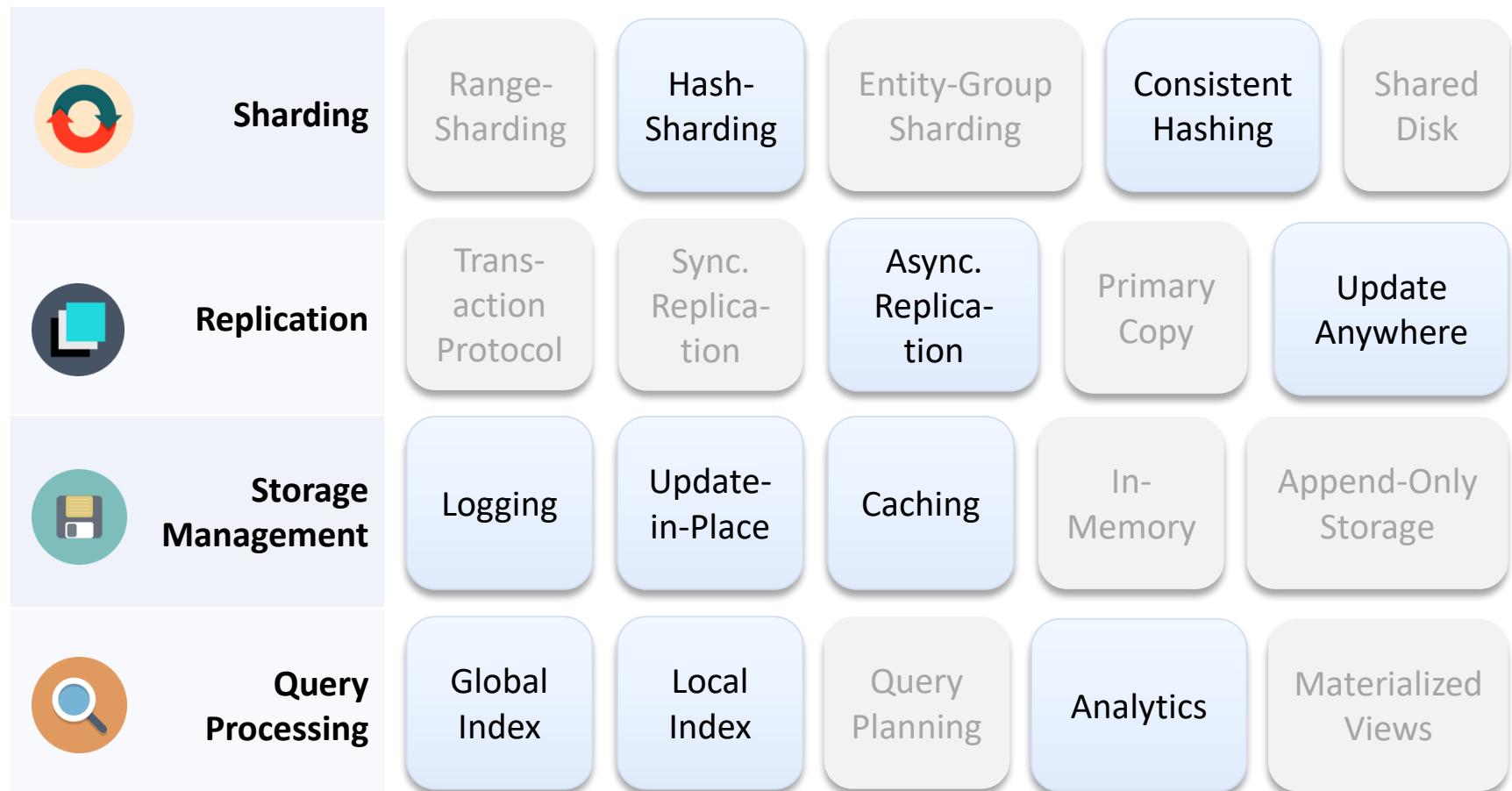
- ▶ Open-Source Dynamo-Implementation
- ▶ Extends Dynamo:
 - Keys are grouped to **Buckets**
 - KV-pairs may have **metadata** and **links**
 - Map-Reduce support
 - Secondary Indices, Update Hooks, Solr Integration
 - Option for **strongly consistent** buckets (experimental)
 - Riak CS: S3-like file storage, Riak TS: time-series database

Riak
Model:
Key-Value
License:
Apache 2
Written in:
Erlang und C



Dynamo and Riak

Classification





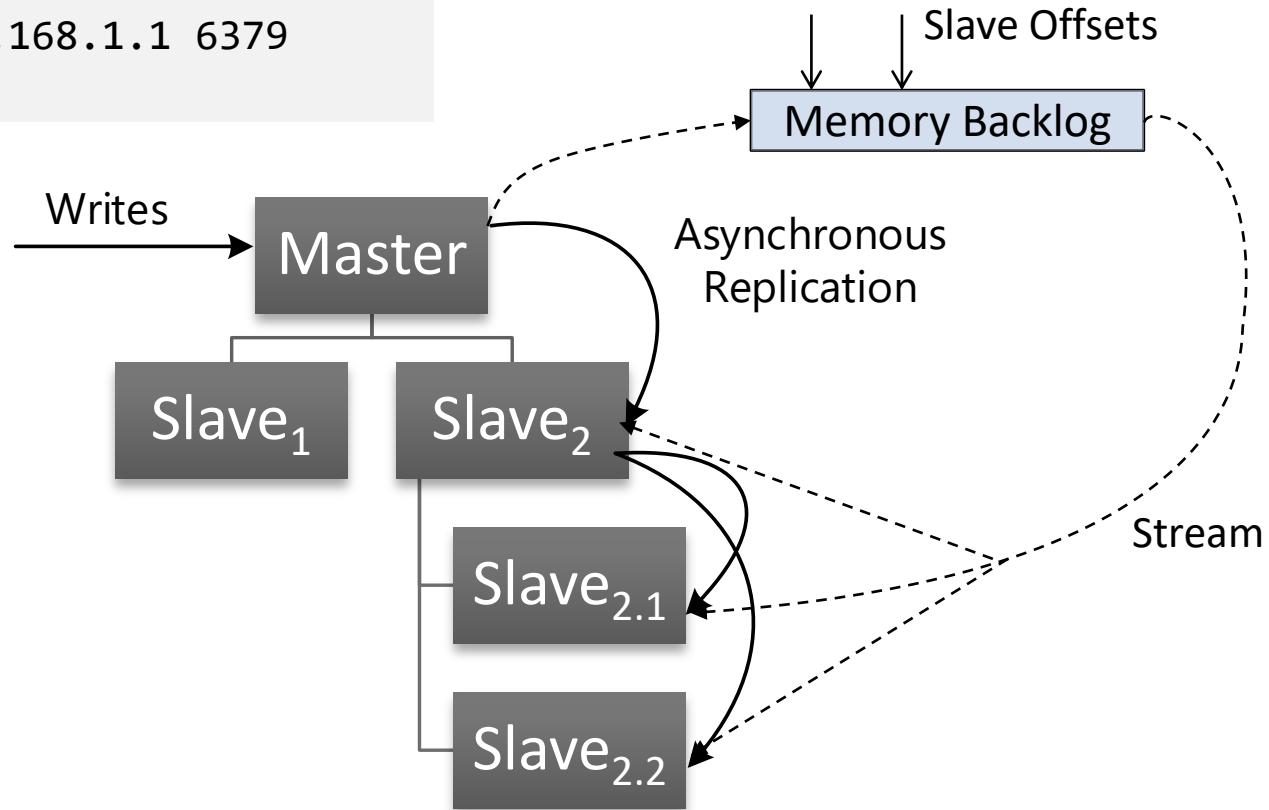
Redis (CA)

- ▶ Remote Dictionary Server
- ▶ In-Memory Key-Value Store
- ▶ Asynchronous Master-Slave Replication
- ▶ Data model: rich data structures stored under key
- ▶ Tunable persistence: logging and snapshots
- ▶ Single-threaded event-loop design (similar to Node.js)
- ▶ Optimistic batch transactions (*Multi blocks*)
- ▶ Very high performance: >100k ops/sec per node
- ▶ Redis Cluster adds sharding

Redis
Model:
Key-Value
License:
BSD
Written in:
C

Master-Slave Replication

```
> SLAVEOF 192.168.1.1 6379  
< +OK
```



Data structures

▶ String, List, Set, Hash, Sorted Set

String

web:index → "<html><head>..."

Set

users:2:friends → {23, 76, 233, 11}

List

users:2:inbox → [234, 3466, 86, 55]

Hash

users:2:settings → Theme → "dark", cookies → "false"

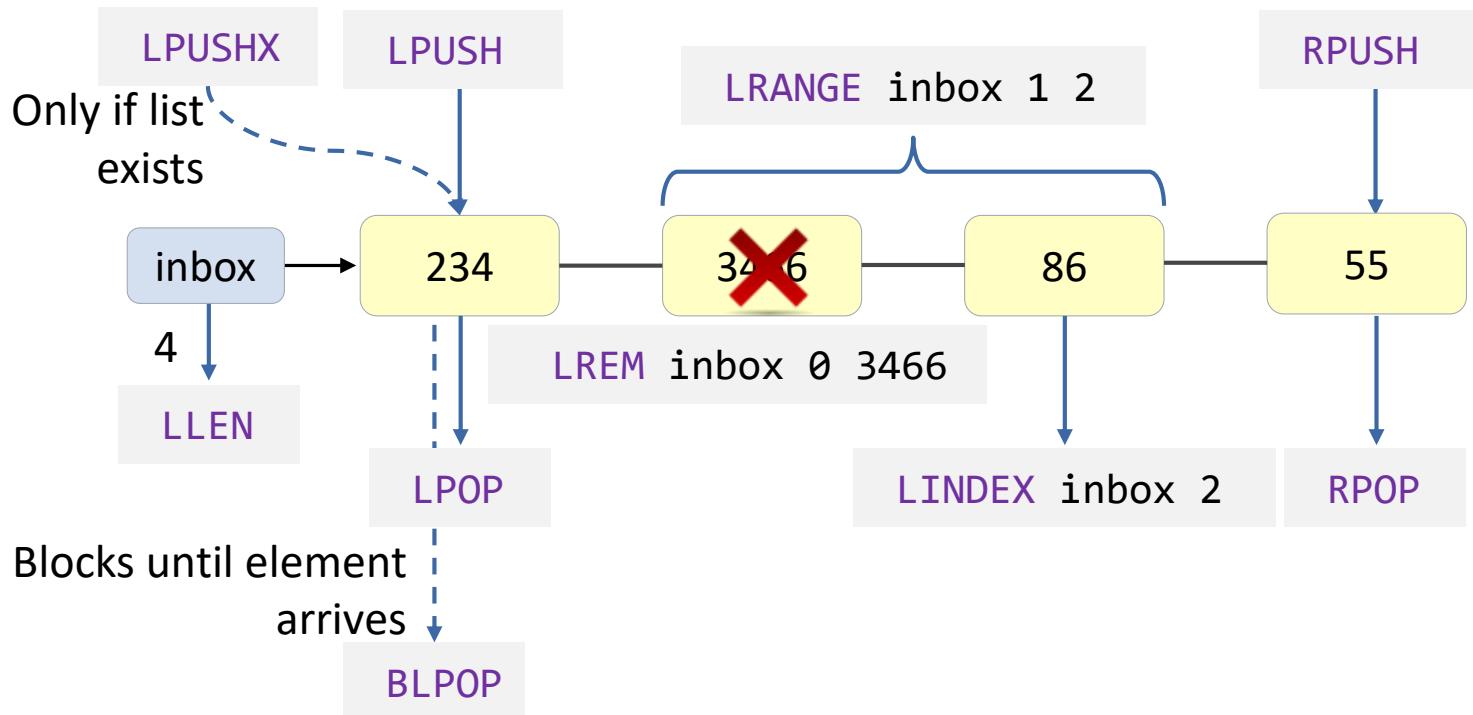
Sorted Set

top-posters → 466 → "2", 344 → "16"

