



# NoSQL and Big Data

**Lecture**

Felix Gessert, Norbert Ritter

(22.05.2017)

# Outline



Foundations: Big Data,  
Scalability, Availability



The 4 Classes of NoSQL  
Databases



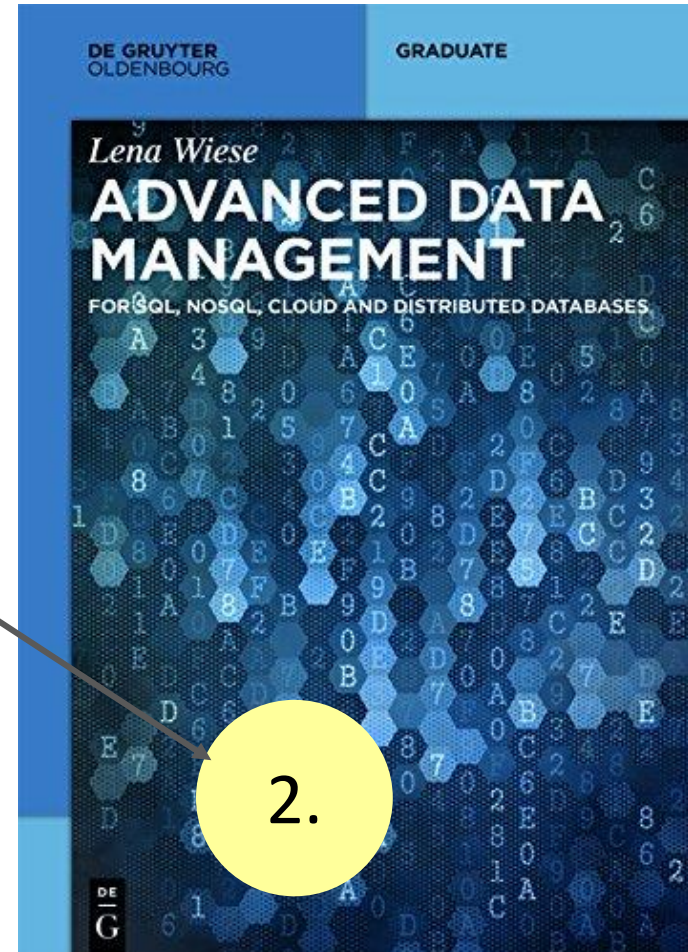
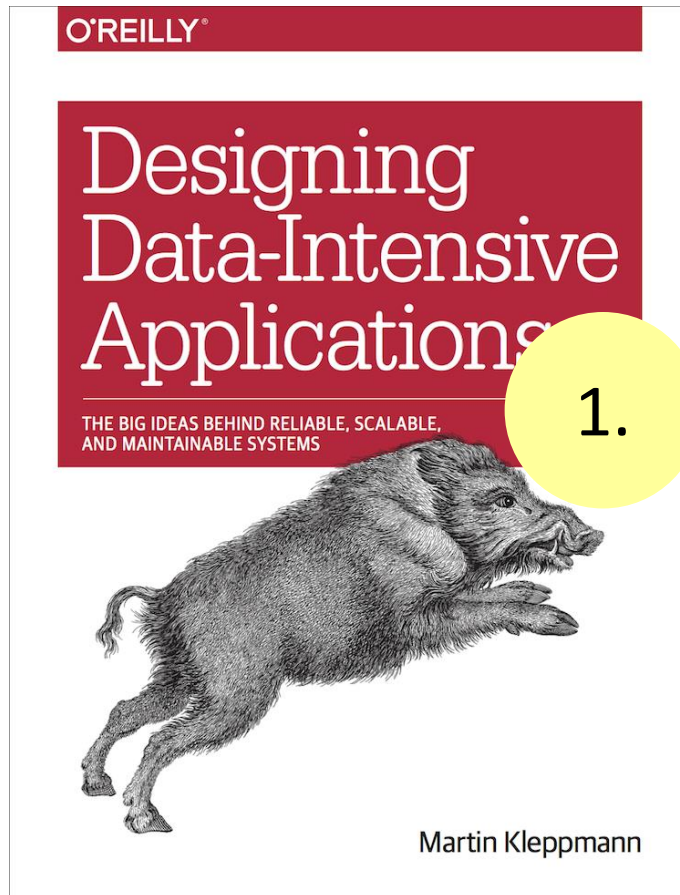
NoSQL Examples: concrete  
Architectures, Systems, APIs



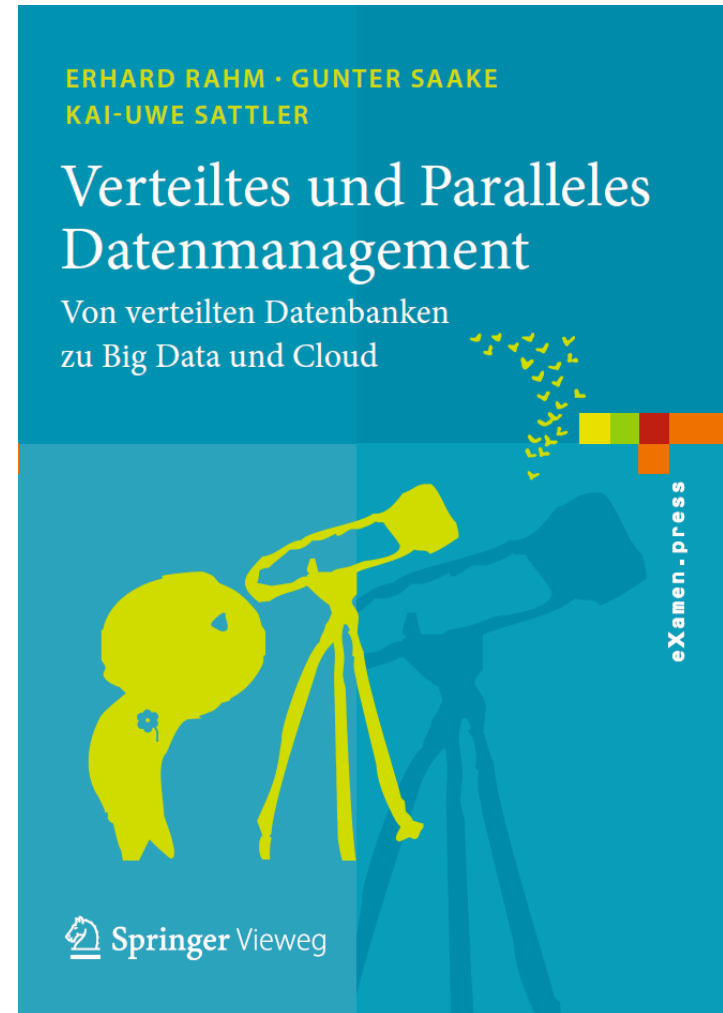
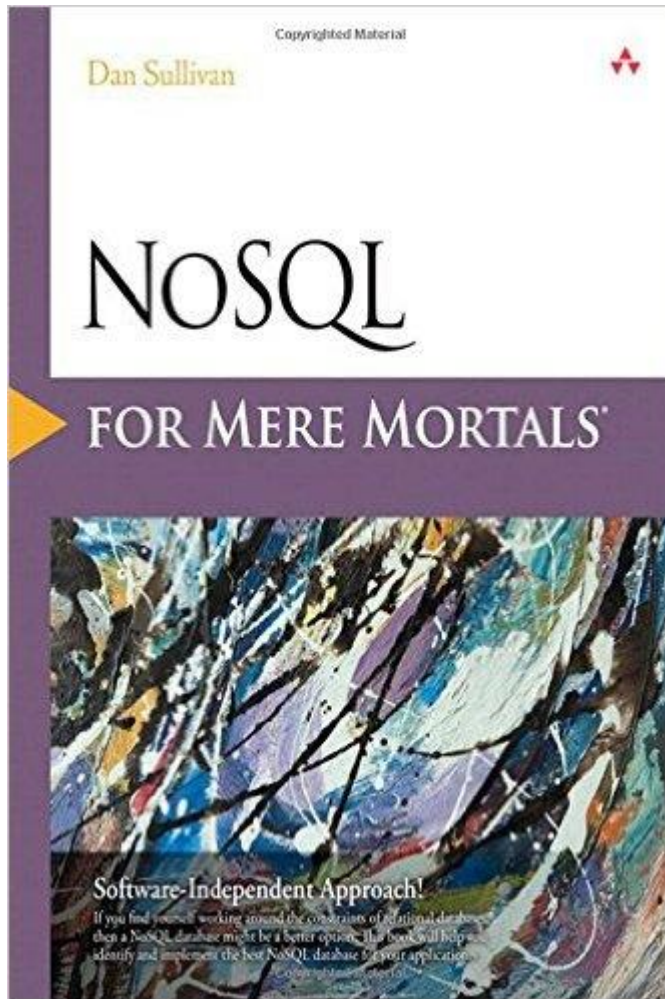
Cloud Databases

- Literature
- Motivation
  - Big Data
  - NoSQL
- CAP Theorem
- 2-Phase-Commit
- NoSQL Triangle
- ACID vs BASE

# Recommended Literature: NoSQL



# Other Literature



# Recommended Literature: Blogs

**Martin Kleppmann**

<https://martin.kleppmann.com/>



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



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

**NoSQL Weekly**

<http://www.nosqlweekly.com/>



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

**High Scalability**

<http://highscalability.com/>

# The Database Explosion

## Sweetspots



*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

# The Database Explosion

## Cloud-Database Sweetspots



*Realtime BaaS*

Communication and  
collaboration



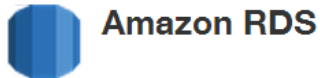
*Wide-Column Store*

Very large tables



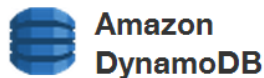
*Managed NoSQL*

Full-Text Search



*Managed RDBMS*

General-purpose  
ACID transactions



*Wide-Column Store*

Massive user-  
generated content



*Object Store*

Massive File  
Storage



*Managed Cache*

Caching and  
transient storage



*Backend-as-a-Service*

Small Websites  
and Apps

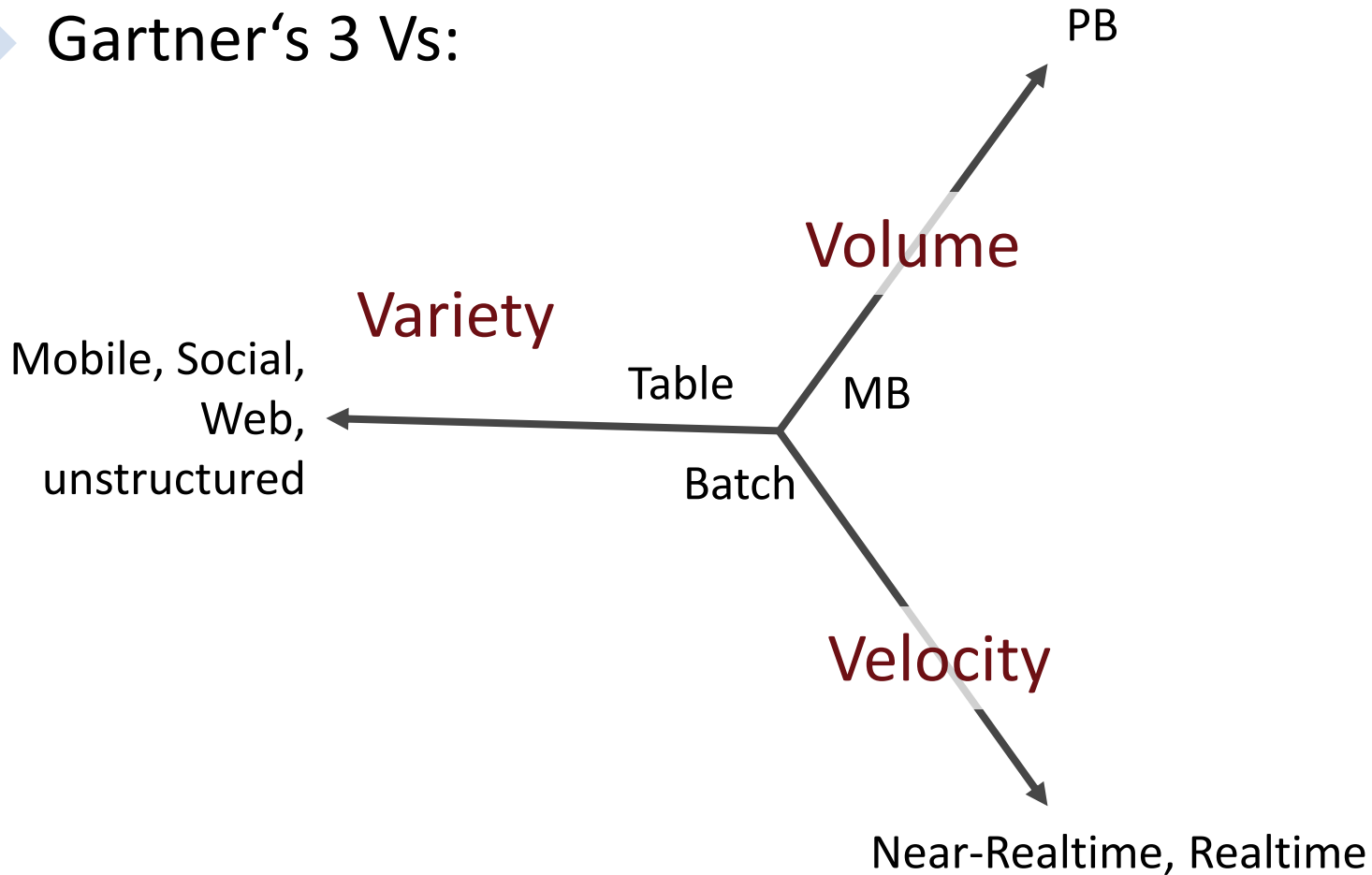


*Hadoop-as-a-Service*

Big Data Analytics

# Motivation Big Data

- ▶ Gartner's 3 Vs:





# 4Vs

## 40 ZETTABYTES

[ 43 TRILLION GIGABYTES ]

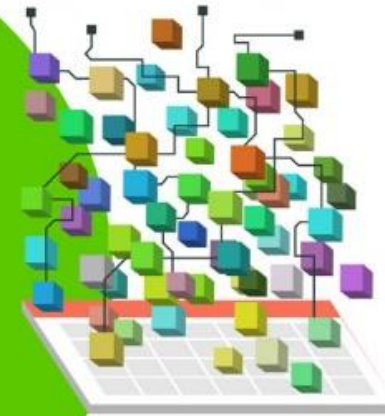
of data will be created by 2020, an increase of 300 times from 2005



## It's estimated that 2.5 QUINTILLION BYTES

[ 2.3 TRILLION GIGABYTES ]

of data are created each day



**6 BILLION PEOPLE**  
have cell phones



**WORLD POPULATION: 7 BILLION**

## Volume SCALE OF DATA



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

**1 TB OF TRADE INFORMATION**

during each trading session



Modern cars have close to

**100 SENSORS**

that monitor items such as fuel level and tire pressure

**Velocity**  
ANALYSIS OF  
STREAMING DATA

By 2016, it is projected there will be

**18.9 BILLION NETWORK CONNECTIONS**

– almost 2.5 connections per person on earth



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

# 4Vs

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

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]