Pub/Sub

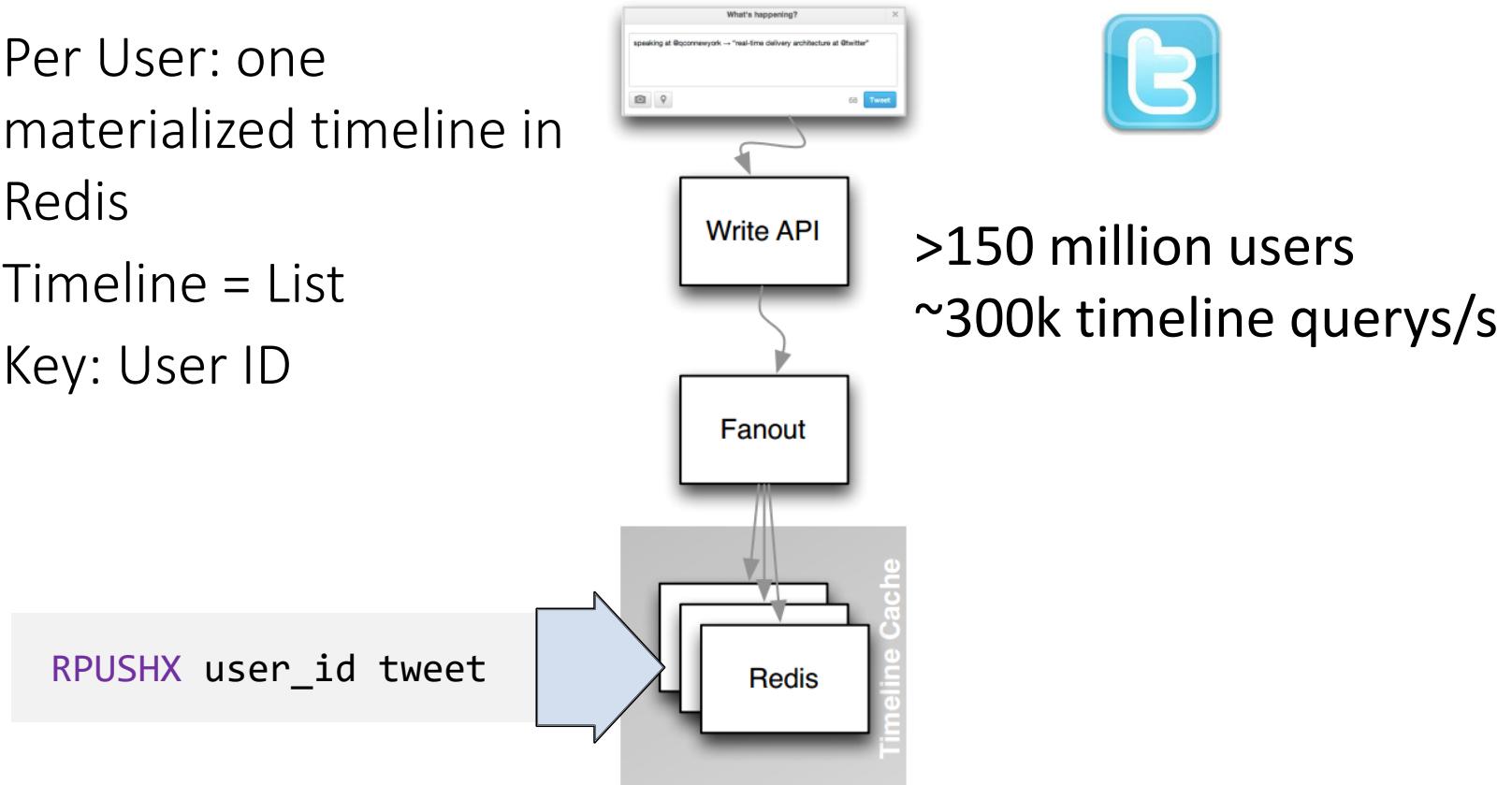
users:2:notifs → "{event: 'comment posted', time : ...}"

Example: List Data Structure



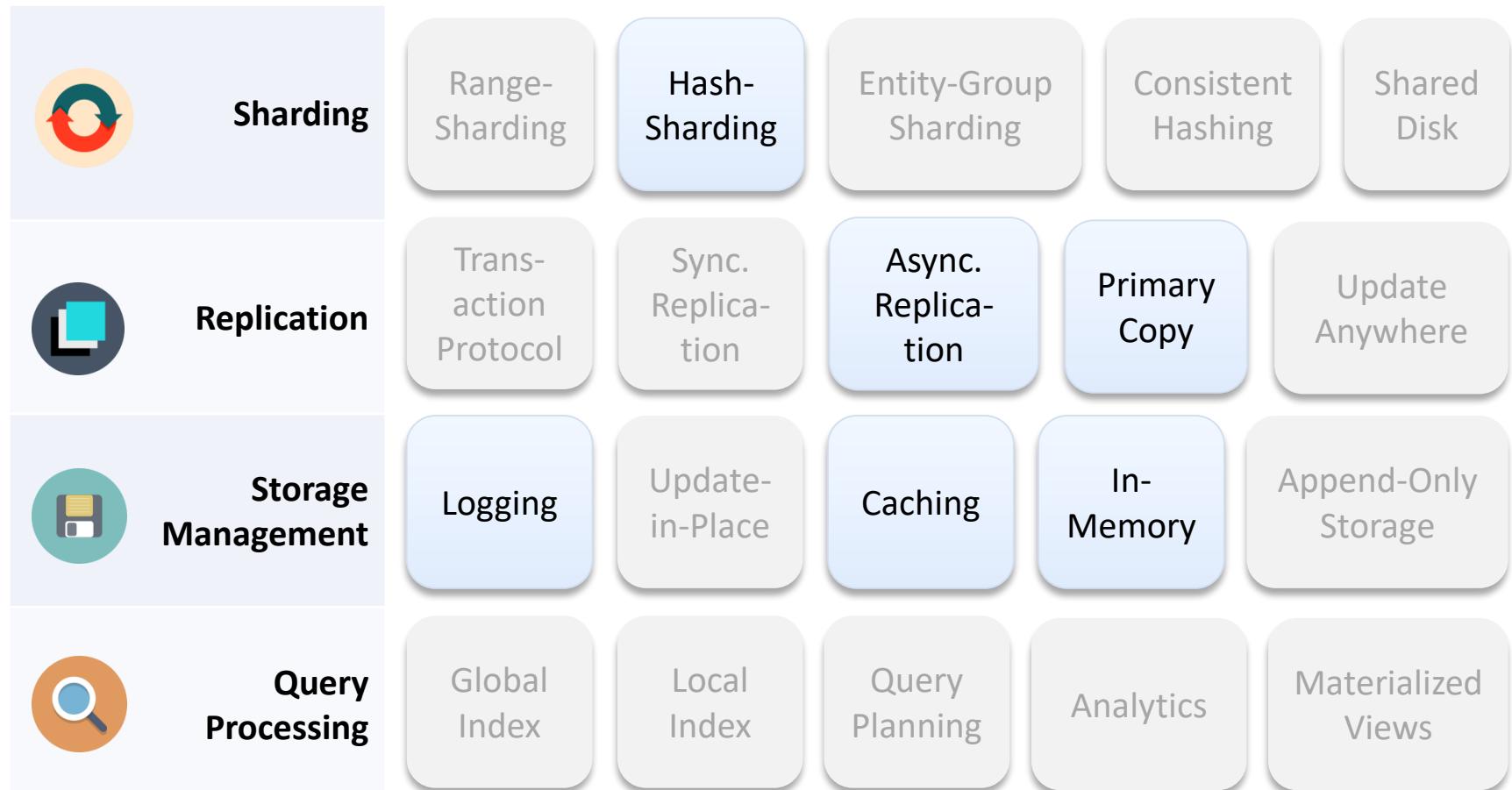
Example Redis Use-Case: Twitter

- ▶ Per User: one materialized timeline in Redis
- ▶ Timeline = List
- ▶ Key: User ID



Classification: Redis

Techniques



Google BigTable (CP)

- ▶ Published by Google in 2006
- ▶ Original purpose: storing the Google search index

A Bigtable is a sparse,
distributed, persistent
multidimensional sorted map.

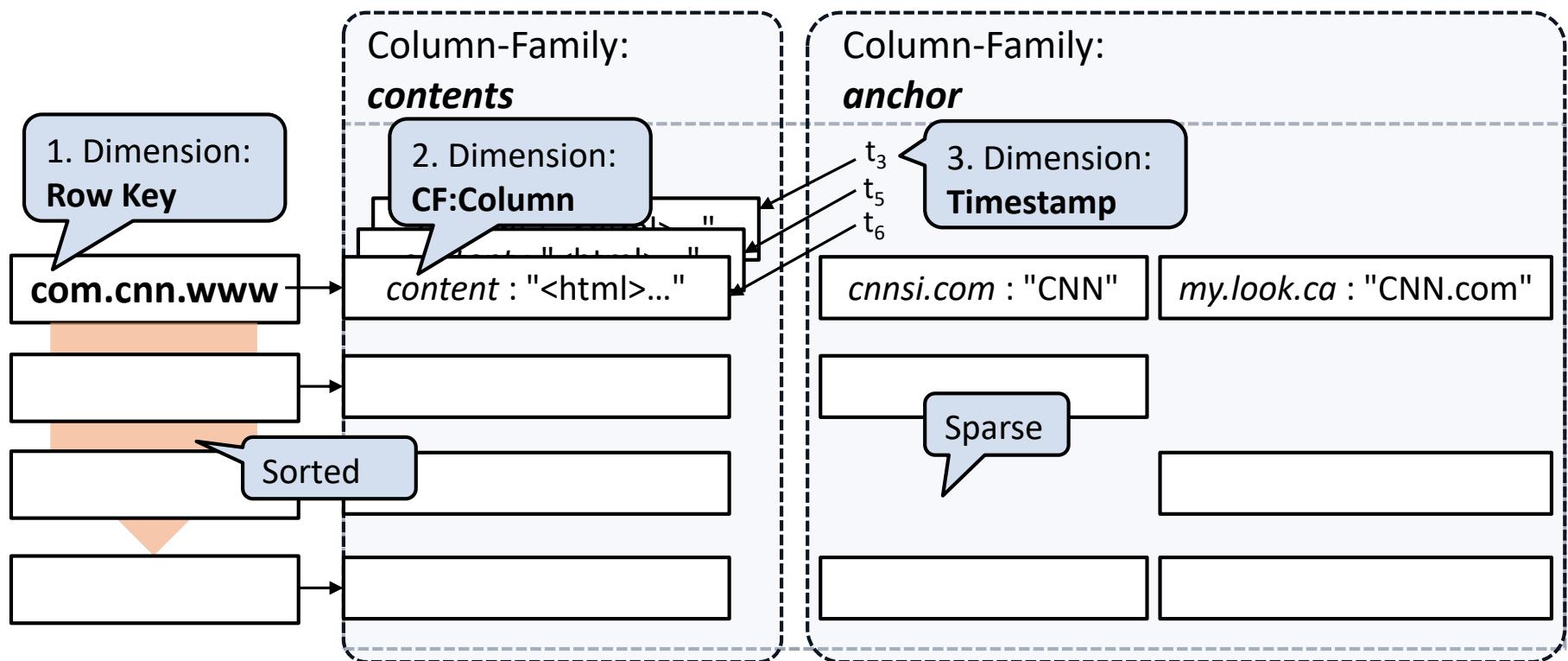
- ▶ Data model also used in: **HBase**, **Cassandra**, **HyperTable**, **Accumulo**



Chang, Fay, et al. "Bigtable: A distributed storage system for structured data."

Wide-Column Data Modelling

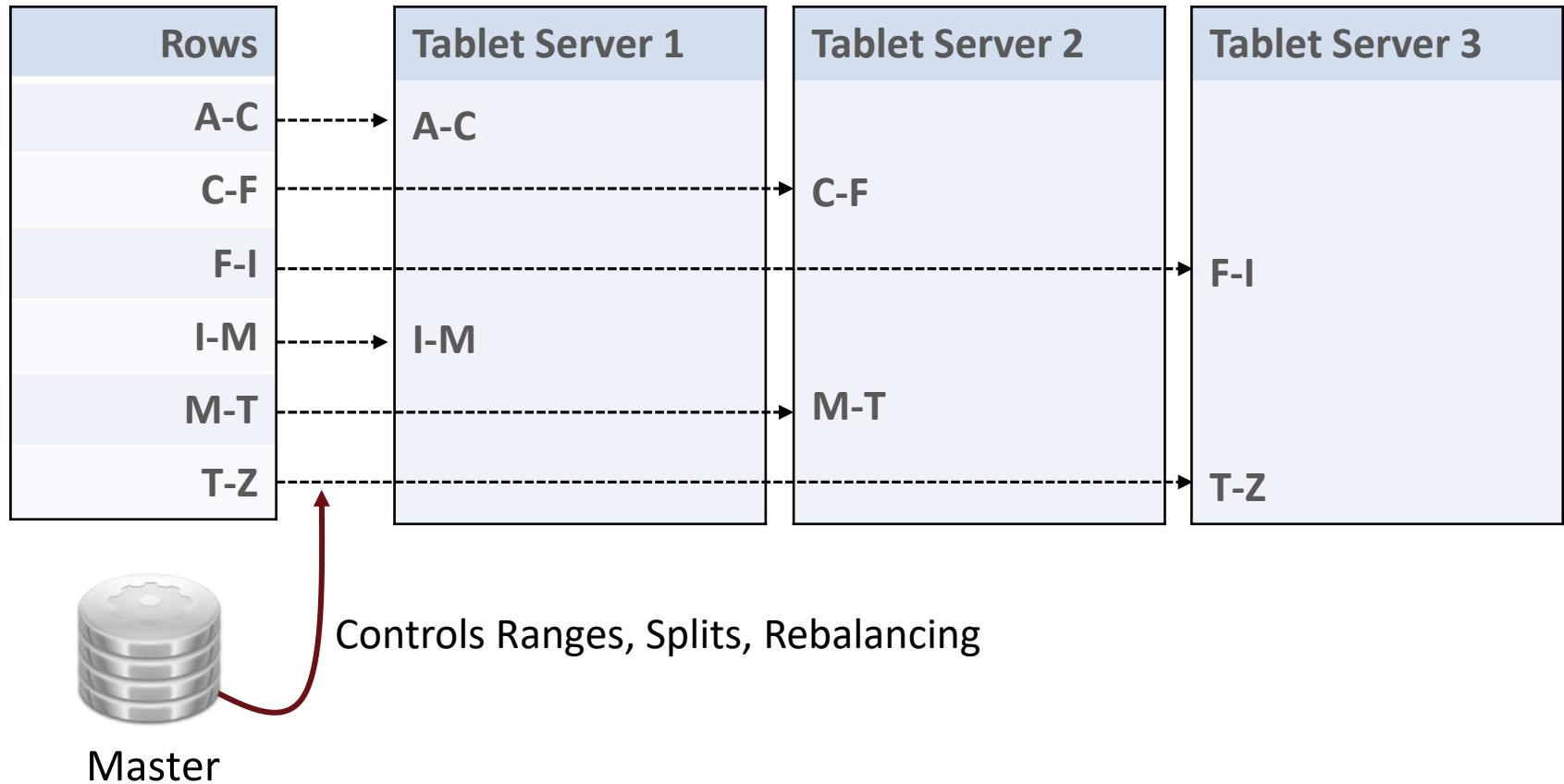
- ▶ Storage of crawled web-sites („Webtable“):



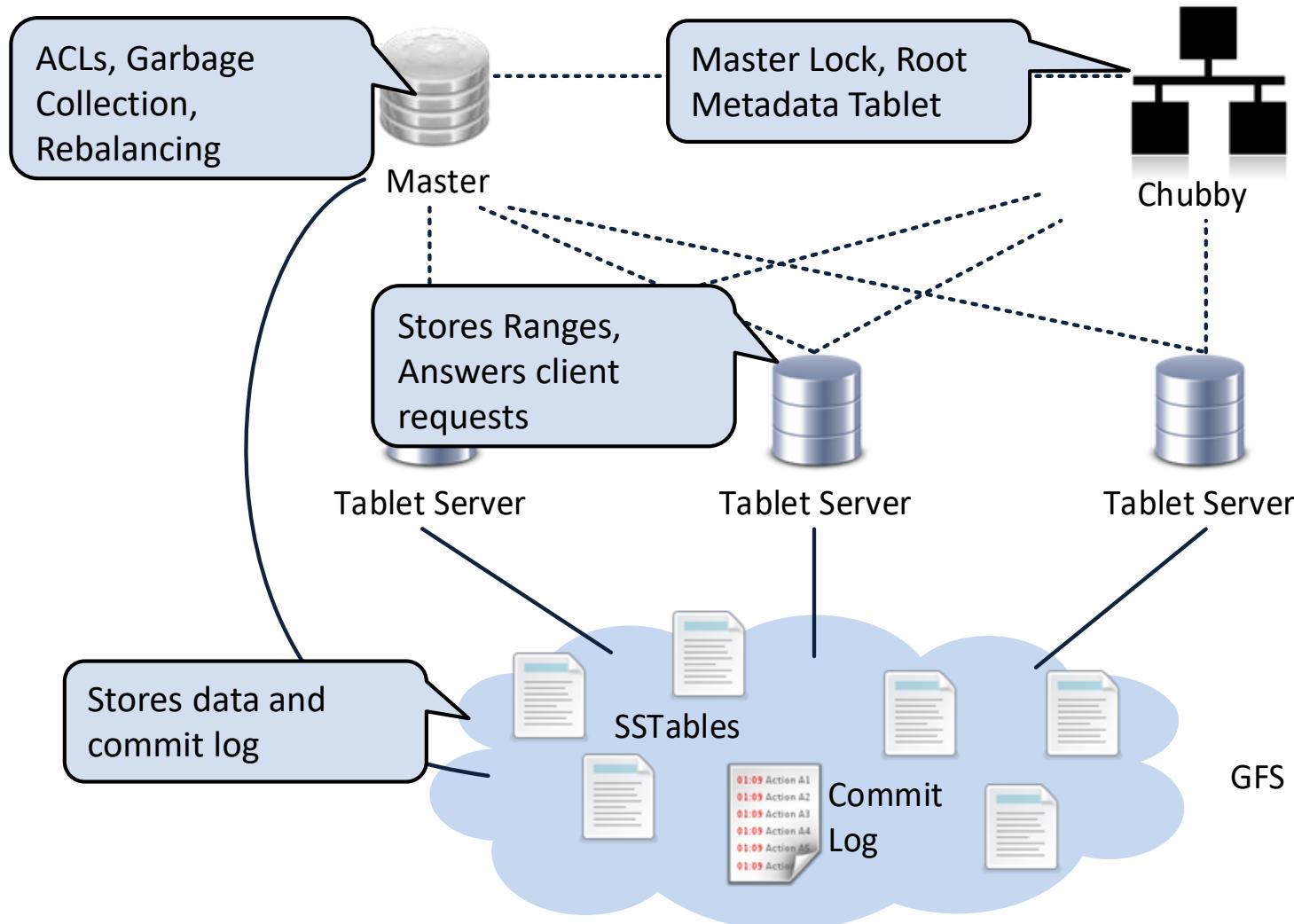
Range-based Sharding

BigTable Tablets

Tablet: Range partition of ordered records

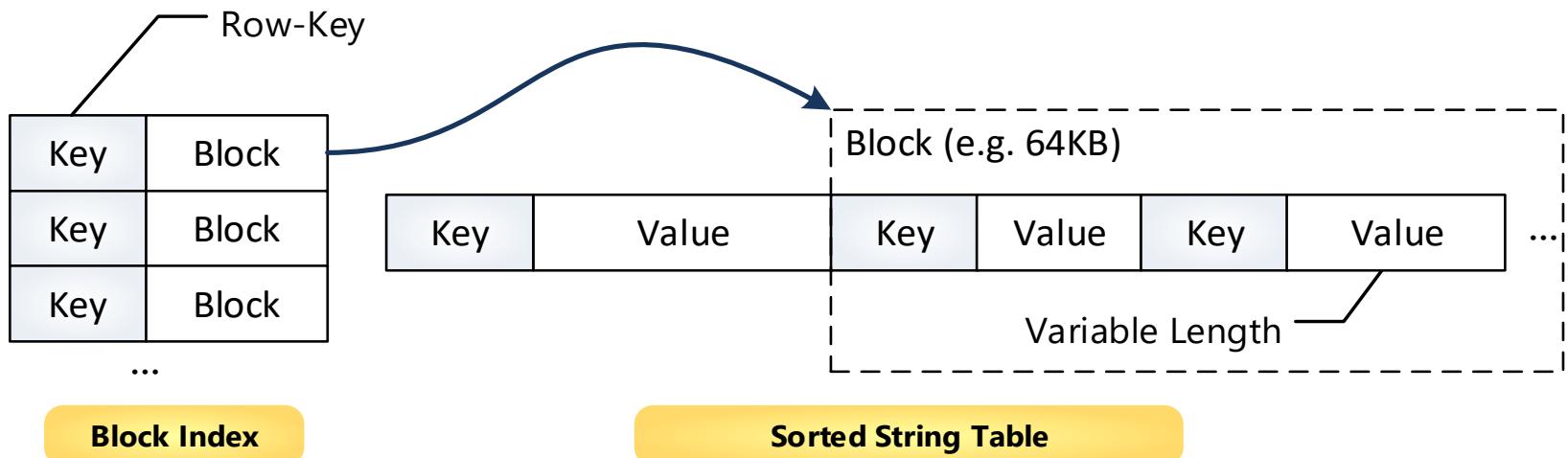


Architecture



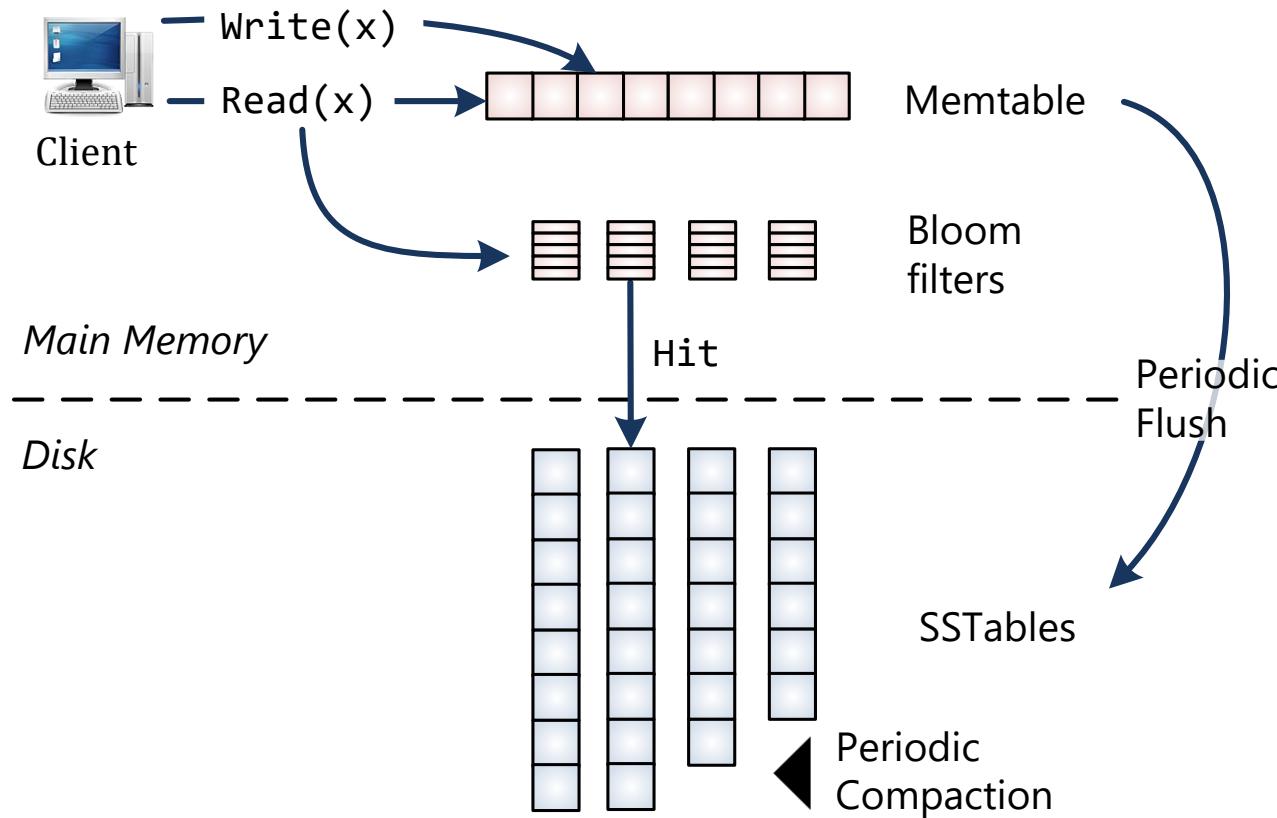
Storage: Sorted-String Tables

- ▶ **Goal:** Append-Only IO when writing (no disk seeks)
- ▶ Achieved through: **Log-Structured Merge Trees**
- ▶ **Writes** go to an in-memory *memtable* that is periodically persisted as an *SSTable* as well as a *commit log*
- ▶ **Reads** query memtable and all SSTables