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



**Variety**  
DIFFERENT FORMS OF DATA

**30 BILLION PIECES OF CONTENT**

are shared on Facebook every month



**400 MILLION TWEETS**

are sent per day by about 200 million monthly active users



# 4Vs

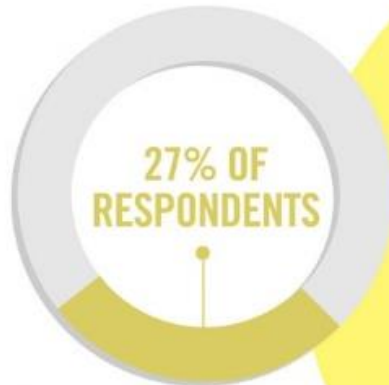
## 1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



Poor data quality costs the US economy around

**\$3.1 TRILLION A YEAR**

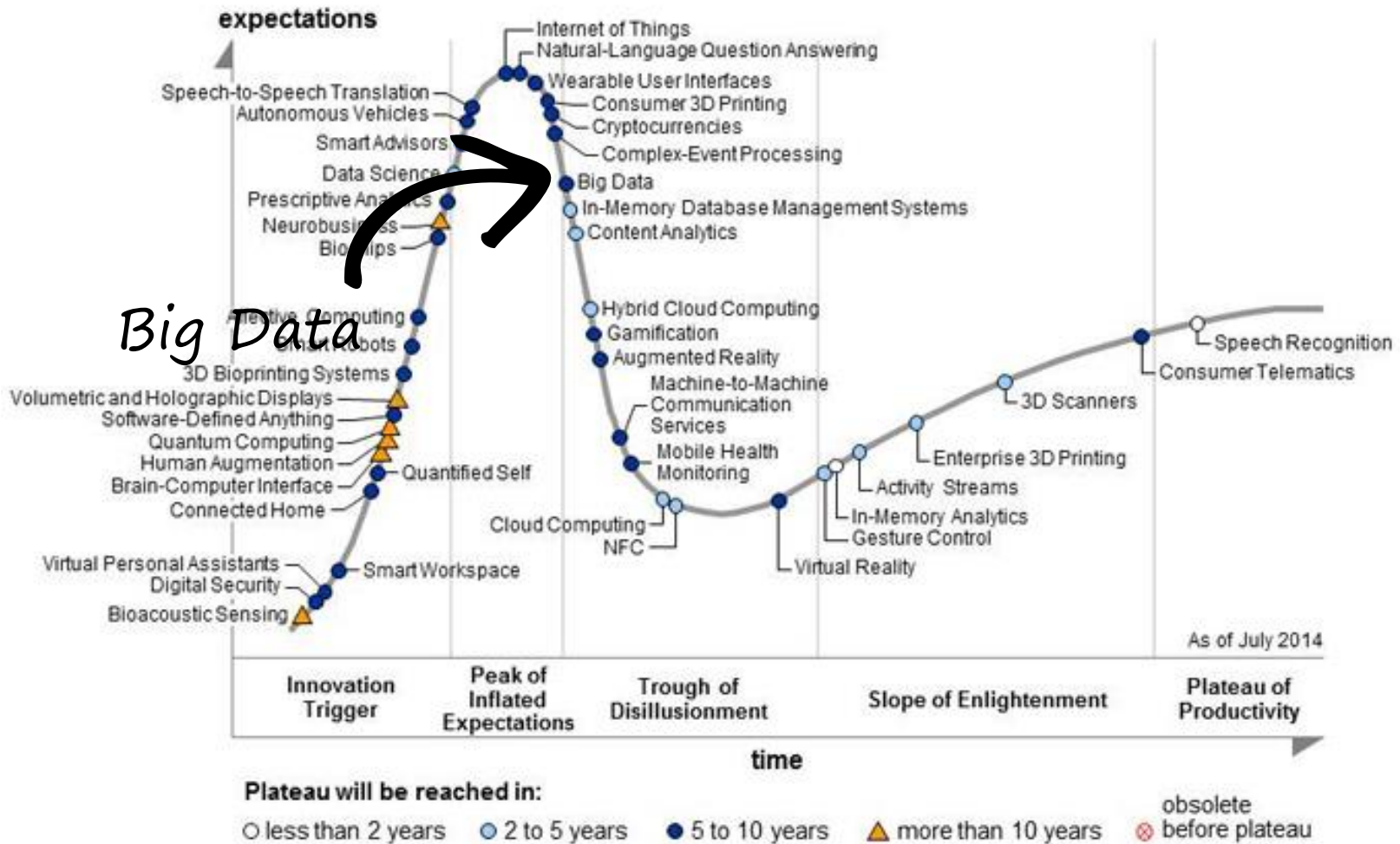


in one survey were unsure of how much of their data was inaccurate

**Veracity**  
UNCERTAINTY OF DATA



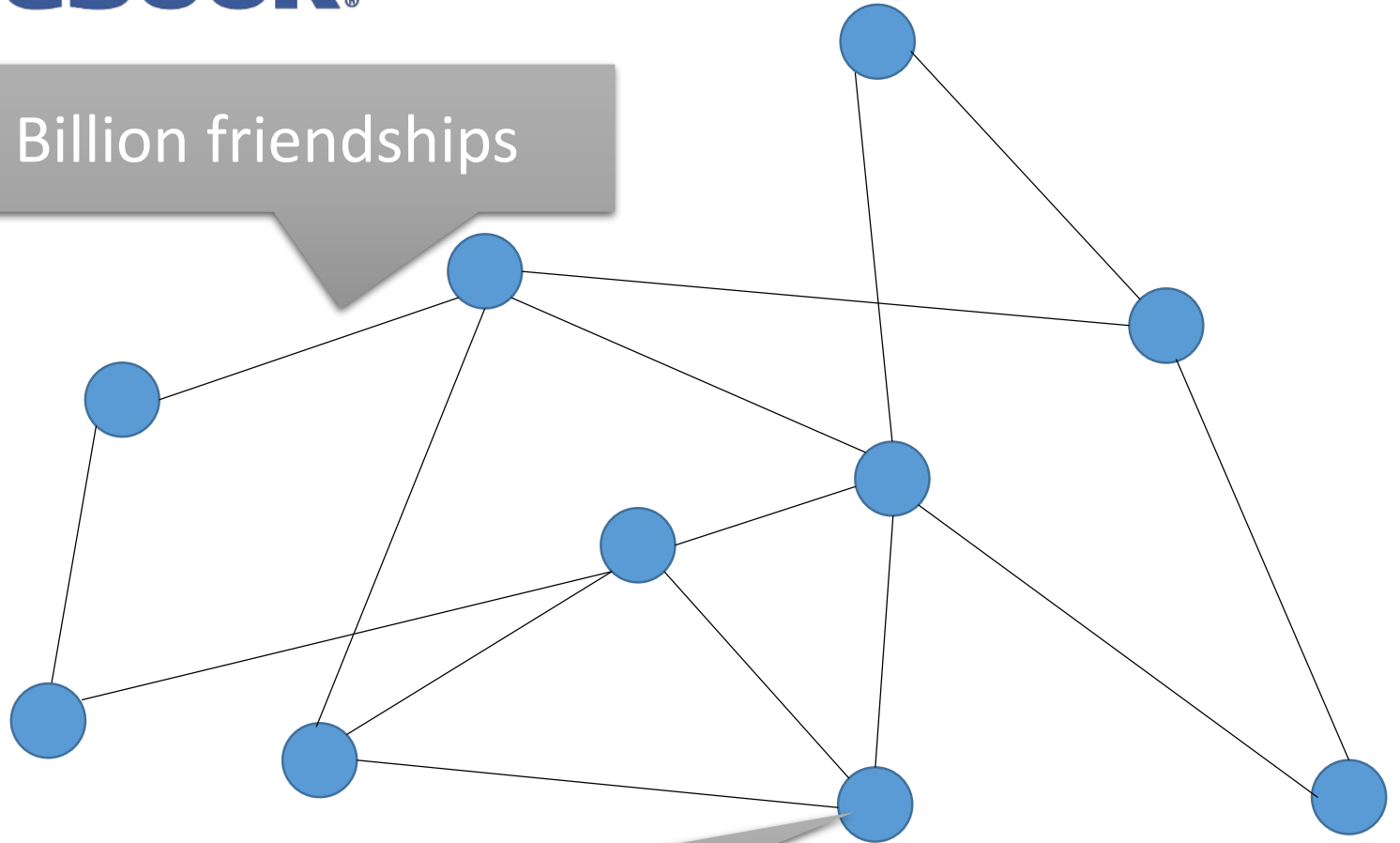
# Hype Cycle



# facebook.

1.4 Billion Members

140 Billion friendships



290 Billion Photos

# Data flood

**Climate Research:** The *Deutsche Klimarechenzentrum (DKRZ)* stores **60 PB** climate data.

**Archiving:** The Internet Archive stores **10 PB** of archived websites.

**Gaming:** World of Warcraft needs **1.3 PB** for storing the game state. Steam delivers **over 30 PB** of data per month.

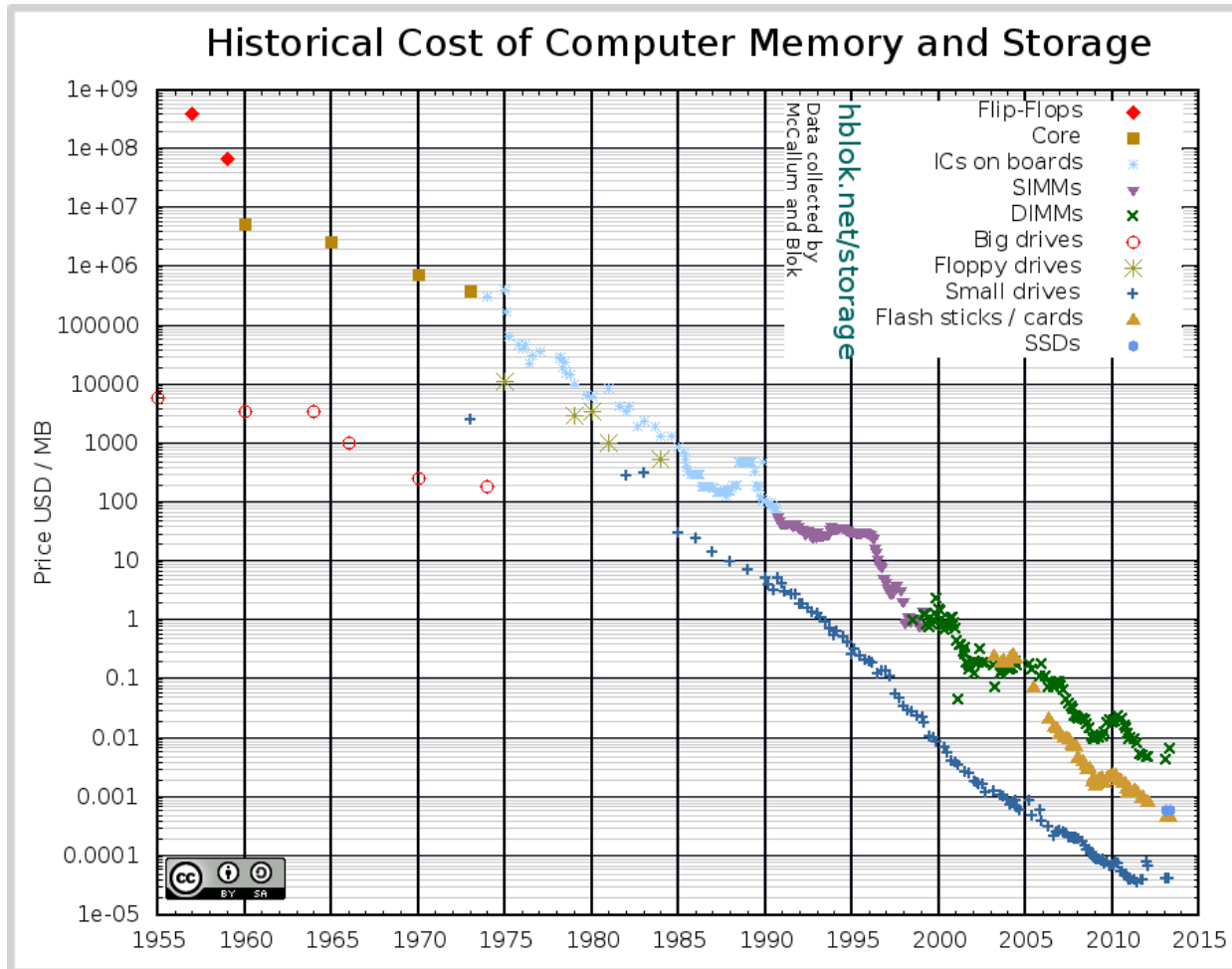
**Movies:** The CGI-effects in Avatar (2009) needed **over 1 PB** storage for rendering.

**Supercomputing:** The Blue Waters Supercomputer is planned to have a storage capacity of **500 PB**.

**Particle Physics:** In search of the Higgs-Boson CERN gathered **200 PB** of data.

**Email:** In May 2013 Microsoft announced that for the migration of Hotmail to outlook.com **over 150 PB** user data were transferred.

# Dropping storage costs





# Big Data Defined

- ▶ Big Data has two sides:

## Big Data

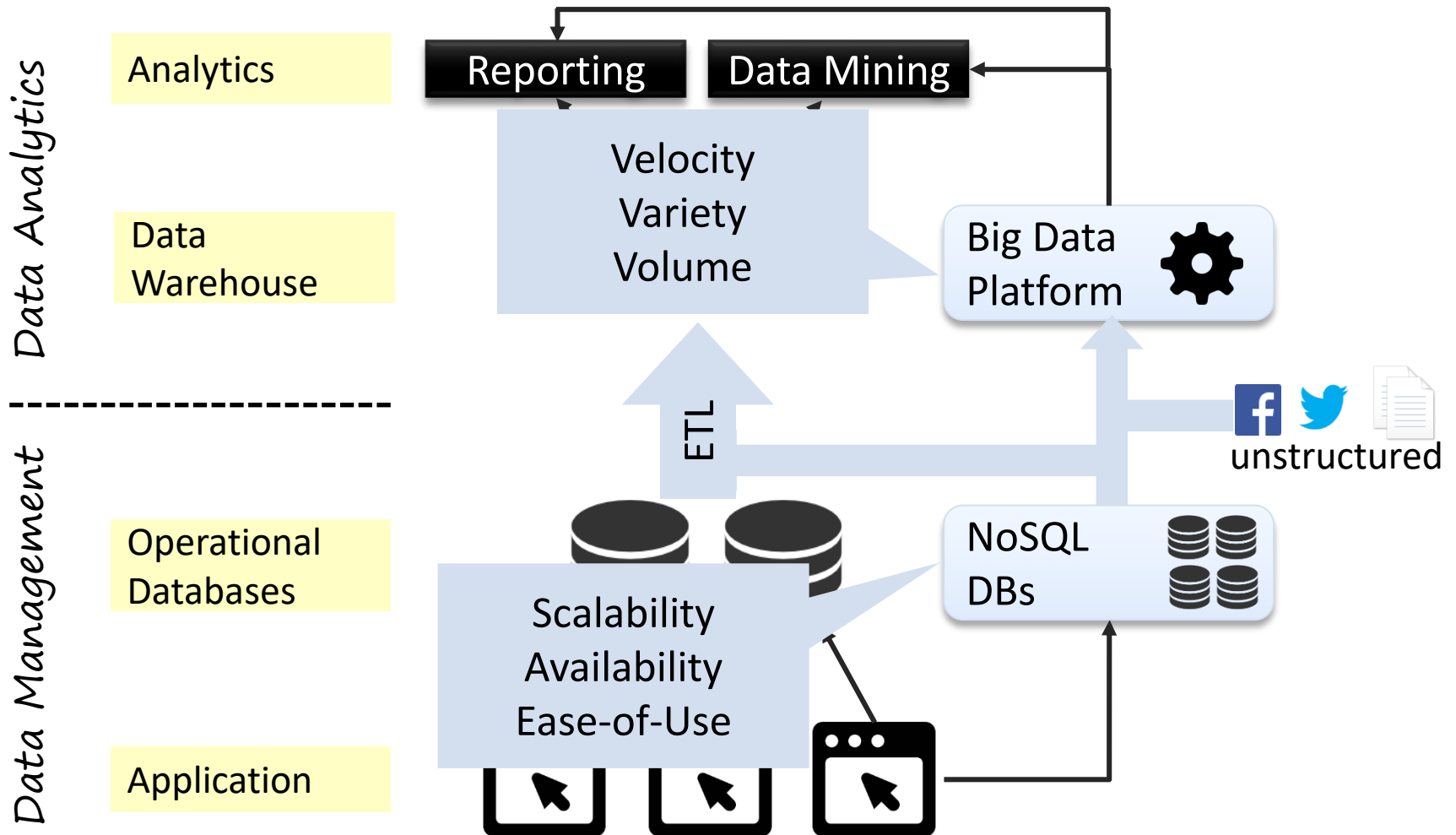
### Big Data Management

- OLTP
- Often referred to as "NoSQL"
- e.g. MongoDB, HBase, Cassandra

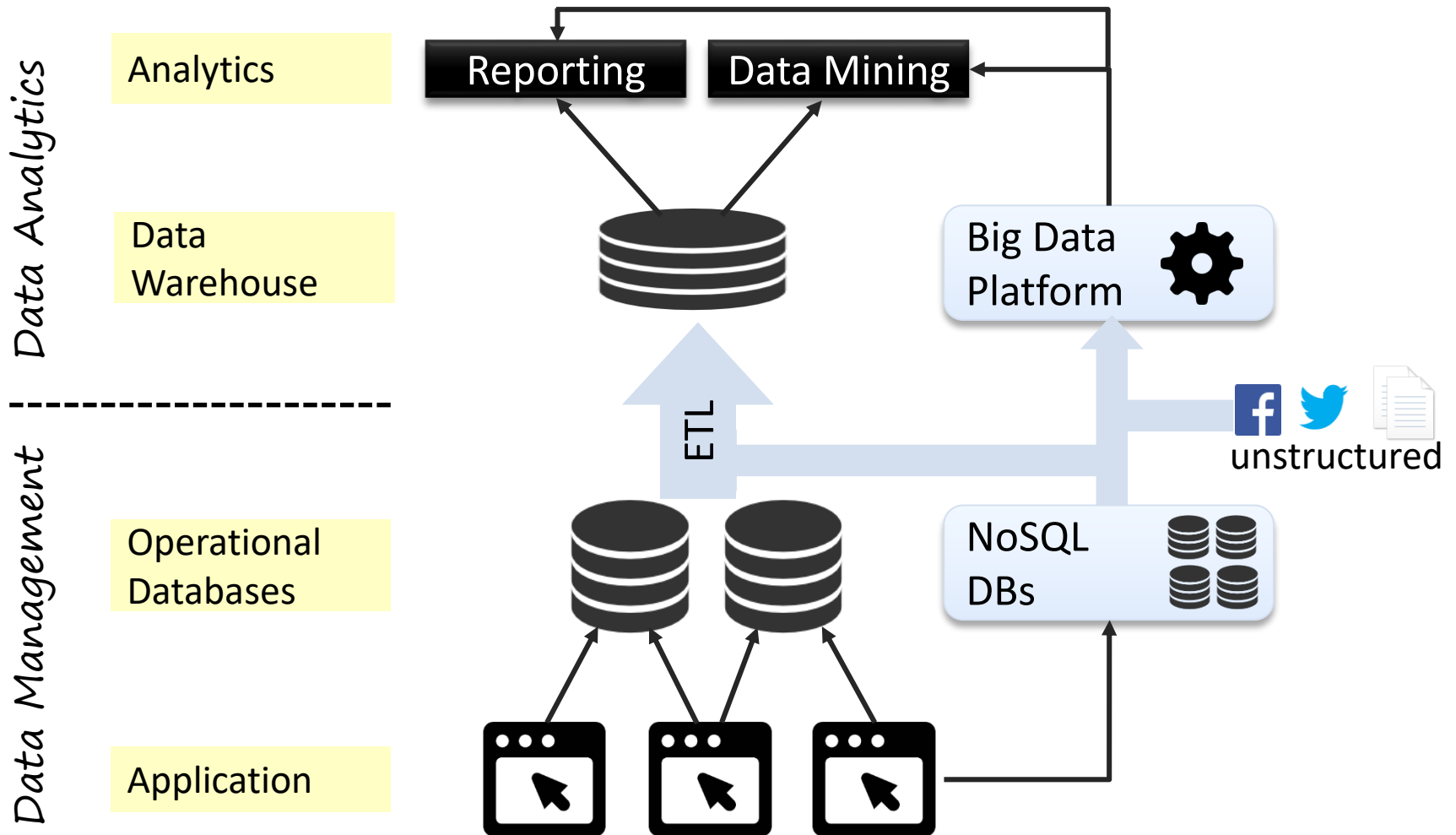
### Big Data Analytics

- OLAP
- Often referred to as „Big Data“
- e.g. Hadoop, Storm

# Architectural Change



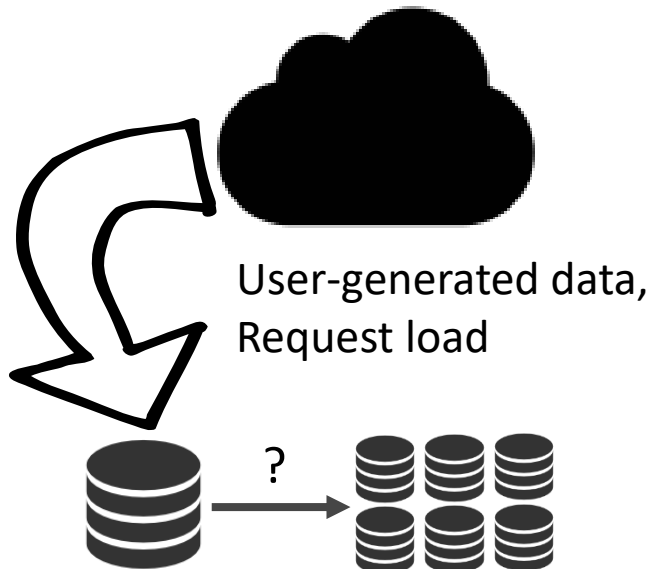
# Architectural Change



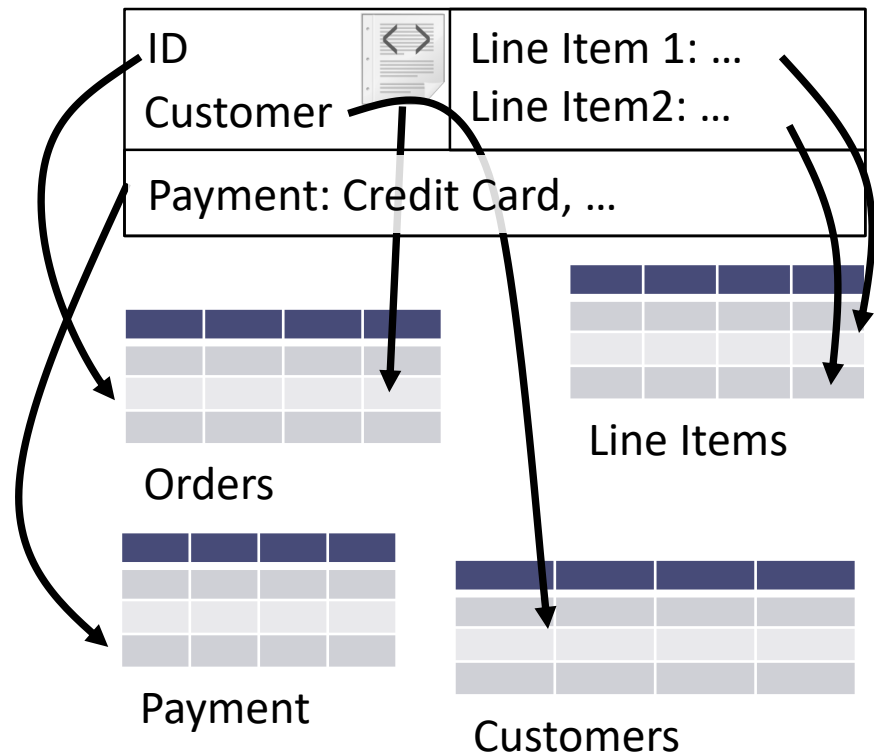
# NoSQL Databases

- ▶ Two main motivations:

## Scalability



## Impedance Mismatch



# 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 client, Replication: HDFS Replication, Write in: Java, Concurrency: 1, Misc: Links: 1 Books (1, 2, 3)

**Cassandra:** massively scalable, partitioned row store, distributed architecture, linear scale performance, no single points of failure, read/write support across multiple data centers & cloud availability zones. API / Query Method: CQL and Thrift, Replication: pccr-to-pccr, Write in: Java, Concurrency: tunable consistency, Misc: built-in data compression, MapReduce support, primary/secondary indexes, security features. Links: [Documentation](#), [Privacy Policy](#)

**HyperTable:** API: Thrift (Java, PHP, Perl, Python, Ruby, C#), Protocol: Thrift, Query Method: HQL, native Thrift API, Replication: HDFS Replication, Concurrency: MVCC, Consistency Model: Fully consistent, Misc: High performance C++ implementation of Google's Bigtable. [Commercial Support](#)

**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, iterative compression, and a server-side programming mechanism that can modify key/value pairs at various points in the data management process.

**Amazon SimpleDB:** Misc: not open source / part of AWS, [Docs](#) (will be outperformed by DynamoDB?)

**Cloudata:** Google's Bigtable clone. Misc: [Website](#)

**Cloudera:** Professional Software & Services based on Hadoop

**HPCC:** from [Lovelace](#), [Info](#), [Wiki](#)

**Stratosphere:** (research system) massive parallel & flexible execution, high generalization and extension ([paper](#), [paper](#), [documentation](#), [Google](#), [IBM](#))

---

**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, Misc: indexing, GridFS, [Framework](#) & [Commercial License](#), Links: [Talk](#), [Wiki](#), [Company](#)

**Elasticsearch:** API: REST and many languages, Protocol: REST, Query Method: via JSOM, Replication: Sharding, automatic and configurable, Write in: Java, Misc: schema mapping, multi-tenancy with arbitrary indexes, Company and Support: [Elastic](#)

**Couchbase Server:** API: Memcached API+protocol (Binary and ASPI), [Most Languages](#), Protocol: Memcached, REST interface for cluster conf & management, Write in: C/C++ & Erlang, clustering, Replication: Pccr to Pccr, fully consistent, Misc: transparent topology changes during operation, provides memcached-compatible caching buckets, commercially supported version available. Links: [Intro](#), [Website](#)

**CouchDB:** API: JSON, Protocol: REST, Query Method: MapReduce of JavaScript Funcs, Replication: Master Master, Write in: Erlang, Concurrency: MVCC, Misc: Links: [3 CouchDB books](#), [Couch Lounge](#) (partitioning / clustering), [Dr. Dobbs](#)

**Redis:** API: protobuf-based, Query Method: untyped, chainable query language (find, JOINs, sub-queries, MapReduce, GroupMapReduce), Replication: Sync and Async, Master Slave with portable acknowledgment, Streaming, geo-set, range-based, Write in: C++, Concurrency: MVCC, Misc: log-structured storage engine with concurrent incremental merge compact

**RavenDB:** .NET solution, Provides HTTP/JSON access, LINK queries & Sharding supported. [Wiki](#)

**MarkLogic Server:** (Research-Commercial) API: JSON, XML, Java Protocols: HTTP, REST, Query Method: Full Text Search, XPath, XQuery, Range, Geospatial, Write in: C++, Concurrency: Sharded-nothing cluster, MVCC, Misc: Prolog-like, circular, ACID transactions & auto-sharding, failover, master slave replication, secure with ADAS, developer Community: [MarkLogic](#)

**Clustertopline Server:** (Research-Commercial) API: XML, PHP, Java, .NET, Protocols: HTTP, REST, native TCP/IP, Query Method: full text search, XML, range and XPath queries, Write in: C++ Concurrency: ACID, compliant, transactional, multi-master cluster, Misc: Petabyte-scalable document store and full text search engine, Information ranking, Replication: Cloudable

**ThruDB:** (please help provide more facts) Uses Apache Thrift to integrate multiple backend databases as BerkeleyDB, Disk, MySQL, etc.

**Terrastore:** API: Java & http, Protocol: http, Language: Java, Query: Range queries, Predicates, Replication: Partitioned with consistent hashing, Consistency: Per-record strict consistency, Misc: Based on Terracotta

**LasDB:** lightweight open source document database written in Java for high performance, runs in memory, supports InnoDB, API: JSON, Java Query Method: REST OData Style Query language, Java fluent Query API, Concurrency: Atomic document writes, Indexes: eventually consistent indexes

**RaptorDB:** JSON based, Document store database with COMPLEX, not MAP functions and automatic hybrid schema mapping and LINQ query filters

**SiigoDB:** A Document Store on top of SQL-Server.

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

**CloudDB:** SDB with: BSON, Protocol: C++, Query Method: dynamic queries and map/reduce, Drivers: Java, C++, PHP, Misc: ACID compliant, Full shell console over people @ engine, cloned requirements are submitted by users, [see issues](#), [Twitter](#), [fb](#), [ask](#), [stackoverflow](#)

# Big Data Analytics

- ▶ **Idea:** make existing massive, unstructured data amounts usable

sources

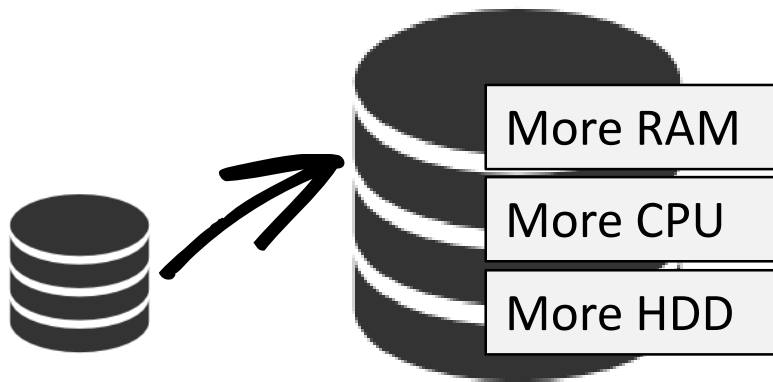


Analyst, Data Scientist,  
Software Developer

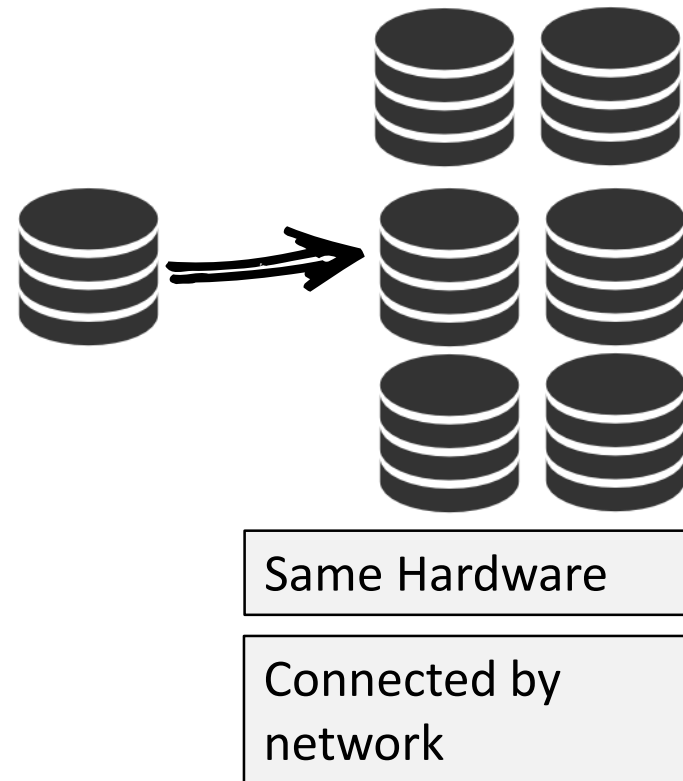
- Structured data (DBs)
  - Log files
  - Documents, Texts, Tables
  - Images, Videos
  - Sensor data
  - Social Media, Data Services
- 
- Statistics, Cubes, Reports
  - Recommender
  - Classifiers, Clustering
  - Knowledge

# Scale-up vs Scale-out

**Scale-Up** (*vertical* scaling):

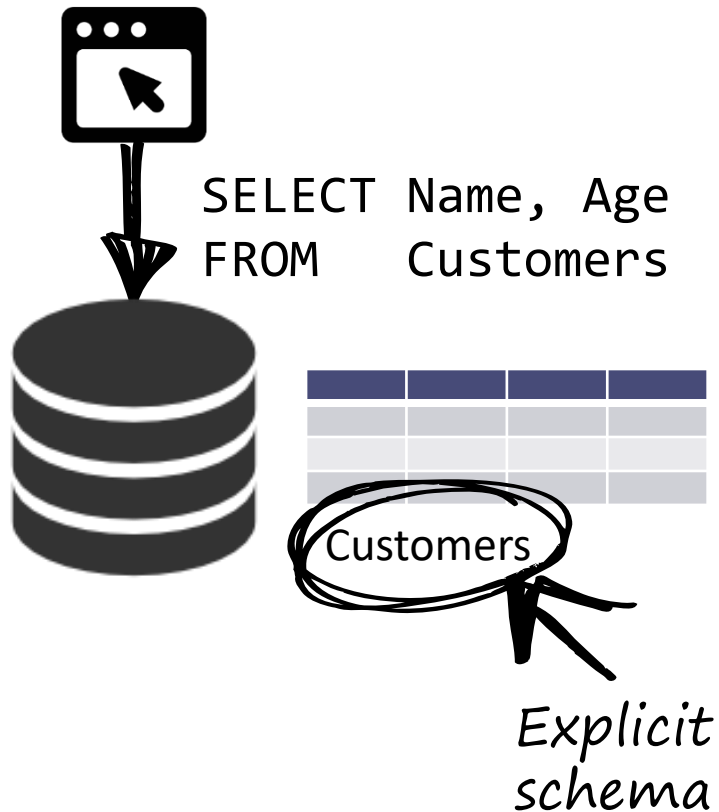


**Scale-Out** (*horizontal* scaling):

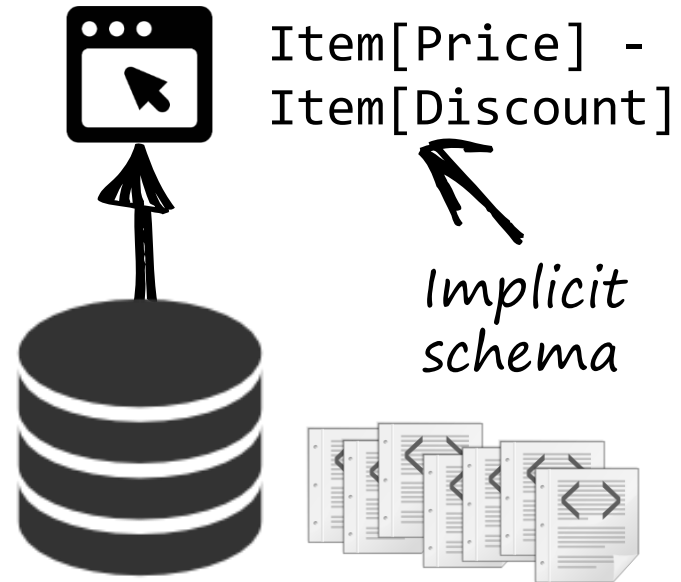


# Schema-free Data Modeling

RDBMS:



NoSQL DB:





# Paradigm Shift



Commercial DBMS

Specialized DB hardware  
(Oracle Exadata, etc.)

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

Highly Available Storage (SAN,  
RAID, etc.)



Open-Source DBMS

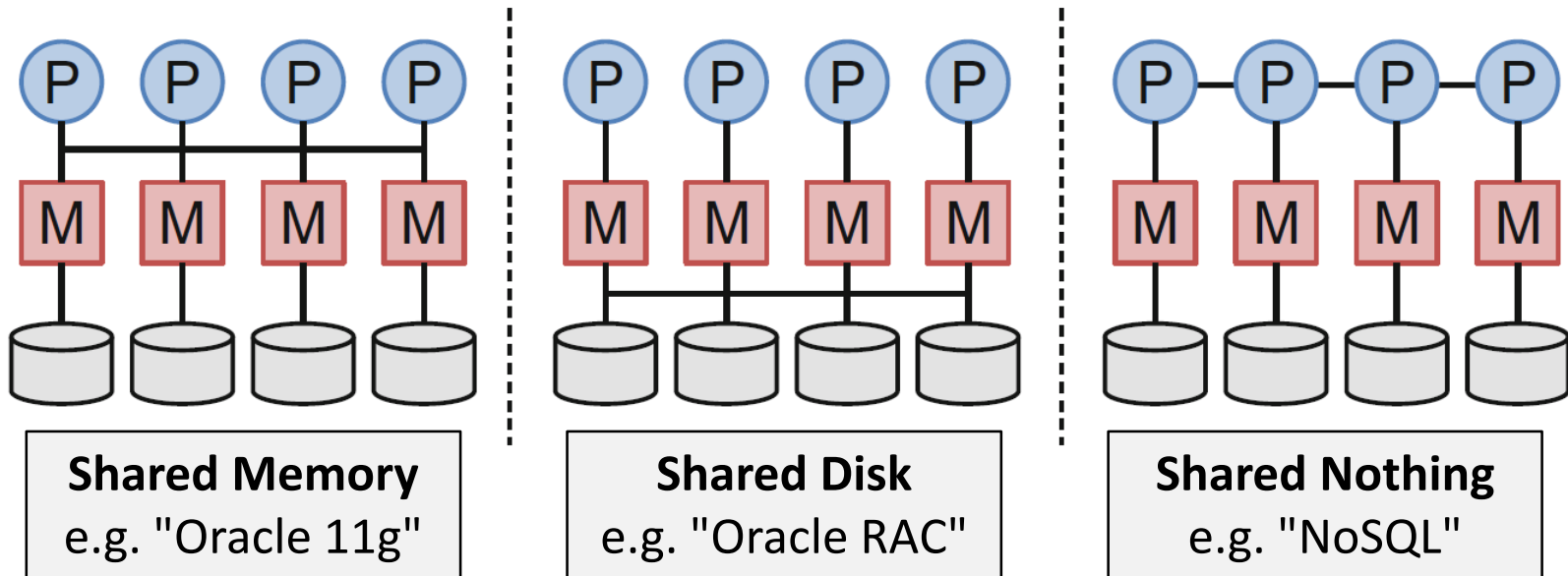
Commodity hardware

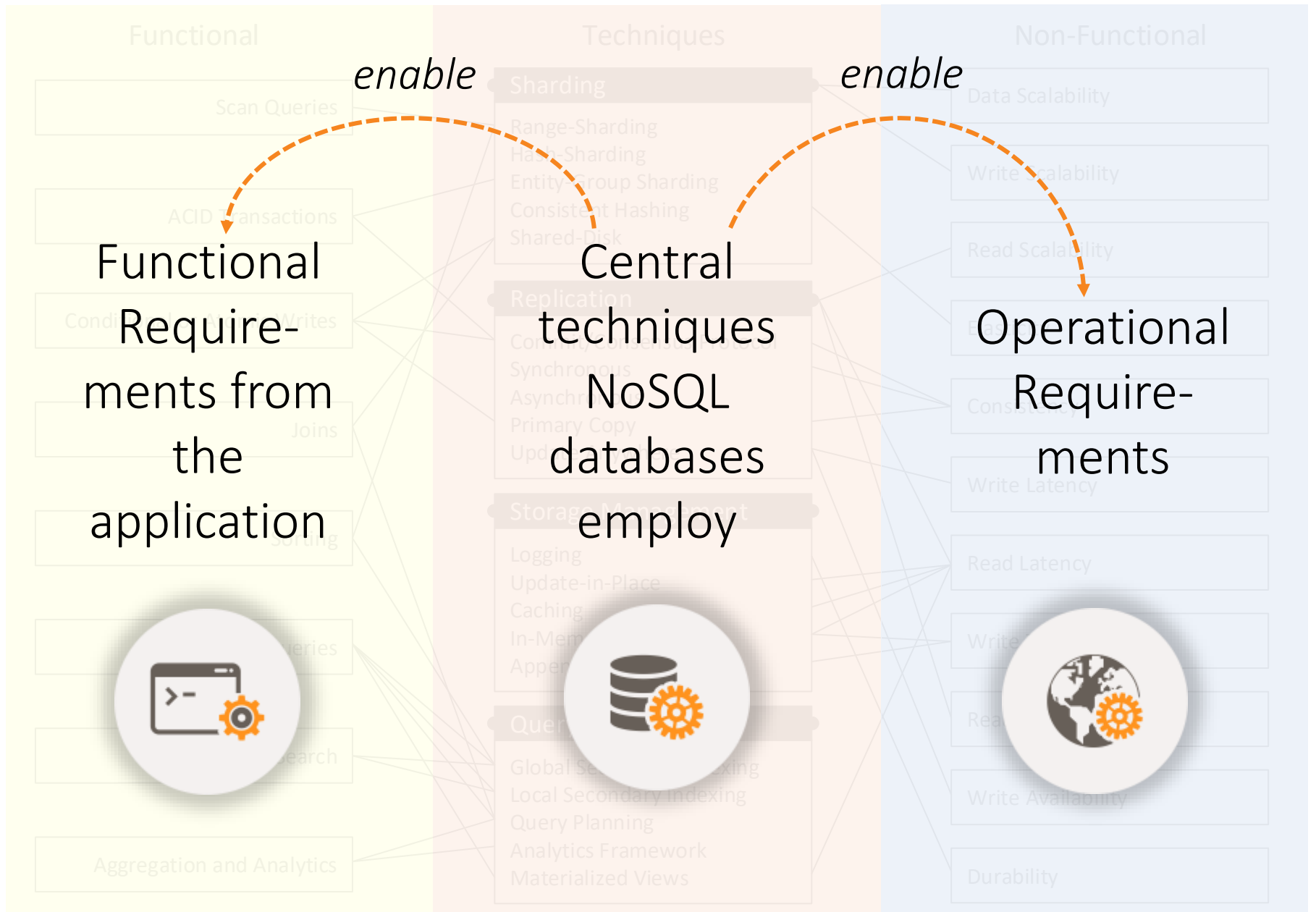
Commodity network  
(Ethernet, etc.)

Commodity drives (standard  
HDDs, JBOD)

# Paradigm Shift

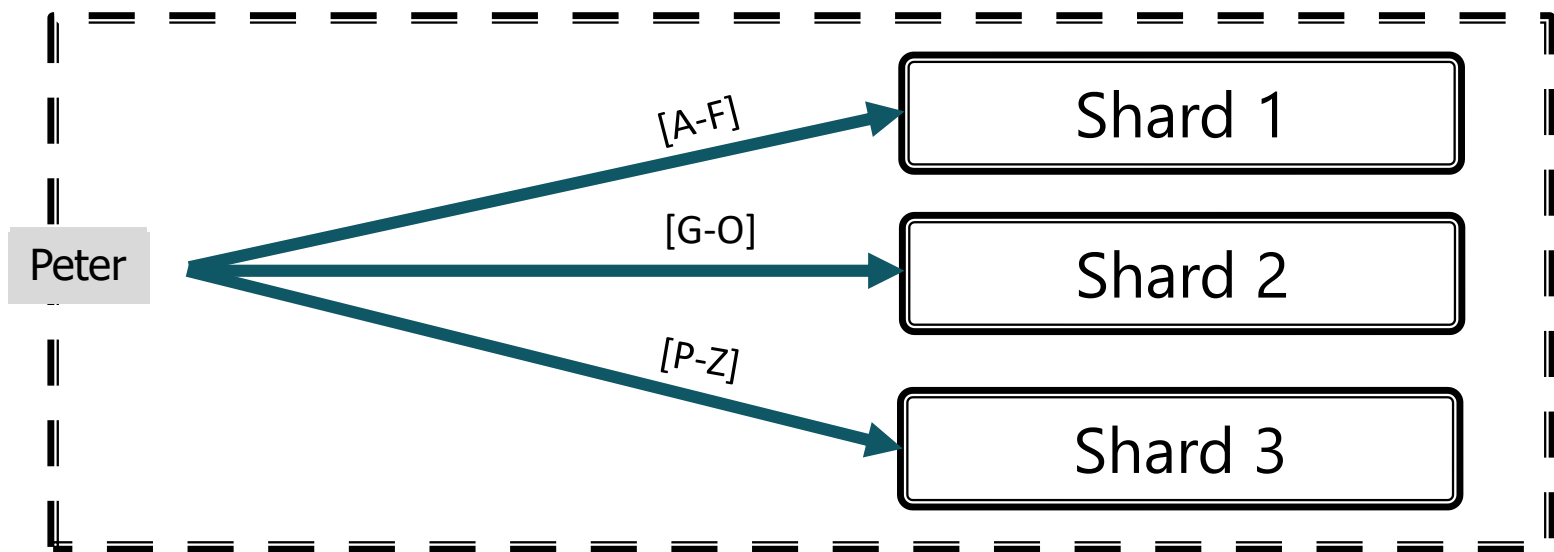
- ▶ Shift towards distributed computing architectures