Storage: Optimization

- ▶ Writes: Buffered in-memory in **Memtable**
- ▶ SSTable disk access optimized by Bloom filters



Apache HBase (CP)

- ▶ Open-Source Implementation of BigTable
- ▶ Hadoop-Integration
 - Data source for Map-Reduce
 - Uses Zookeeper and HDFS
- ▶ Data modelling challenges: key design, tall vs wide
 - **Row Key:** only access key (no indices) → key design important
 - **Tall:** good for scans
 - **Wide:** good for gets, consistent (*single-row atomicity*)
- ▶ No typing: application handles serialization
- ▶ Interface: REST, Avro, Thrift

HBase
Model:
Wide-Column
License:
Apache 2
Written in:
Java

Example: Facebook Insights



MD5(Reversed Domain) + Reversed Domain + URL-ID									Row Key
6PM Total	6PM Male	...	01.01 Total	01.01 Male	...	Total	Male	...	
10	7		100	65		567	567		

CF:Daily

CF:Monthly

CF>All

Atomic HBase Counter

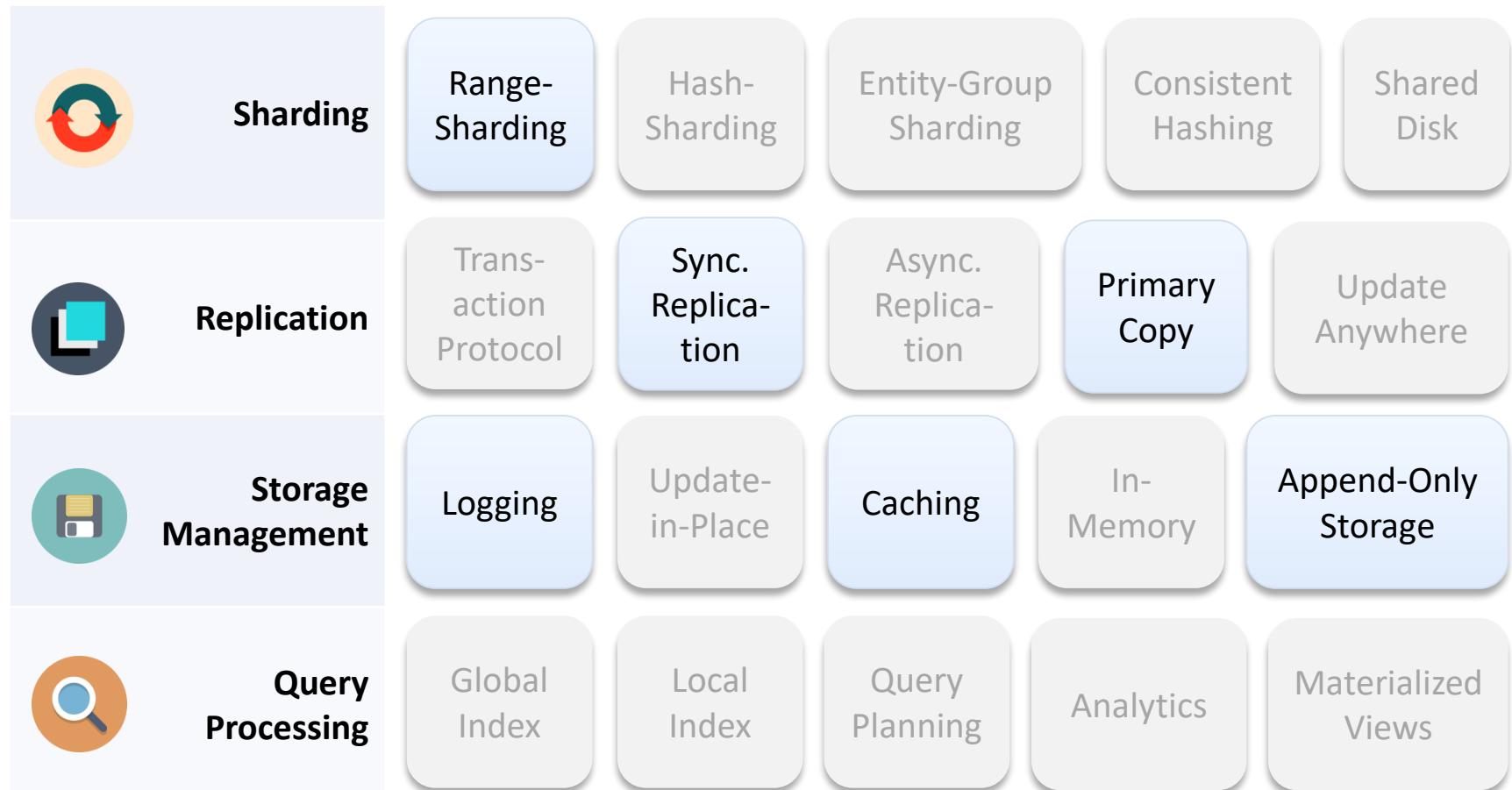
TTL – automatic deletion of old rows



Lars George: "Advanced HBase Schema Design"

Classification: HBase

Techniques



Apache Cassandra (AP)

- ▶ Published 2007 by Facebook
- ▶ **Idea:**
 - BigTable's wide-column data model
 - Dynamo ring for replication and sharding
- ▶ Cassandra Query Language (CQL): SQL-like query- and DDL-language
- ▶ **Compound indices:** *partition key* (shard key) + *clustering key* (ordered per partition key) → Limited range queries

Cassandra

Model:

Wide-Column

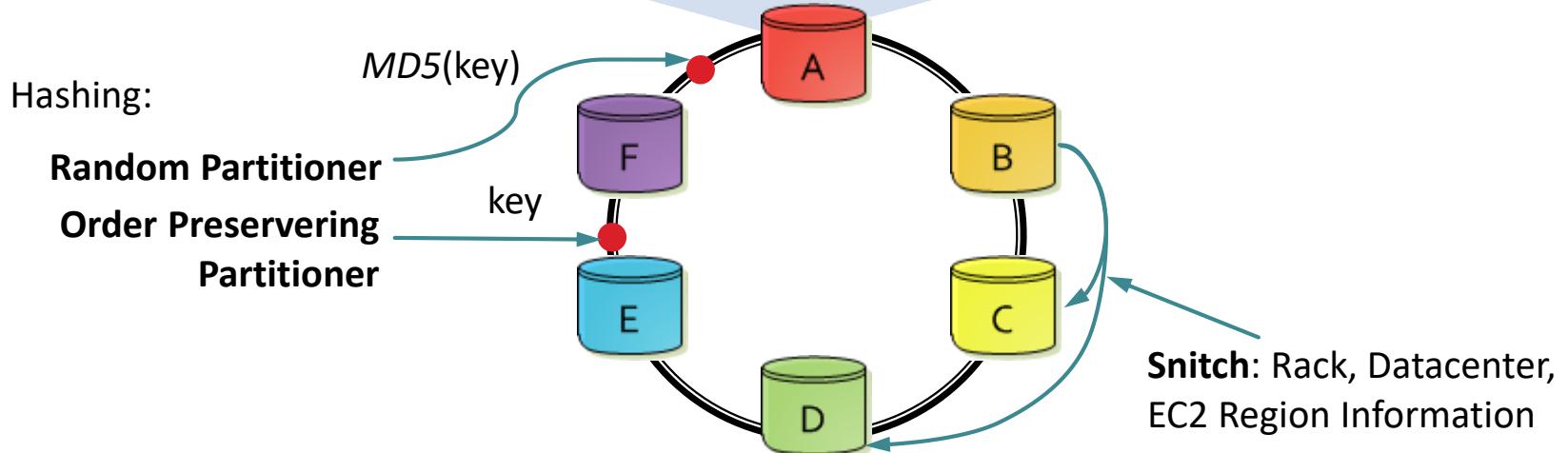
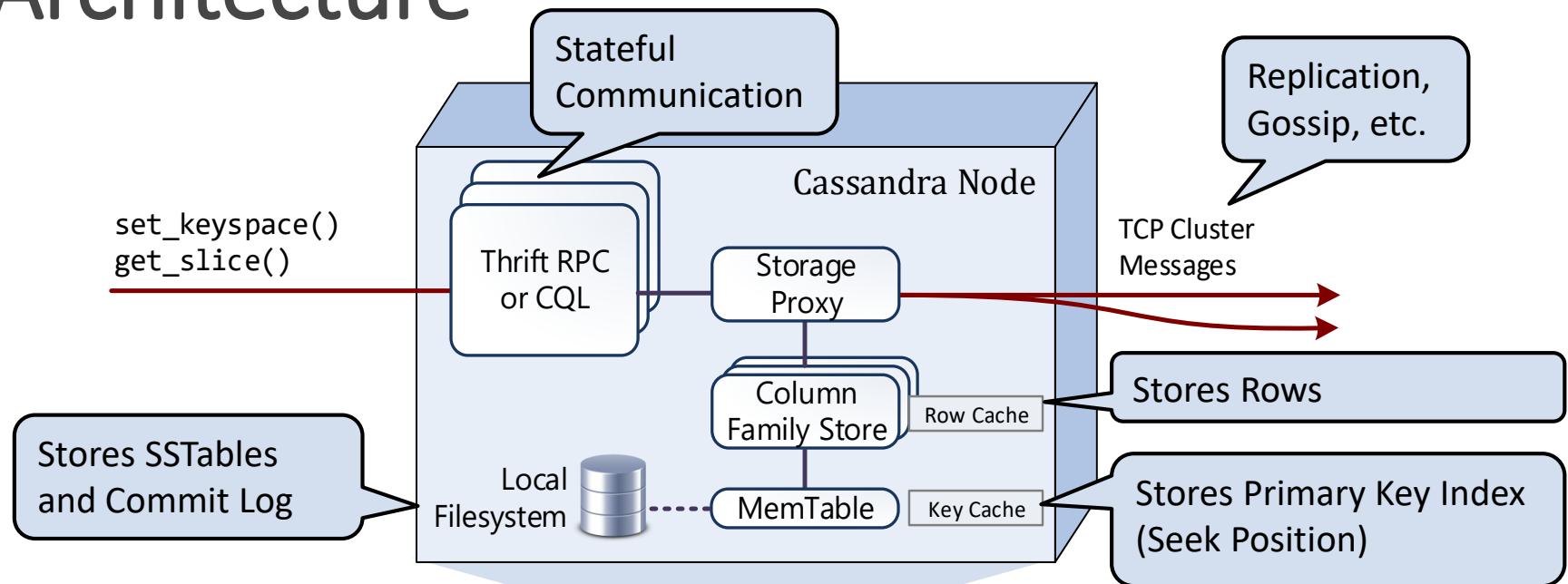
License:

Apache 2

Written in:

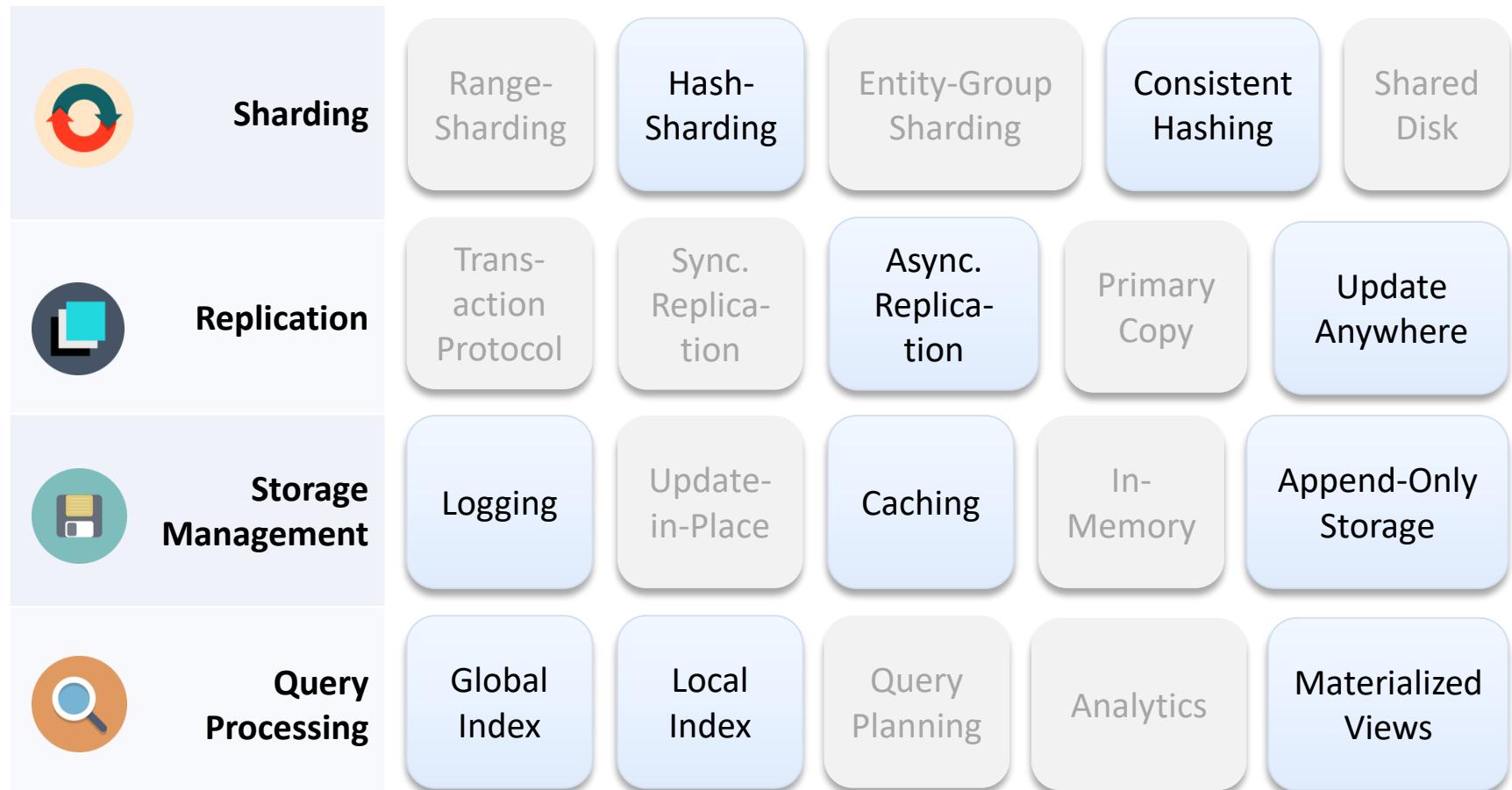
Java

Architecture



Classification: Cassandra

Techniques



MongoDB (CP)

- ▶ From **humongous** \cong gigantic
- ▶ Schema-free document database with tunable consistency
- ▶ Allows complex queries and indexing
- ▶ **Sharding** (either range- or hash-based)
- ▶ **Replication** (either synchronous or asynchronous)
- ▶ Storage Management:
 - **Write-ahead logging** for redos (*journaling*)
 - **Storage Engines**: memory-mapped files, in-memory, Log-structured merge trees (WiredTiger), ...

MongoDB

Model:

Document

License:

GNU AGPL 3.0

Written in:

C++

Data Modelling

```
{  
    "_id" : ObjectId("51a5d316d70beffe74ecc940")  
    title : "Iron Man 3",  
    year : 2013,  
    rating : 7.6,  
    director: "Shane Block",  
    genre : [ "Action",  
              "Adventure",  
              "Sci -Fi"],  
    actors : ["Downey Jr., Robert",  
              "Paltrow , Gwyneth"],  
    tweets : [ {  
        "text": "user" : "Franz Kafka",  
        "coordinates": "text" : "#nowwatching Iron Man 3",  
        "retweet" : false,  
        "date" : ISODate("2013-05-29T13:15:51Z")  
    }]  
}
```

Movie Document

Genre

Actor

Tweet

Denormalisation instead of joins

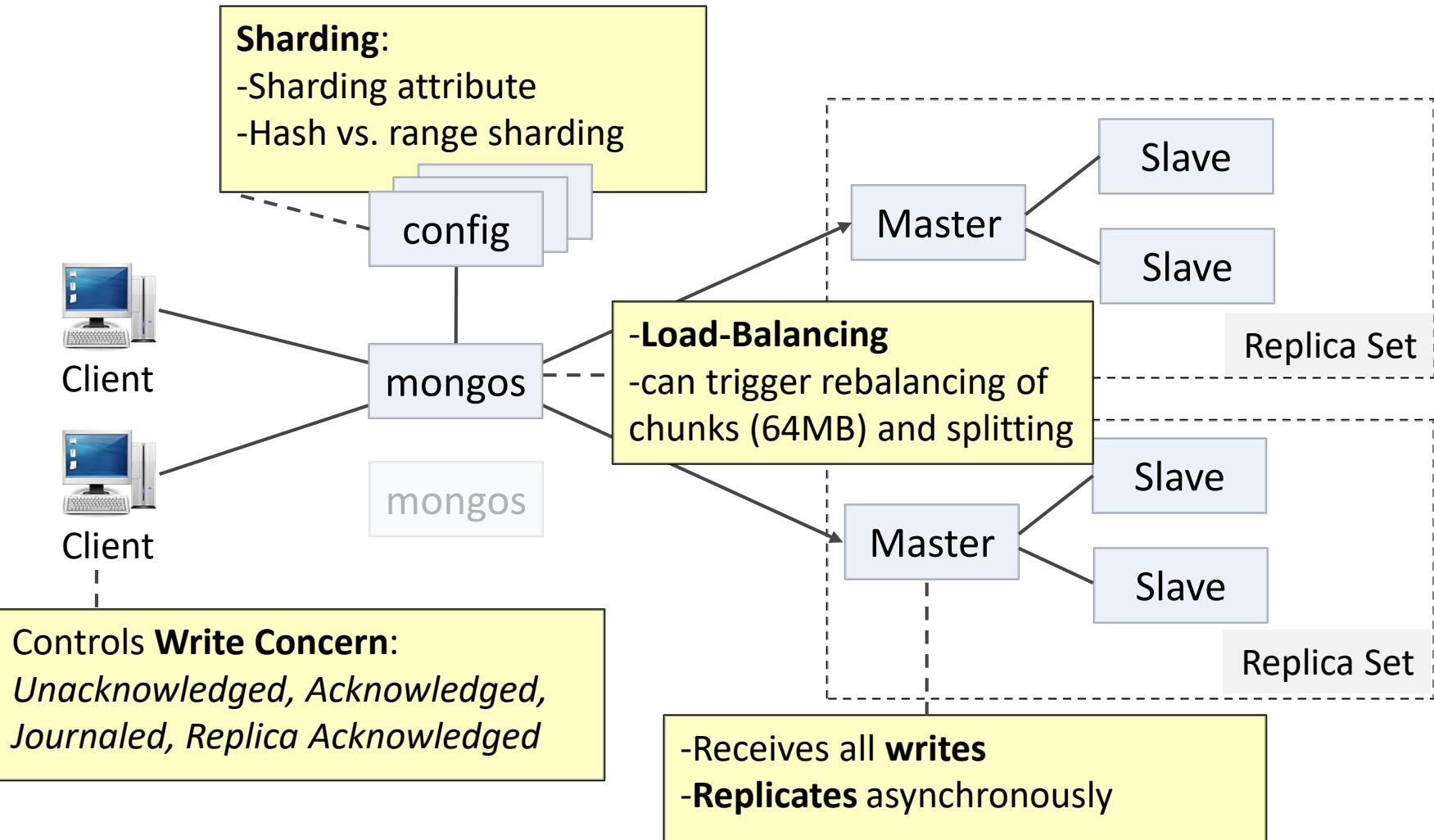
Nesting replaces 1:n and 1:1 relations

Schemafreeness:
Attributes per document

Unit of atomicity:
document

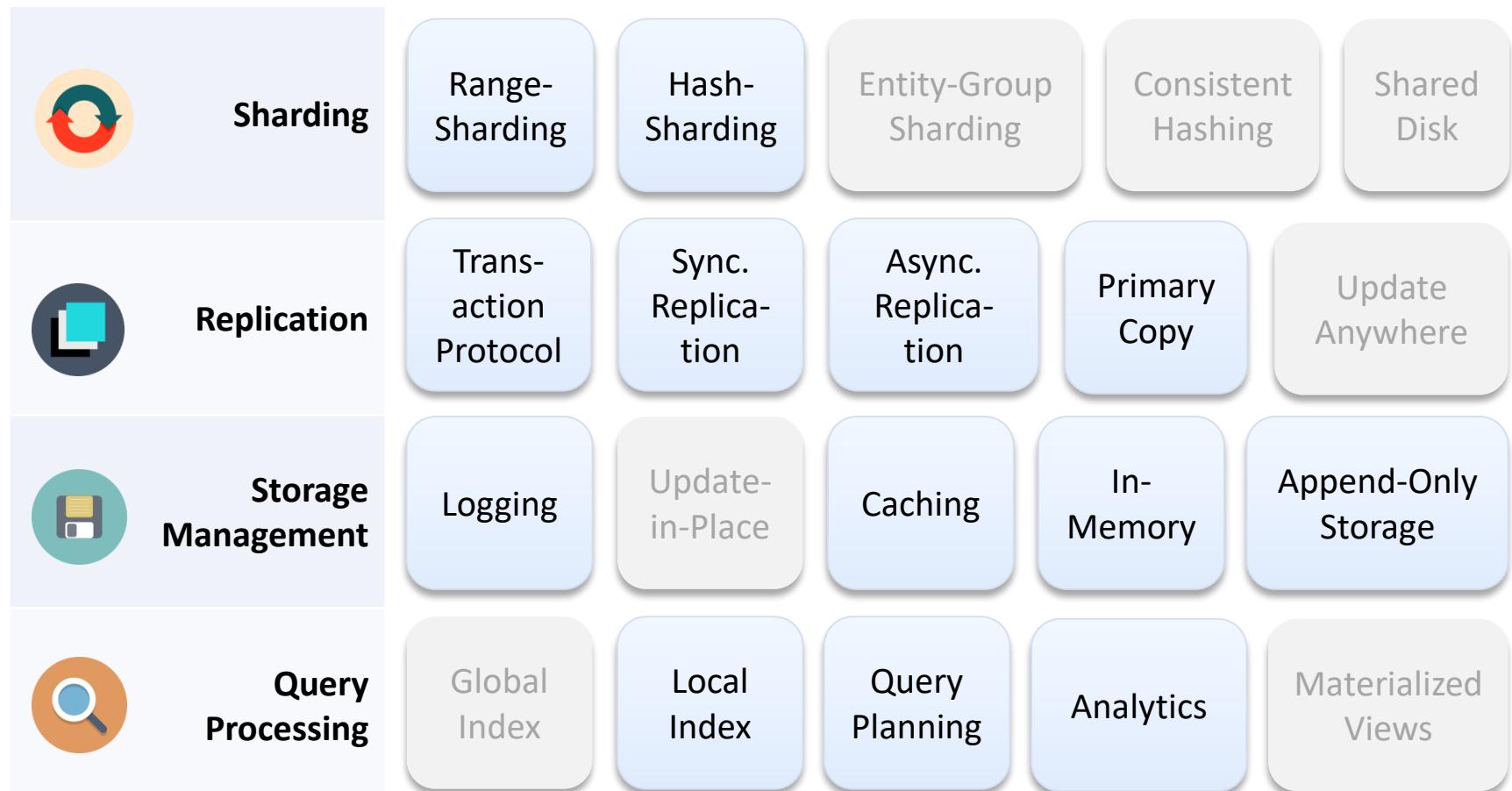
Principles

Sharding und Replication



Classification: MongoDB

Techniques



Other Systems

Graph databases

- ▶ **Neo4j** (ACID, replicated, Query-language)
- ▶ **HypergraphDB** (directed Hypergraph, BerkleyDB-based)
- ▶ **Titan** (distributed, Cassandra-based)
- ▶ **ArangoDB, OrientDB** („multi-model“)
- ▶ **SparkleDB** (RDF-Store, SPARQL)
- ▶ **InfinityDB** (embeddable)
- ▶ **InfiniteGraph** (distributed, low-level API, Objectivity-based)