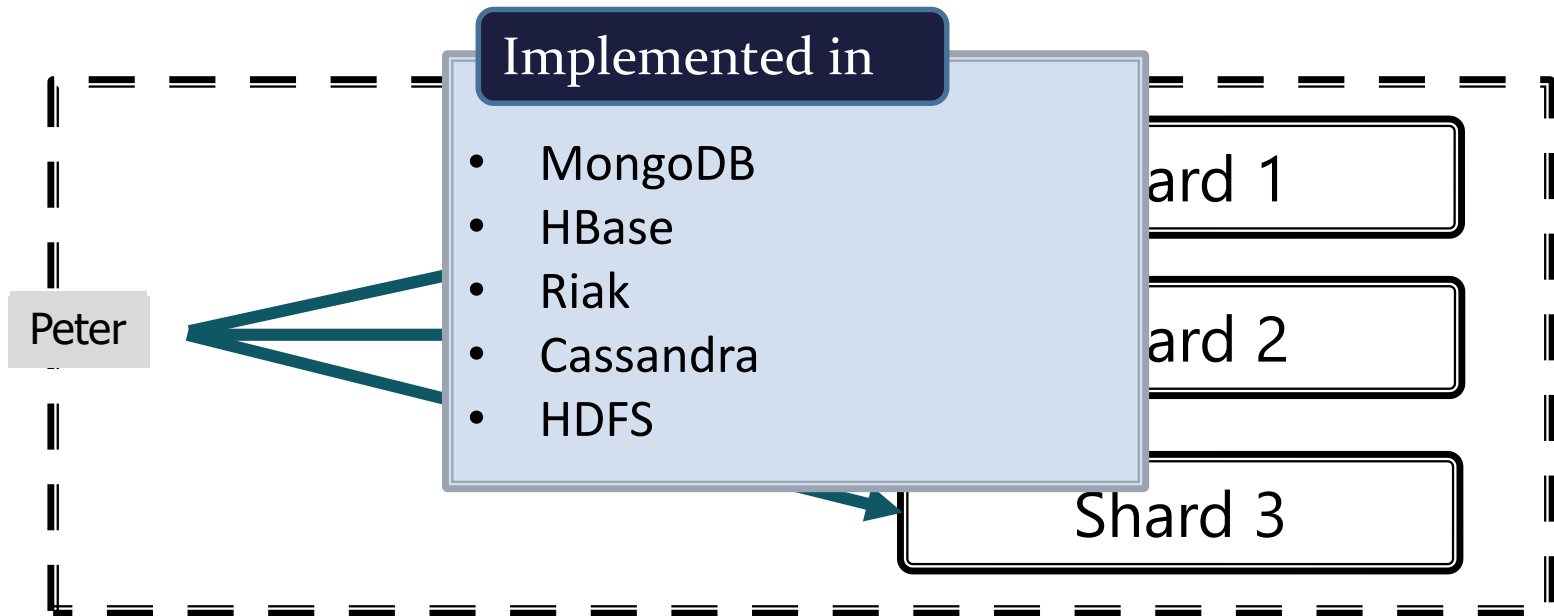
# Sharding (aka Partitioning, Fragmentation)

- ▶ Horizontal distribution of data over server nodes
- ▶ **Partitioning strategies:** Hash-based vs. Range-based
- ▶ **Difficulty:** Multi-Shard-Operations (join, aggregation)



# Sharding (aka Partitioning, Fragmentation)

- ▶ Horizontal distribution of data over server nodes
- ▶ **Partitioning strategies:** Hash-based vs. Range-based
- ▶ **Difficulty:** Multi-Shard-Operations (join, aggregation)



# Sharding

## Hash-based Sharding

- Builds hash of data values (e.g. over the key) to determine a partition (shard) for a data item (tuple)
- **Pro:** Perfectly even distribution
- **Contra:** No data locality – data items are pseudorandomly scattered over partitions

## Range-based Sharding

- Assigns ranges defined over fields (shard keys) to partitions
- **Pro:** Data locality preserved (for shard keys)
- **Contra:** distribution might grow uneven → repartitioning/balancing required

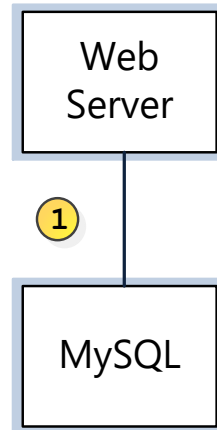
# Traditional Sharding

Example: **Tumblr**

- ▶ Caching
- ▶ Sharding from application

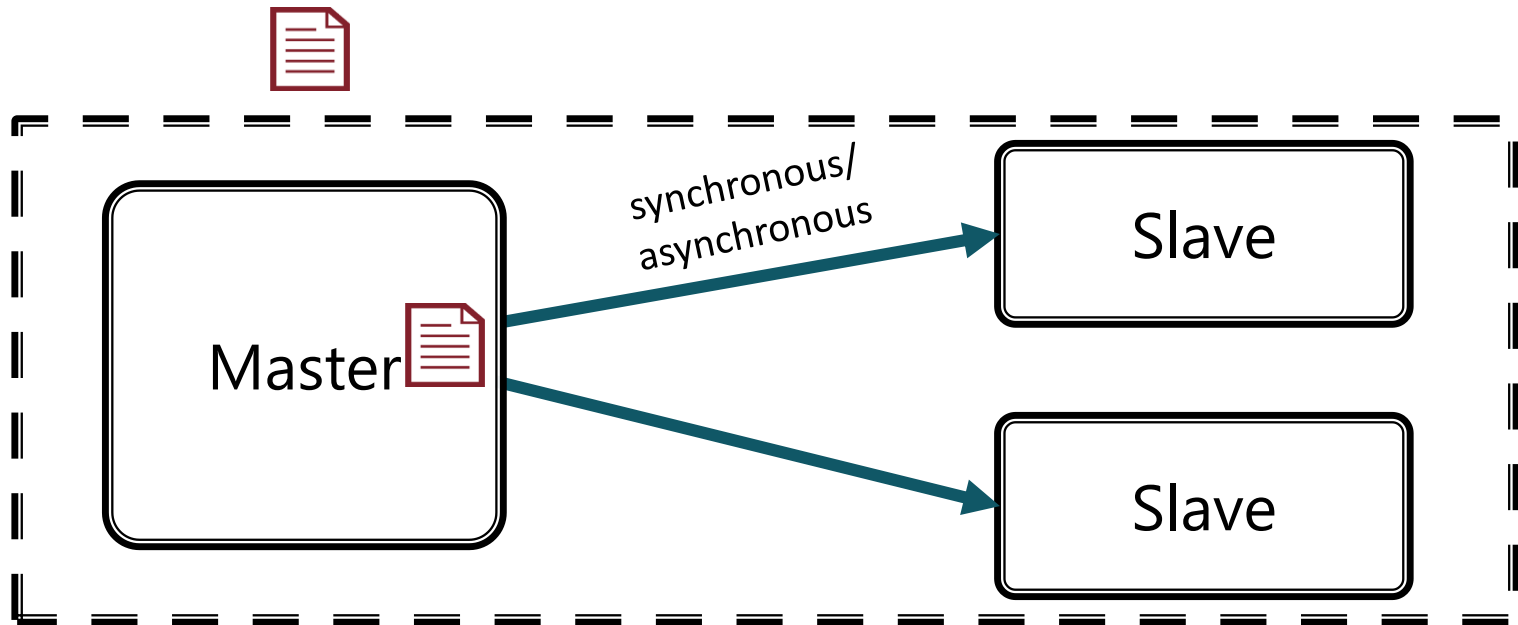
Moved towards:

- ▶ Redis
- ▶ HBase



# Replication

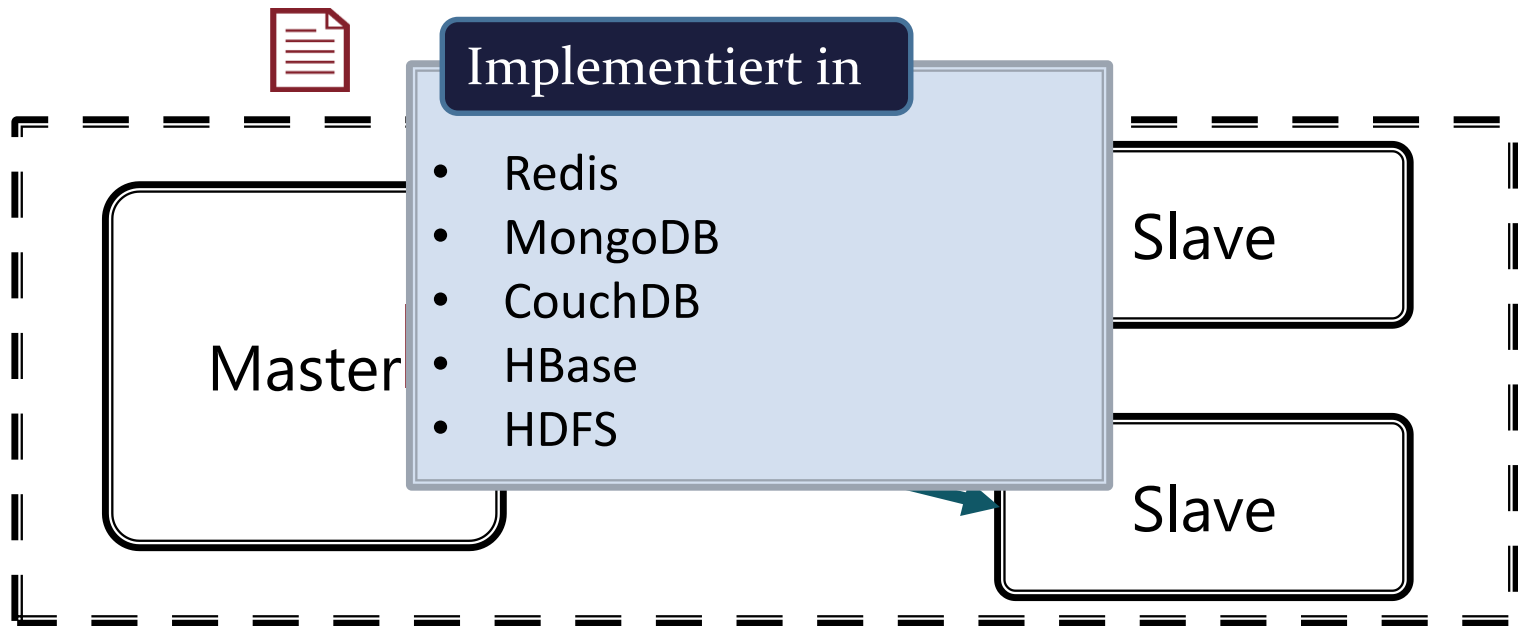
- ▶ Stores  $N$  copies of each data item
- ▶ **Consistency model:** synchronous vs asynchronous
- ▶ **Coordination:** Multi-Master, Master-Slave





# Replication

- ▶ Stores  $N$  copies of each data item
- ▶ **Consistency model:** synchronous vs asynchronous
- ▶ **Coordination:** Multi-Master, Master-Slave



# Replication: consistency models

## Asynchronous

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

## Synchronous

- 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: coordination

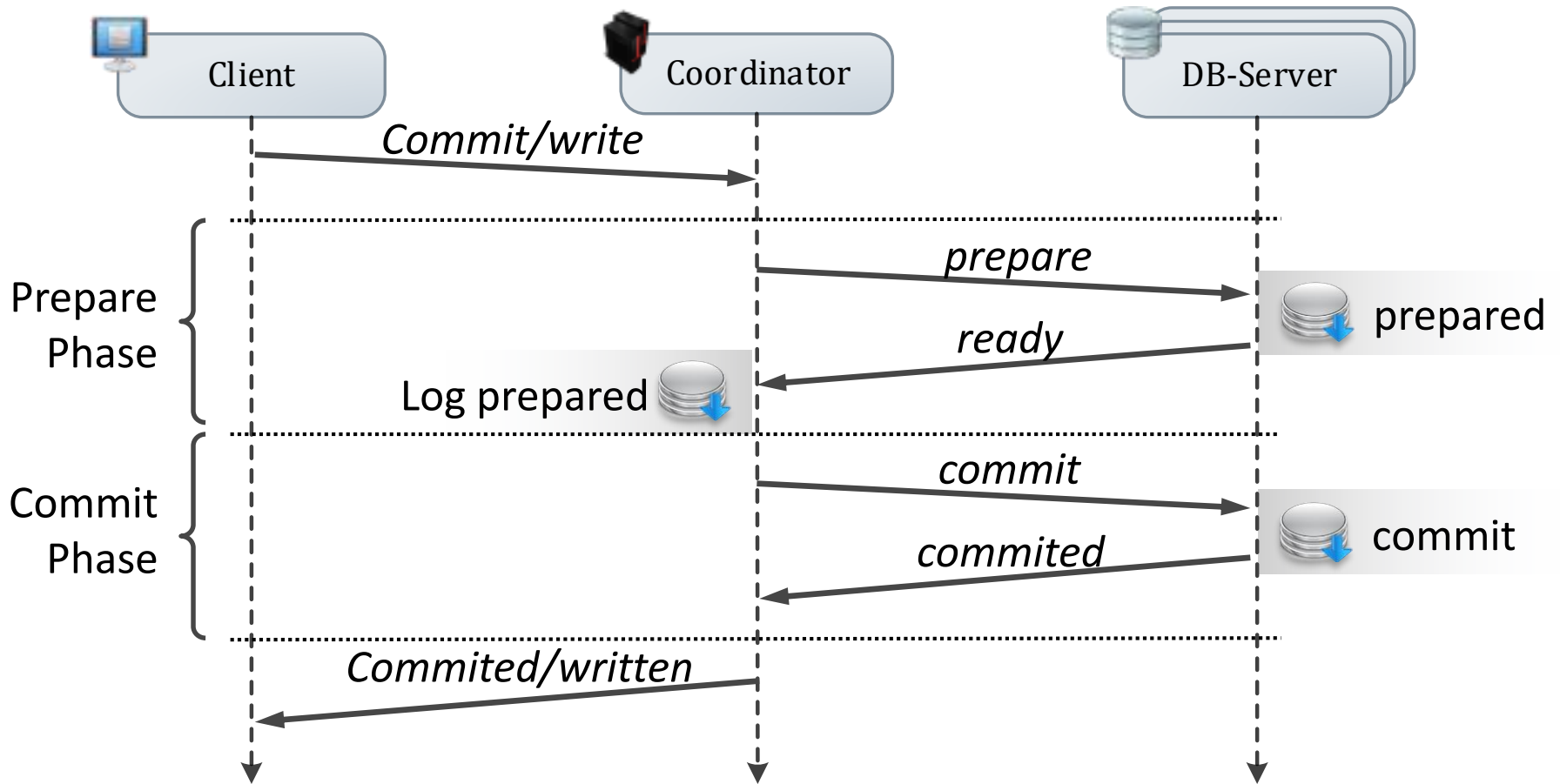
## Master-Slave (*Primary Copy*)

- 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 complicated coordination protocols (e.g. Paxos) or is inconsistent

# Synchronous Replication: 2PC



# Error scenarios (timeouts)

In **INITIAL**:

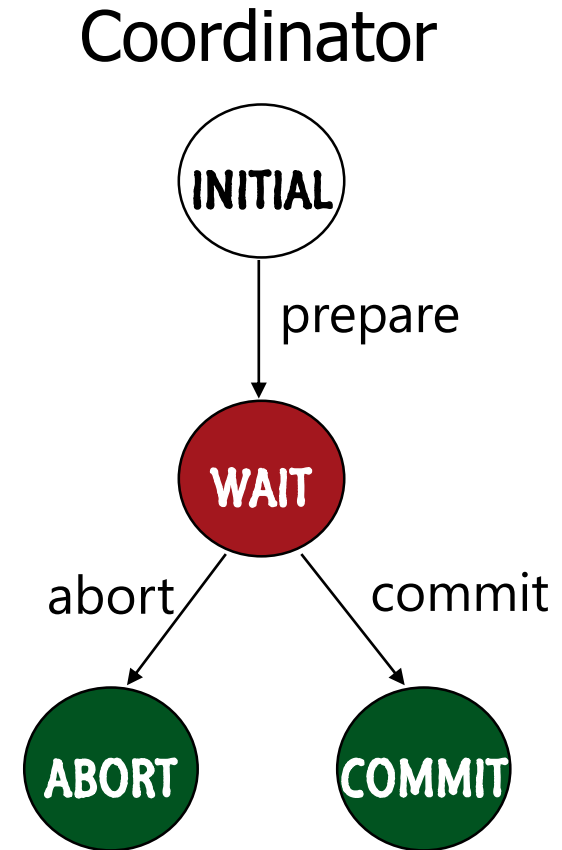
- No consequence

In **WAIT**:

- Abort

In **ABORT** oder **COMMIT**:

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



# Error scenarios (timeouts)

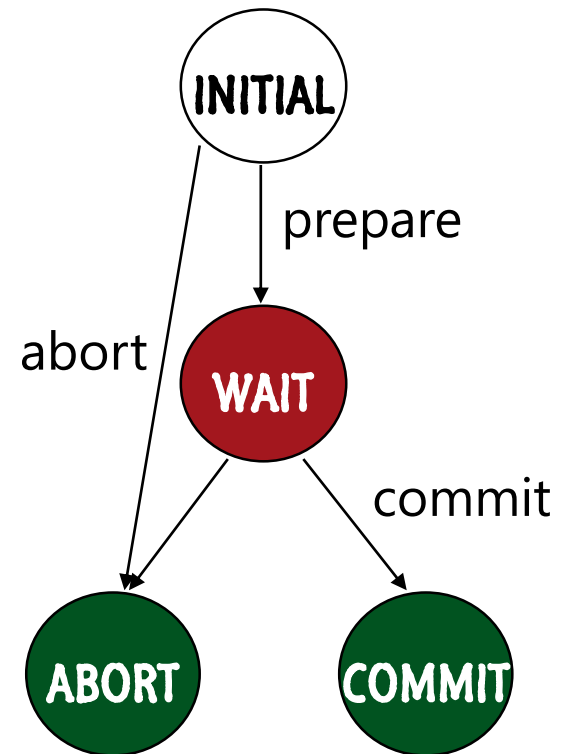
In **INITIAL**:

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

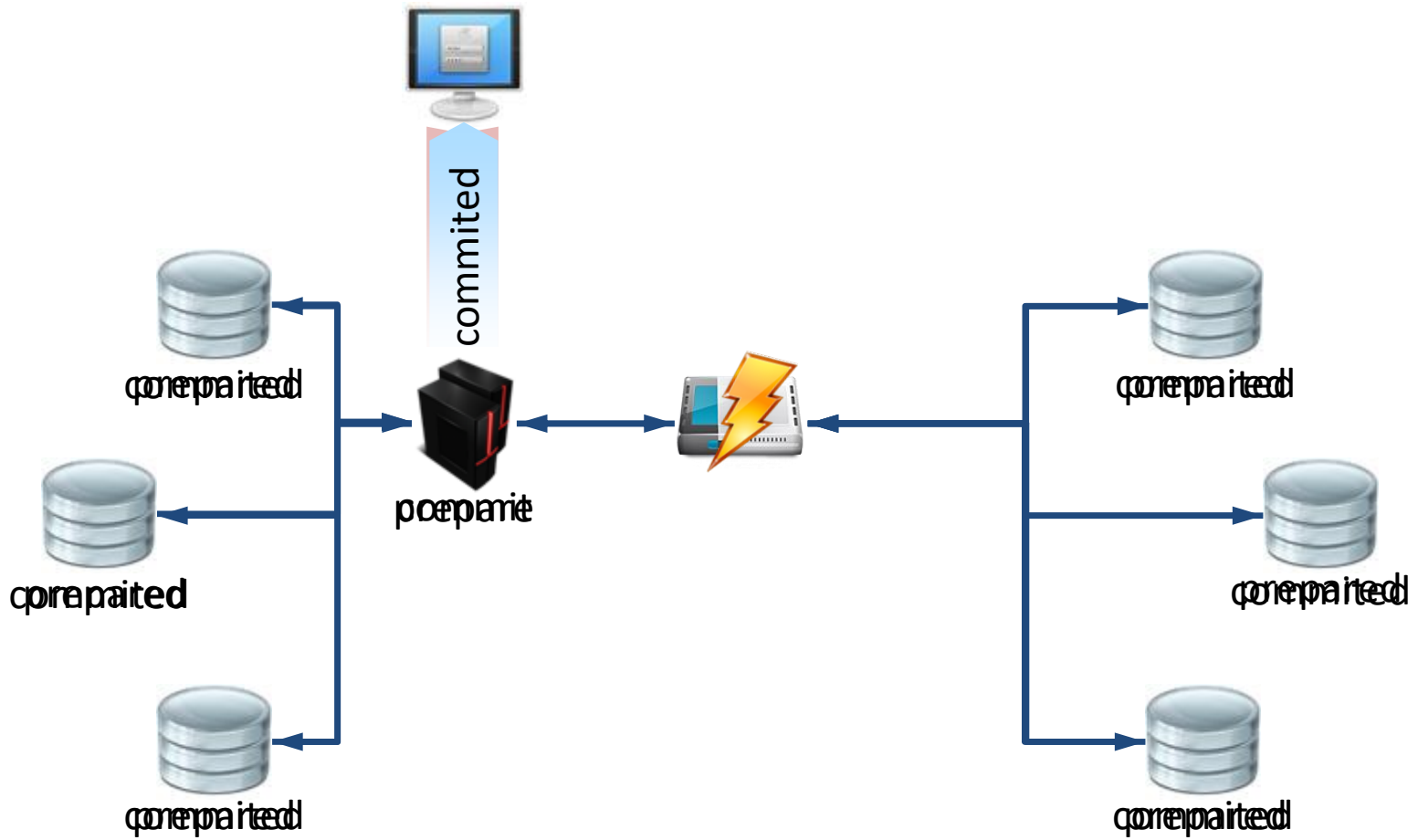
In **WAIT**:

- Wait for message from coordinator

Resource-Manager



# 2PC is not available



# Commit protocols

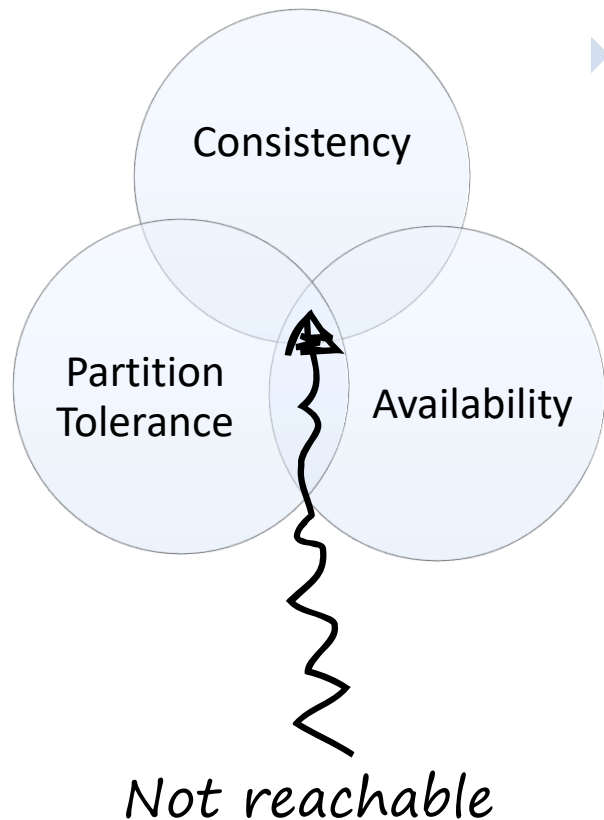
- ▶ Provide atomic propagation of writes or commits
- ▶ Most prominent implementation: 2-Phase-Commit

$R$  *Read-RMs*,  $W$  *Read/Write-RMs*,  $N=R+W$

Commit protocol	Messages	Property
<i>1-Phase-Commit</i>	$2N$	Not always possible
<i>Linear 2PC</i>	$2N+1$	Not parallel
<i>Hierarchical and normal 2PC</i>	$4N-2R$	Might block indefinitely
<i>3-Phase-Commit</i>	$6N-4R$	No consistency guarantee
<i>Paxos-Commit</i>	$3N+2F(N+1)+1$	$F$ failures tolerated
<i>Distributed 2PC</i>	$N^2$	No 2nd phase



# 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

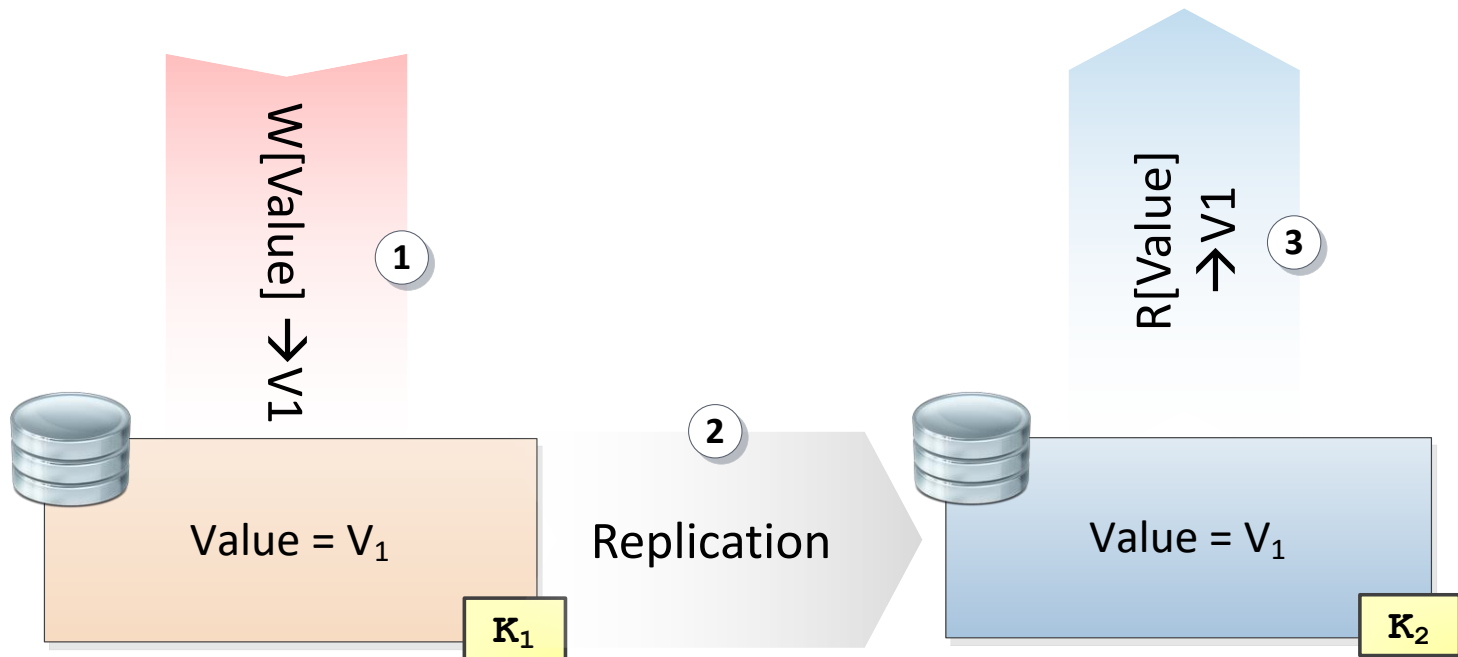
- ▶ Intuition for the impossibility of simultaneously achieving **C**, **A** and **P** at once:



*Value is replicated to two nodes*

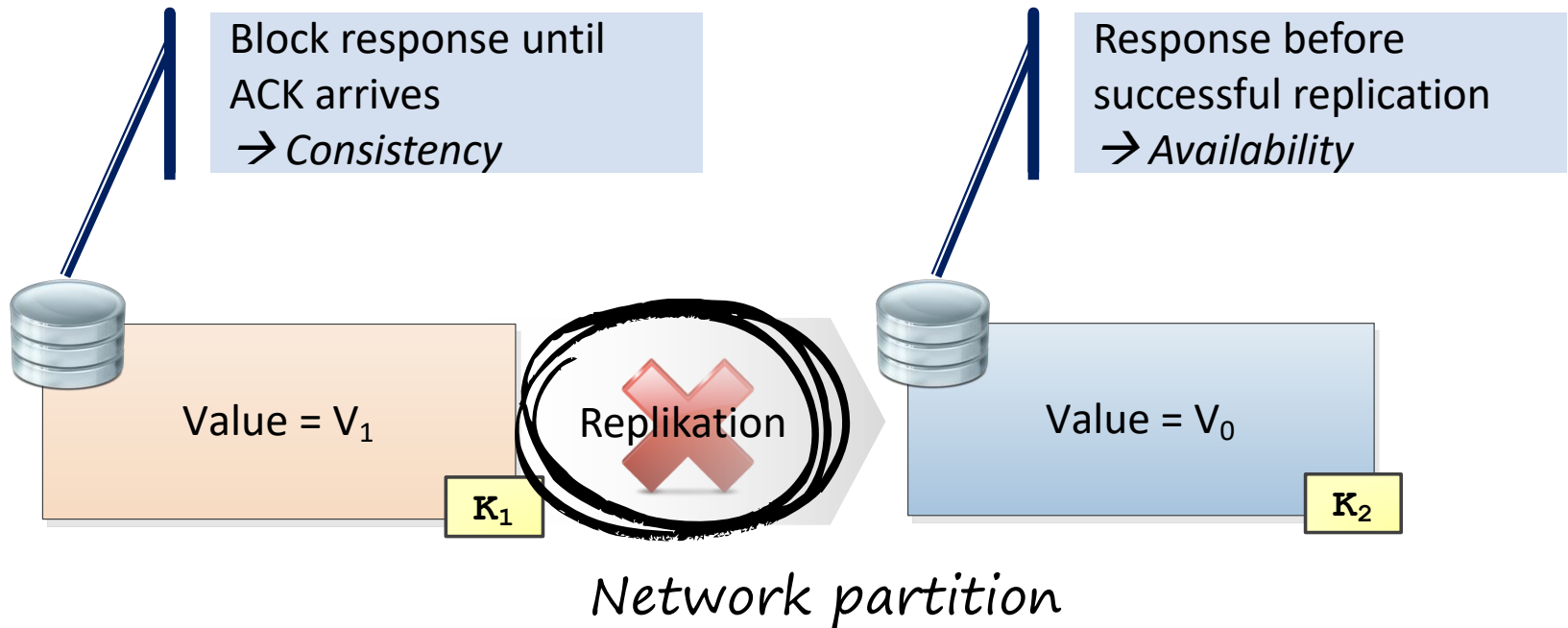
# CAP-Theorem: simplified proof

- ▶ Failure-free reading and writing:

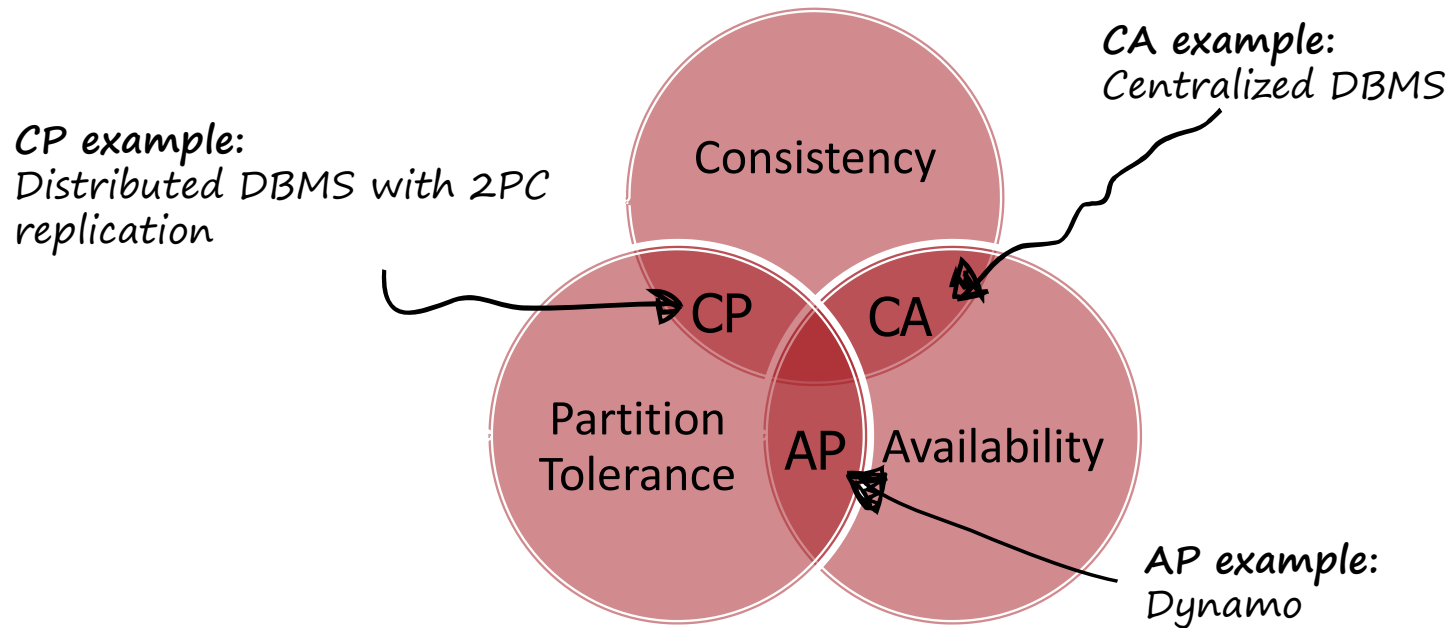


# CAP-Theorem: simplified proof

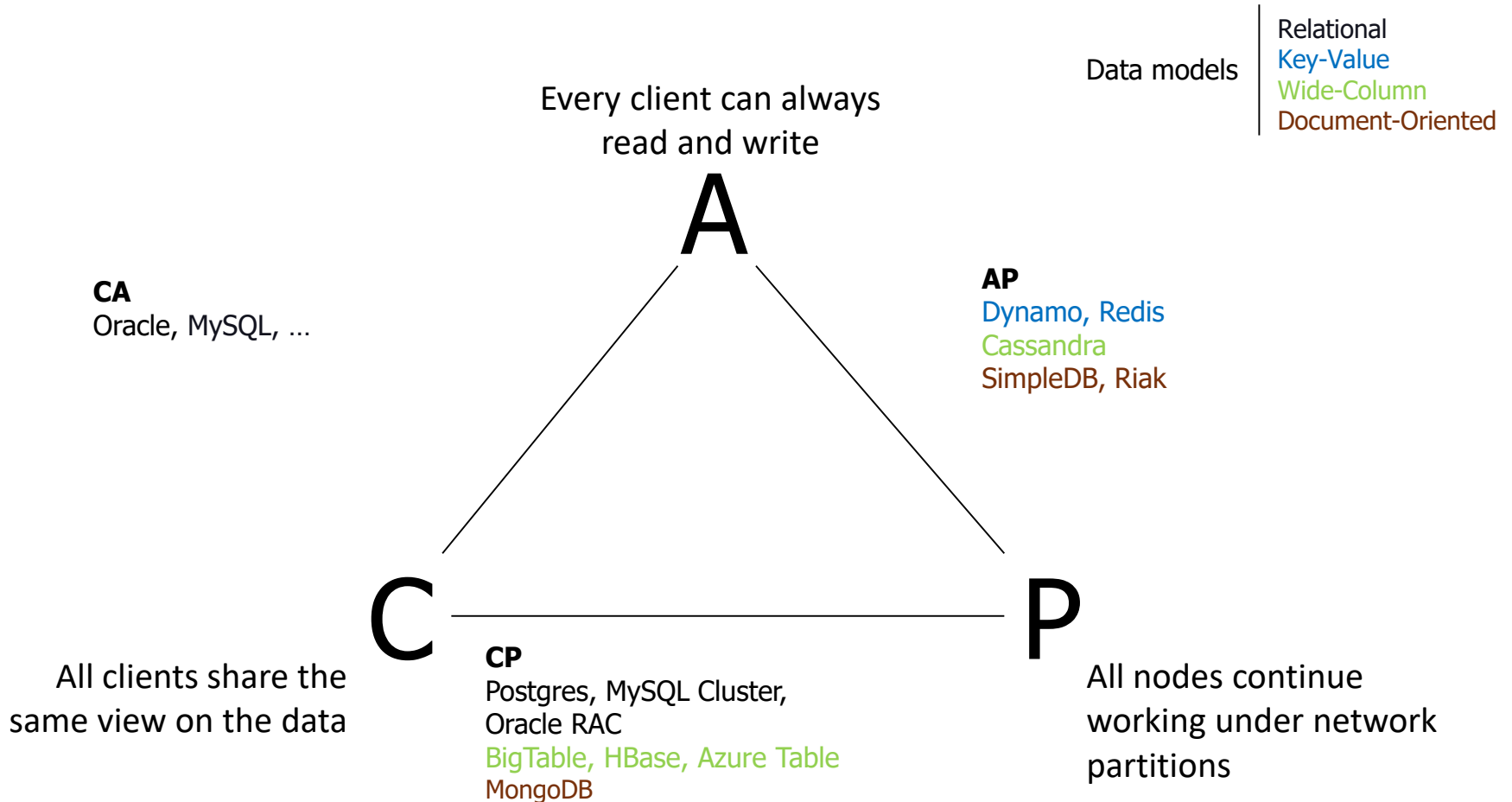
- ▶ **Problem:** when a network partition occurs, either consistency or availability have to be given up



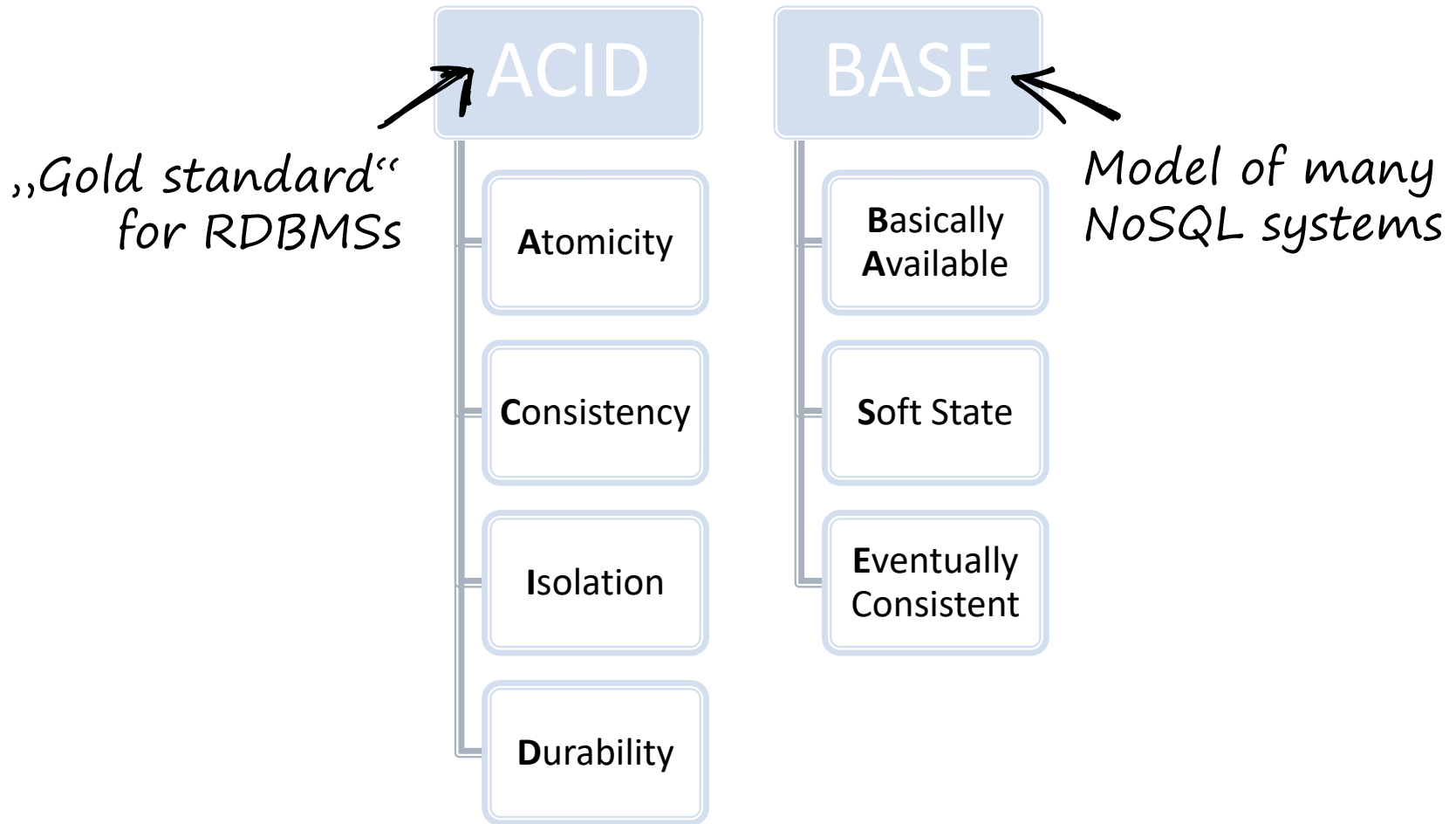
# CAP-Theorem: Wrap-up



# NoSQL triangle

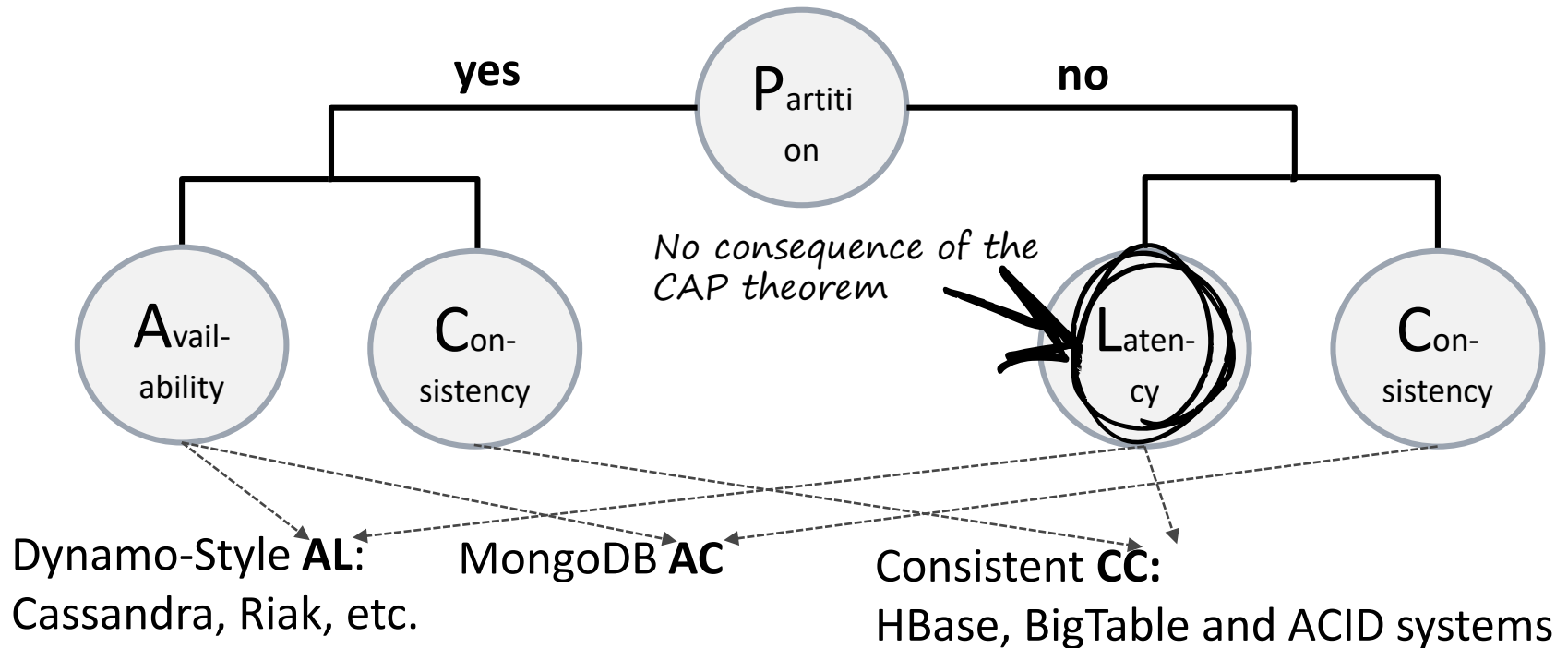


# ACID vs BASE



# PACELC – an alternative CAP formulation

- ▶ **Idea:** Classify system according to their heaviour during network partitions



Abadi, Daniel. "Consistency tradeoffs in modern distributed database system design: CAP is only part of the story."