Other Systems

Key-Value Stores

- ▶ **Aerospike** (SSD-optimized)
- ▶ **Voldemort** (Dynamo-style)
- ▶ **Memcache** (in-memory cache)
- ▶ **LevelDB** (embeddable, LSM-based)
- ▶ **RocksDB** (LevelDB-Fork with Transactions and Column Families)
- ▶ **HyperDex** (Searchable, Hyperspace-Hashing, Transactions)
- ▶ **Oracle NoSQL database** (distributed frontend for BerkleyDB)
- ▶ **HazelCast** (in-memory data-grid based on Java Collections)
- ▶ **FoundationDB** (ACID through Paxos)

Other Systems

Document Stores

- ▶ **CouchDB** (Multi-Master, lazy synchronization)
- ▶ **CouchBase** (distributed Memcache, N1QL~SQL, MR-Views)
- ▶ **RavenDB** (single node, SI transactions)
- ▶ **RethinkDB** (distributed CP, MVCC, joins, aggregates, real-time)
- ▶ **MarkLogic** (XML, distributed 2PC-ACID)
- ▶ **ElasticSearch** (full-text search, scalable, unclear consistency)
- ▶ **Solr** (full-text search)
- ▶ **Azure DocumentDB** (cloud-only, ACID, WAS-based)

Other Systems

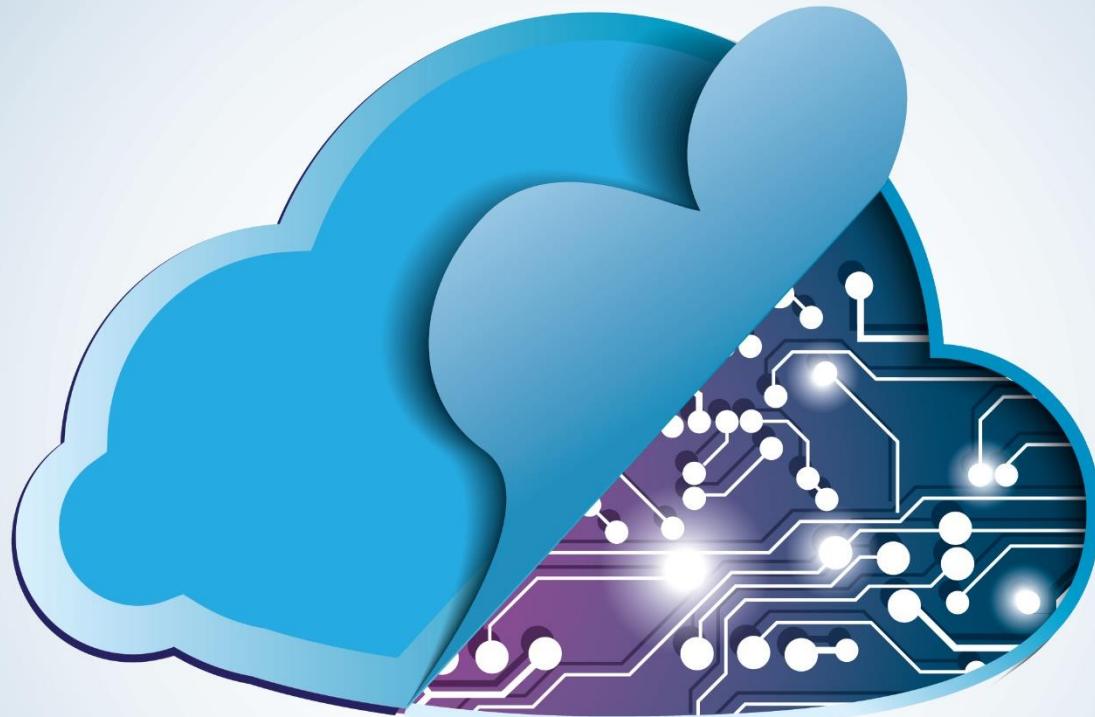
Wide-Column Stores

- ▶ **Accumulo** (BigTable-style, cell-level security)
- ▶ **HyperTable** (BigTable-style, written in C++)

Other Systems

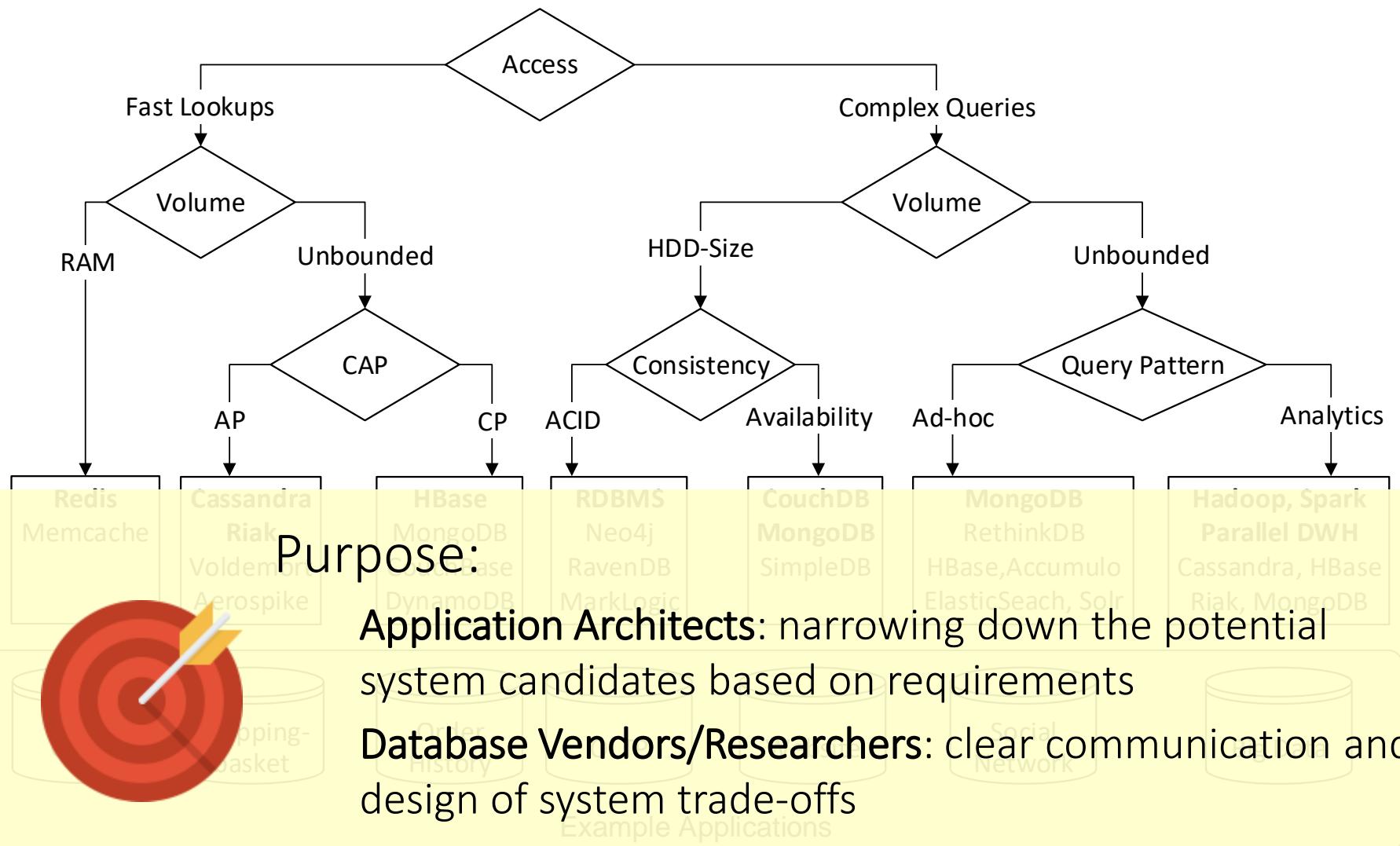
NewSQL Systems

- ▶ **CockroachDB** (Spanner-like, SQL, no joins, transactions)
- ▶ **Crate** (ElasticSearch-based, SQL, no transaction guarantees)
- ▶ **VoltDB** (HStore, ACID, in-memory, uses stored procedures)
- ▶ **Calvin** (log- & Paxos-based ACID transactions)
- ▶ **FaunaDB** (based on Calvin design, by Twitter engineers)
- ▶ **Google F1** (based on Spanner, SQL)
- ▶ **Google Cloud Spanner** (Improved F1 as a service)
- ▶ **Microsoft Cloud SQL Server** (distributed CP, MSSQL-comp.)
- ▶ **MySQL Cluster, Galera Cluster, Percona XtraDB Cluster**
(distributed storage engine for MySQL)



How can the choices for an appropriate system be narrowed down?

NoSQL Decision Tree



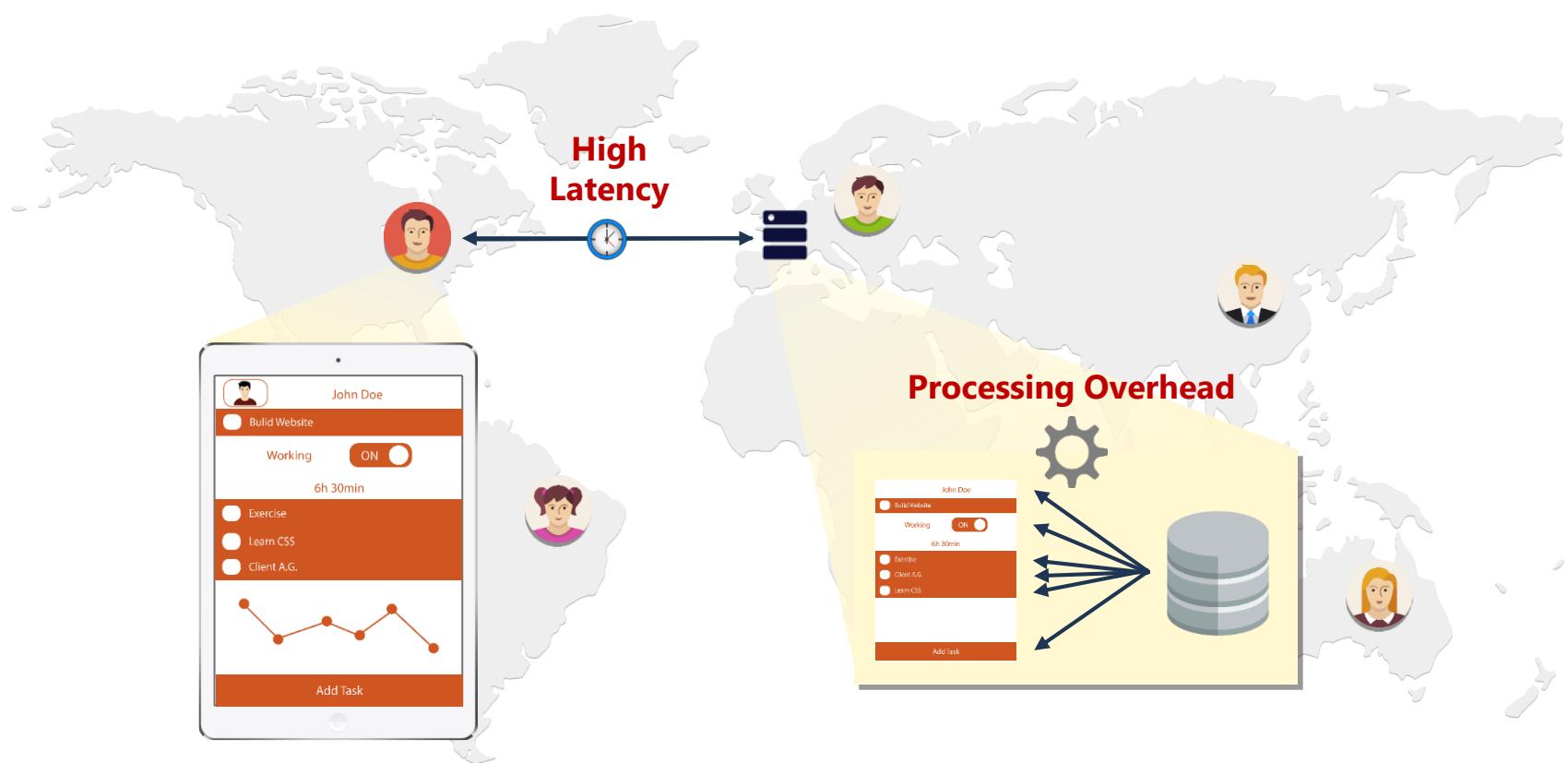
A photograph of a woman with long dark hair, seen from behind, looking out over a body of water towards a city skyline at sunset. The sky is a warm orange and yellow. In the background, several large industrial port cranes are silhouetted against the bright horizon, and the water reflects the colors of the sky.

OUTLOOK

Baqend Research at the University of Hamburg

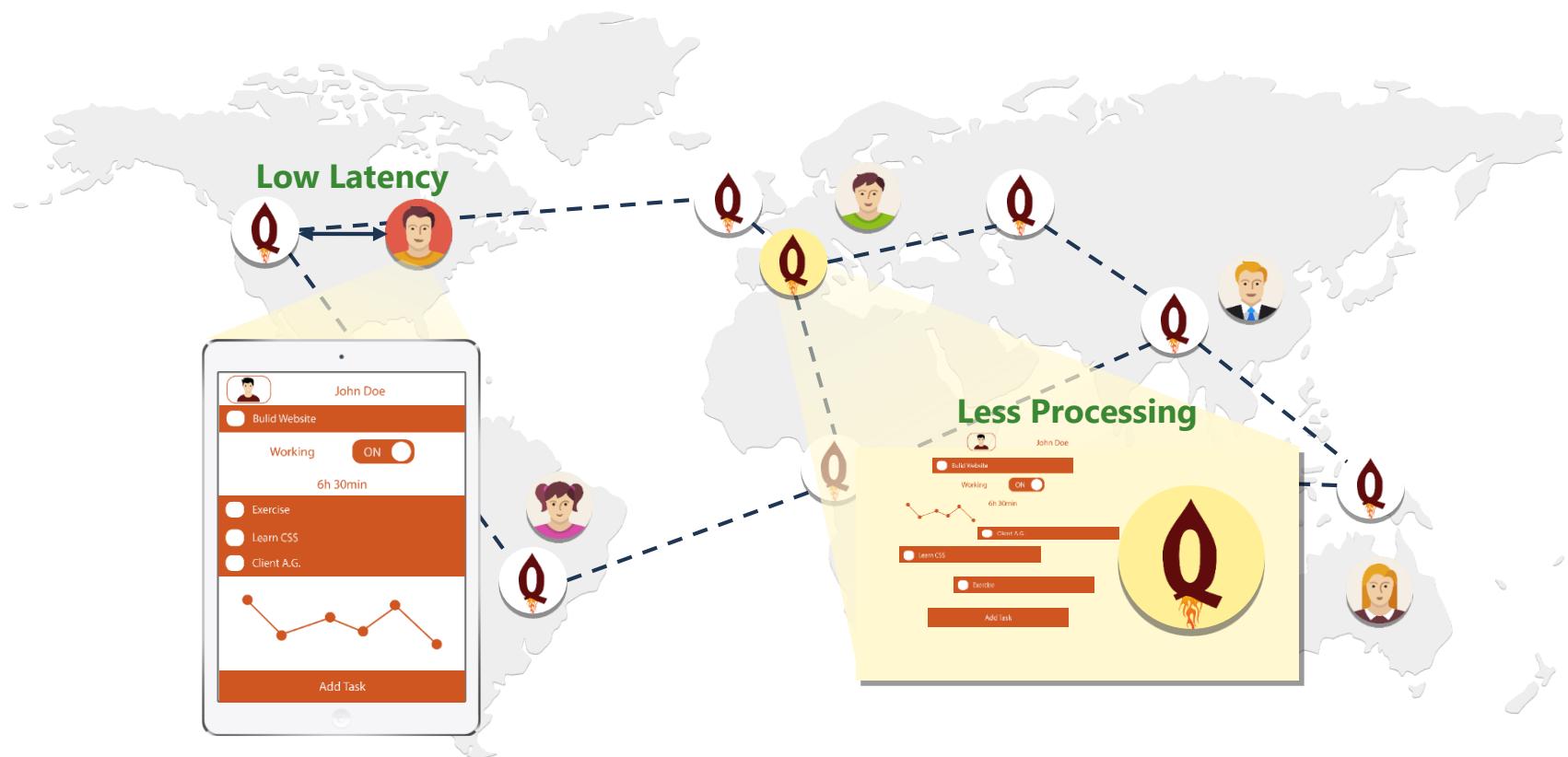
Problem: Slow Websites

Two Bottlenecks: Latency and Processing



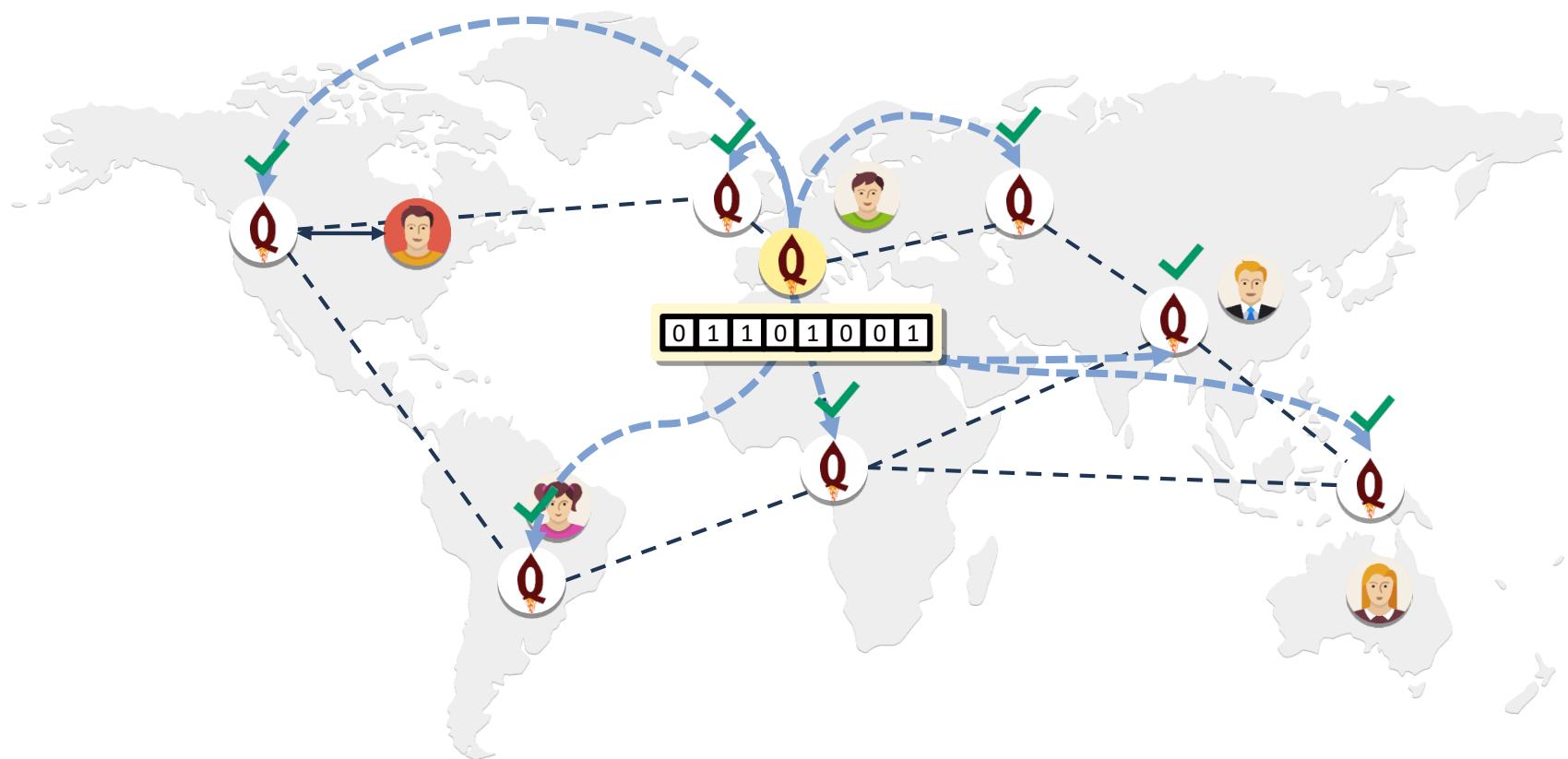
Solution: Global Caching

Fresh Data From Distributed Web Caches



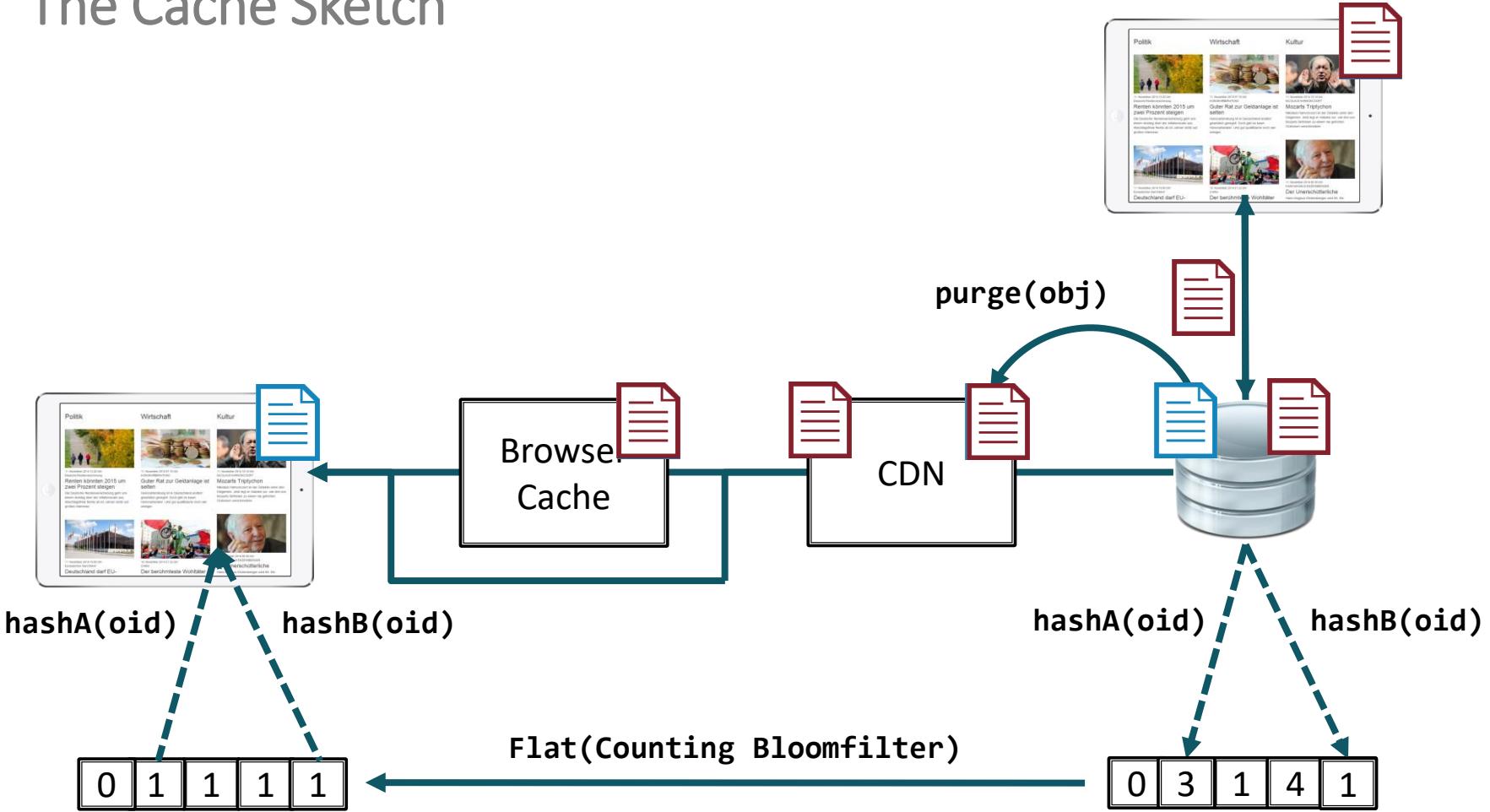
New Caching Algorithms

Solve Consistency Problem



Consistent Web Caching

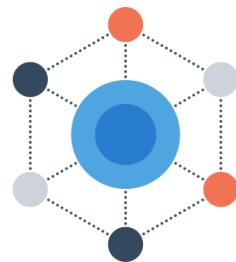
The Cache Sketch



Baqend

Try It Out!

Platform



- Platform for building
(Progressive) Web Apps
- **15x Performance Edge**
- **Faster Development**

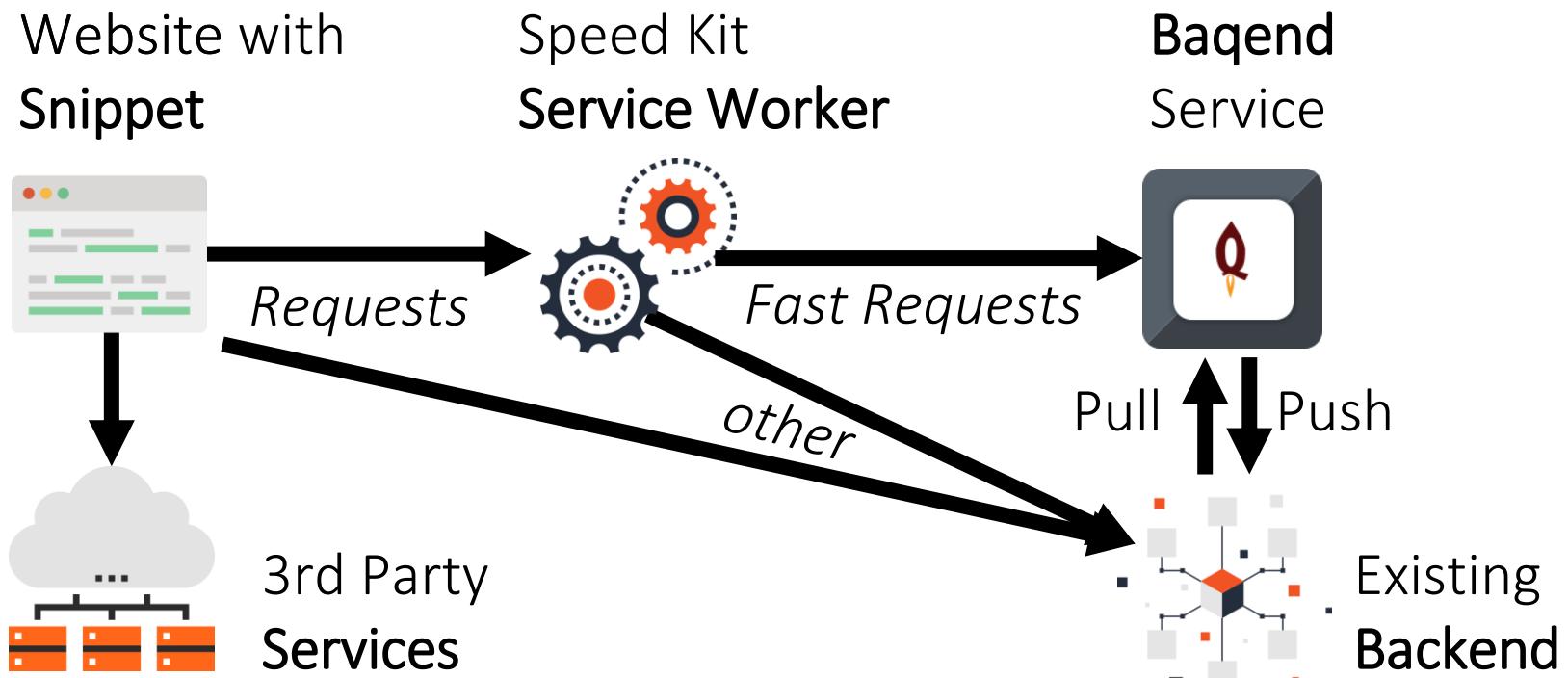
Speed Kit



- Turns Existing Sites
into **PWAs**
- **50-300% Faster Loads**
- **Offline Mode**

Speed Kit

Baqend Caching for Existing Websites



Speed Kit

Measure Your Page Speed!

<https://test.speed-kit.com/>

https://www.baur.de/



You are using Speed Kit 1.12.1

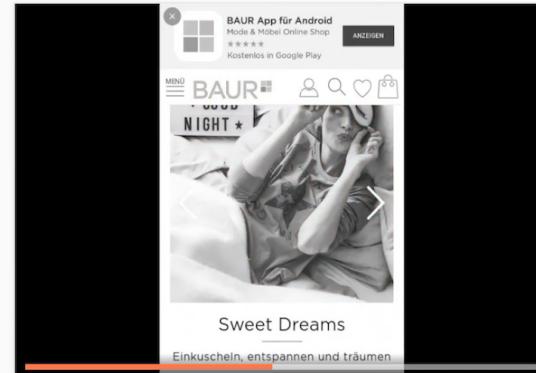
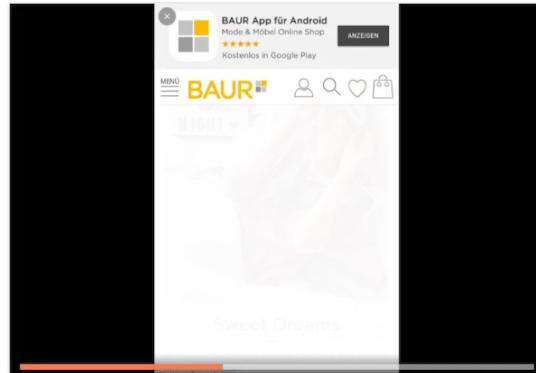
Without Speed Kit

1283 ms

1.9x
Faster

Your Website (Speed Kit 1.12.1)

662 ms



Without Speed Kit

1.3s

Your Website

0.7s

Average

Fast

Open Bachelor and Master Topics:

research.baqend.com



Offene Masterarbeiten

Providing Low Latency for PostgreSQL through the Orestes Caching Middleware

Orestes ist eine datenbankunabhängige Cloud-Middleware für (NoSQL) Datenbanken, die durch Caching extrem geringe Latenz erzielt. Orestes ist so konzipiert, dass durch die Implementierung bestimmter Schnittstellen (z.B. CRUD, Queries, Schemaänderungen) beliebige Datenbanken angebunden und durch eine *Unified REST API* einheitlich angesprochen werden können. Die Datenbank wird dabei automatisch angereichert um Fähigkeiten, die ihr sonst fehlen, u.a. global verteiltes Web-Caching, semi-strukturierte Schemata, Zugriffskontrolle durch ACLs, etc. Ziel der Arbeit ist die Anbindung eines der Systeme DynamoDB, Cassandra, HBase oder PostgreSQL mit einer Studie darüber, welche Eigenschaften sich erzielen lassen und welche aufgrund der Datenbankarchitektur und seiner Schnittstellen nicht erreichbar sind.

Ziele:

- Analyse der Datenbankeigenschaften aus Sicht der verschiedenen Orestes-Schnittstellen für CRUD, Queries, Partial Updates etc.
- Problembeschreibung: Identifikation der fehlenden Fähigkeiten der Datenbank und Studie, welche sich durch Orestes automatisch erzielen lassen.
- **Konsistenzmodell:** Welche Auswirkung haben die Konsistenzgarantien der Datenbank für Clients und wie spielen sie mit Caching zusammen?

Open Student Positions at Baqend:

baqend.com/jobs.html



Customer Success Manager

Help customers in leveraging Baqend technology.

Student Position: [Download PDF](#)



UX/Web Designer

Create user interfaces people love.

Full Time: [Download PDF](#)
Internship: [Download PDF](#)



JavaScript Engineer

High-performance, modern JavaScript development.

Full Time: [Download PDF](#)
Internship: [Download PDF](#)



DevOps Engineer

Build scalable and reliable cloud services.

Full Time: [Download PDF](#)
Internship: [Download PDF](#)



Growth Hacker

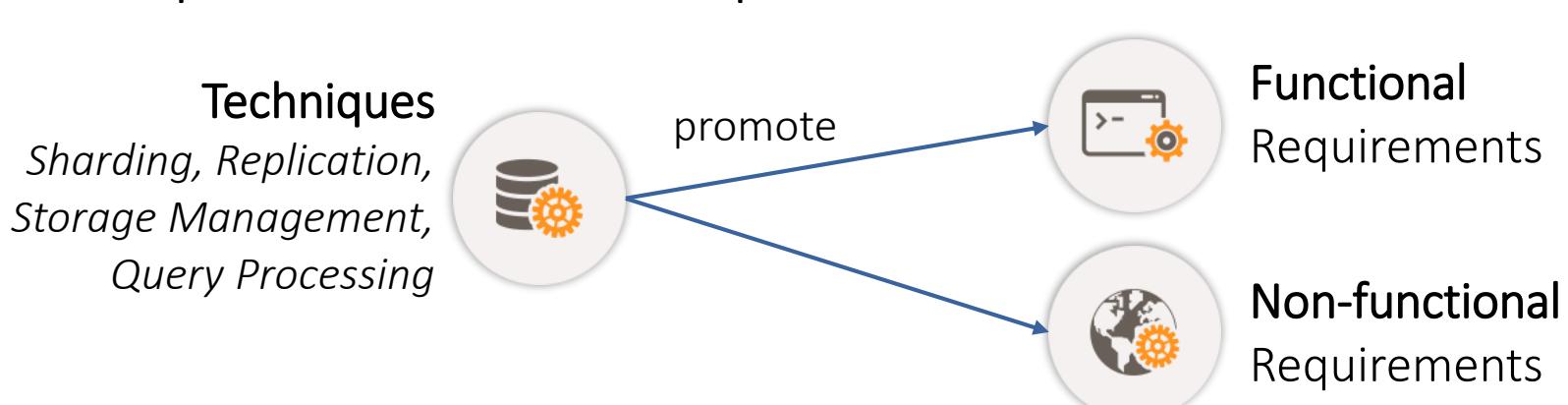
Versatile marketing for growing our startup.

Full Time: [Download PDF](#)
Internship: [Download PDF](#)



Summary (1/2)

- ▶ High-Level NoSQL Categories:
 - ▶ Key-Value, Wide-Column, Document, Graph
 - ▶ Two out of {Consistent, Available, Partition Tolerant}
- ▶ The **NoSQL Toolbox**: systems use similar techniques that promote certain capabilities



- ▶ Decision Tree



Summary (2/2)

- ▶ Variety of different NoSQL systems:
 - ▶ **Dynamo and Riak:**
 - ▶ KV-store with consistent hashing
 - ▶ **Redis:**
 - ▶ replicated, in-memory KV-store
 - ▶ **BigTable, HBase, Cassandra:**
 - ▶ wide-column stores
 - ▶ **MongoDB:**
 - ▶ sharded and replicated document store