# Negative Results

## In Distributed Computing

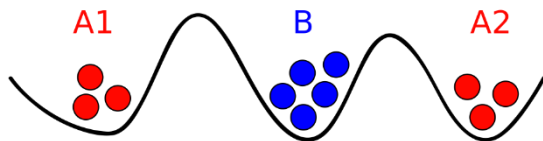
### *Asynchronous Network, Unreliable Channel*

#### 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

# Negative Results

## Consensus Algorithms

- ▶ 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, Cassandra Lightweight Transactions)
  - **Raft** (e.g. etcd service)
  - Zookeeper Atomic Broadcast (**ZAB**)



**Safety  
Properties**




**Liveness  
Property**

# Negative Results

## Correctness/Serializability

Distributed ACID and availability are incompatible:

Write A=1  
Read B



$w_1(a = 1) r_1(b = \perp)$

Write B=1  
Read A



$w_2(b = 1) r_2(a = \perp)$

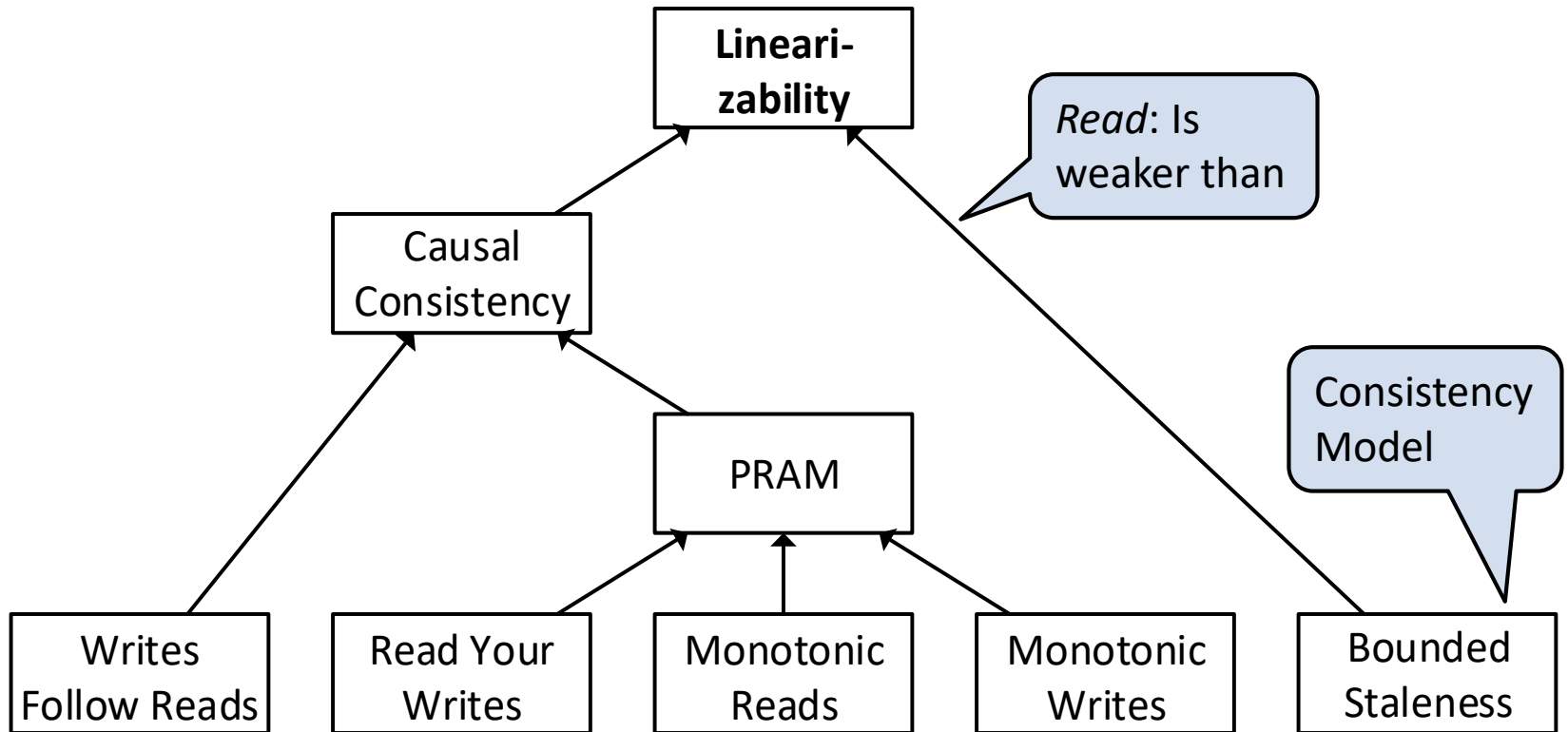
- ▶ Weaker isolation levels are possible:
  - RAMP Transactions (P. Bailis, A. Fekete, A. Ghodsi, J. M. Hellerstein, und I. Stoica, „Scalable Atomic Visibility with RAMP Transactions“, SIGMOD 2014)
- ▶ **Consequence:** trade-offs are important

# Typical Trade-offs

Read Performance	Write Performance
Latency	Durability
Synchronous replication	Asynchronous replication
Row-based	Column-based
Transactions	Availability
REST	RPC
Commodity servers	High-end hardware
Normalisation	Denormalisation
Schemas	Schemafreeness

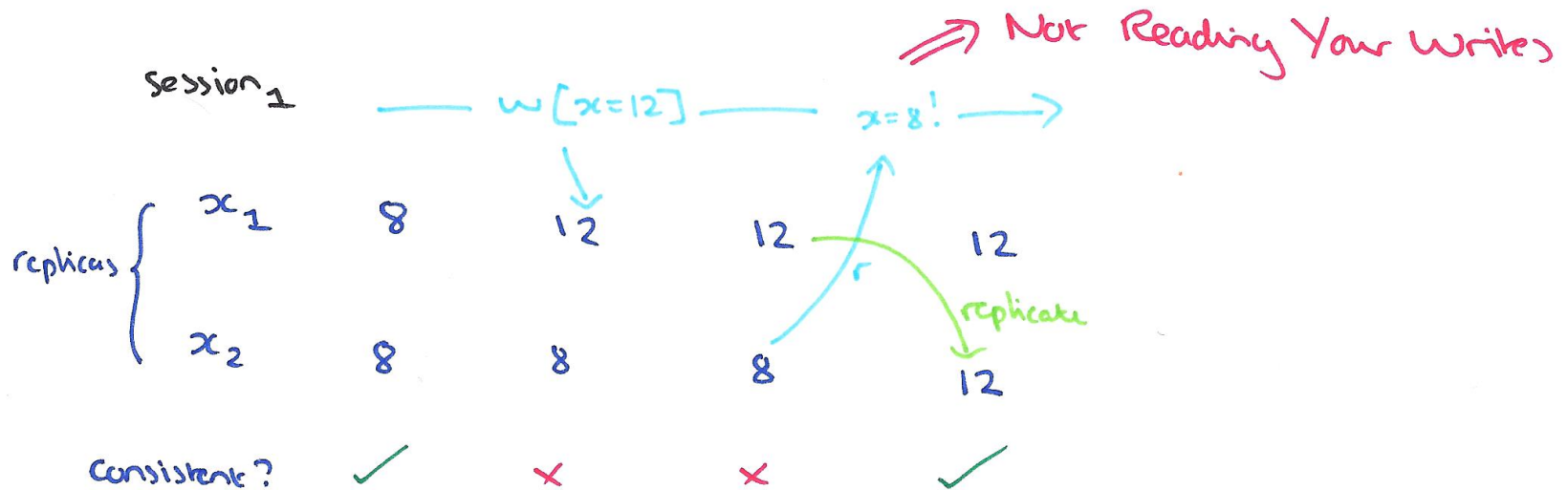
# Client-Centric Consistency Models

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



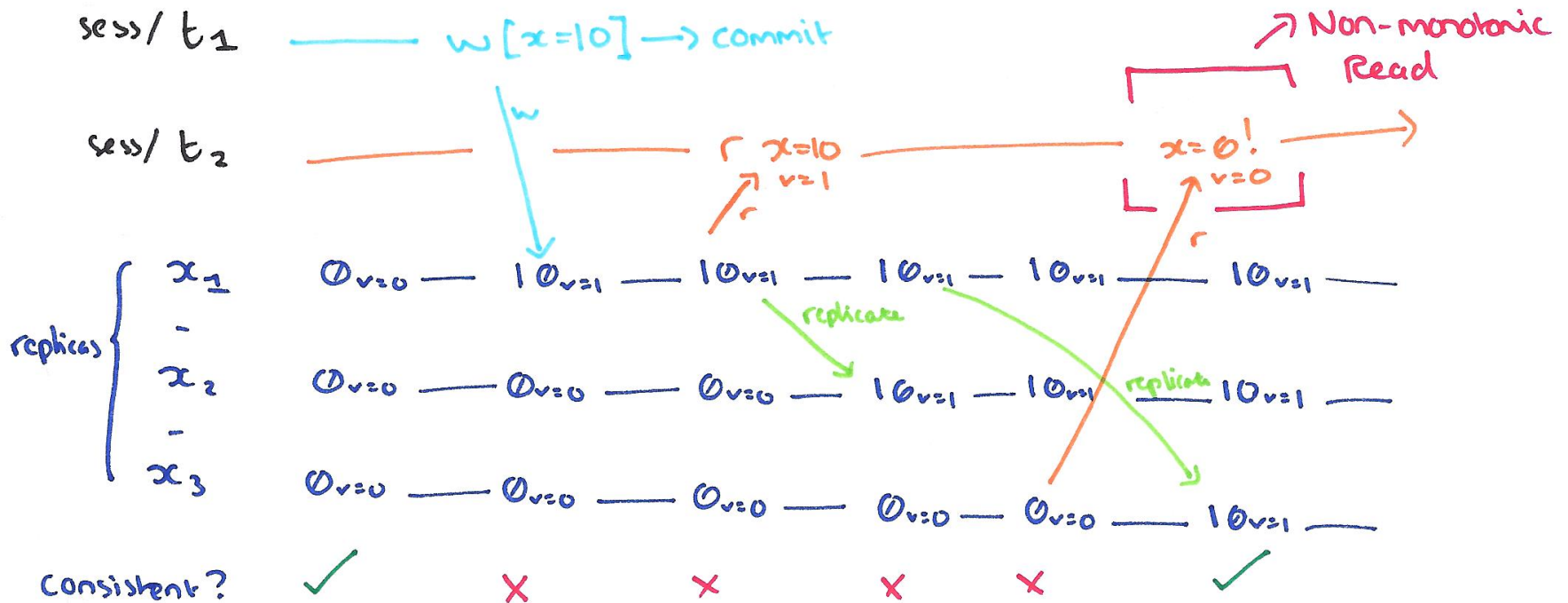
# Read Your Writes (RYW)

**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.



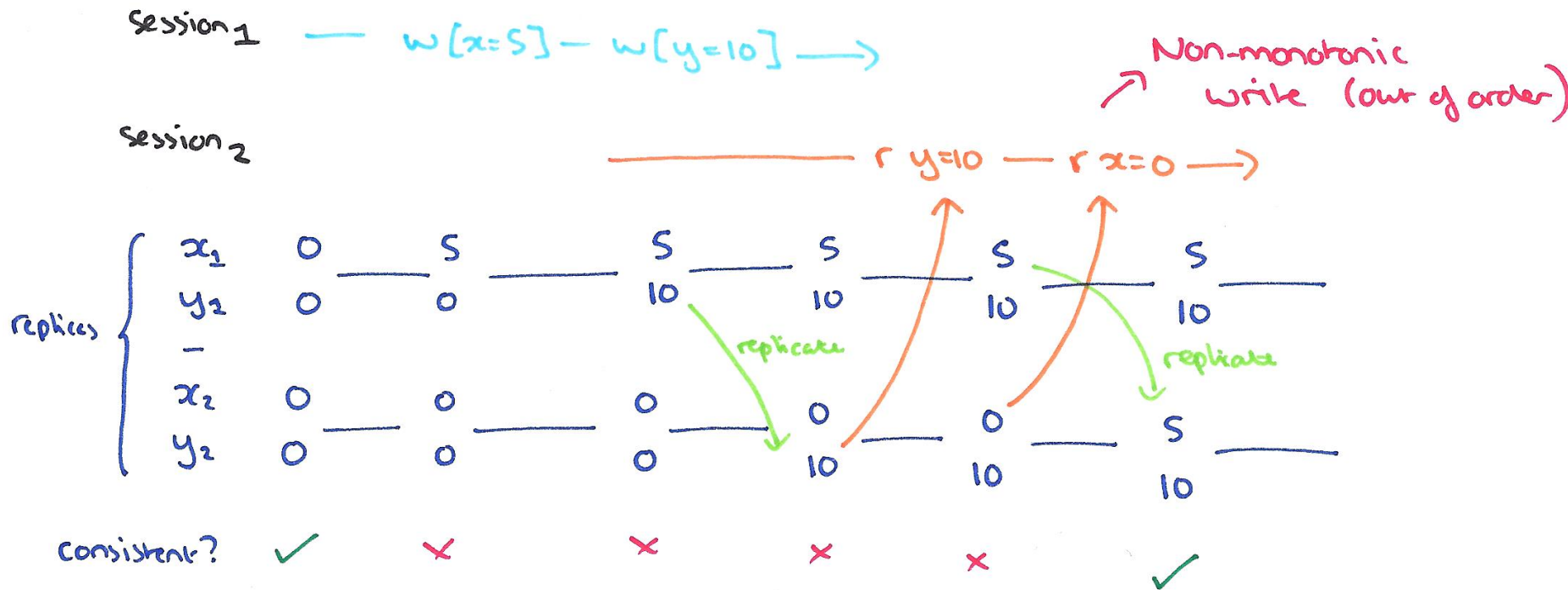
# Monotonic Reads (MR)

**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)

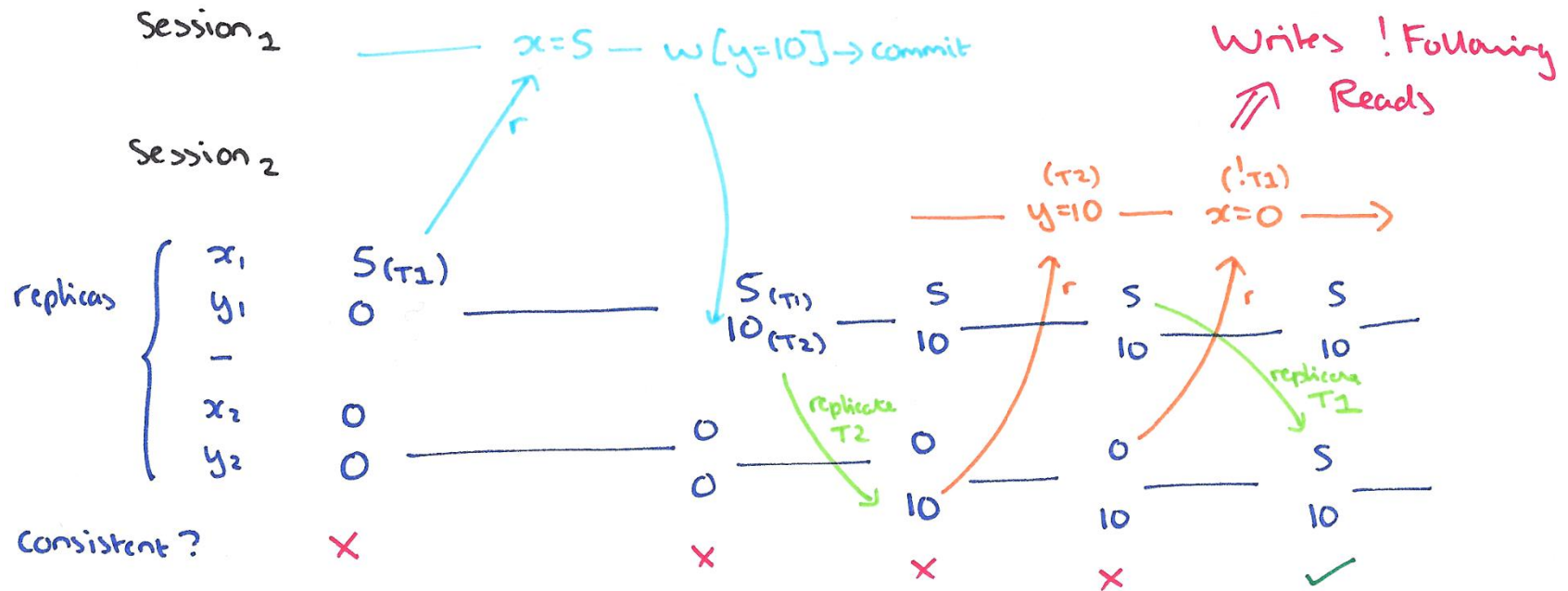
**Definition:** Once a user has written a new value for a data item in a session, any previous write has to be processed before the current one. I.e., the order of writes inside the session is strictly maintained.





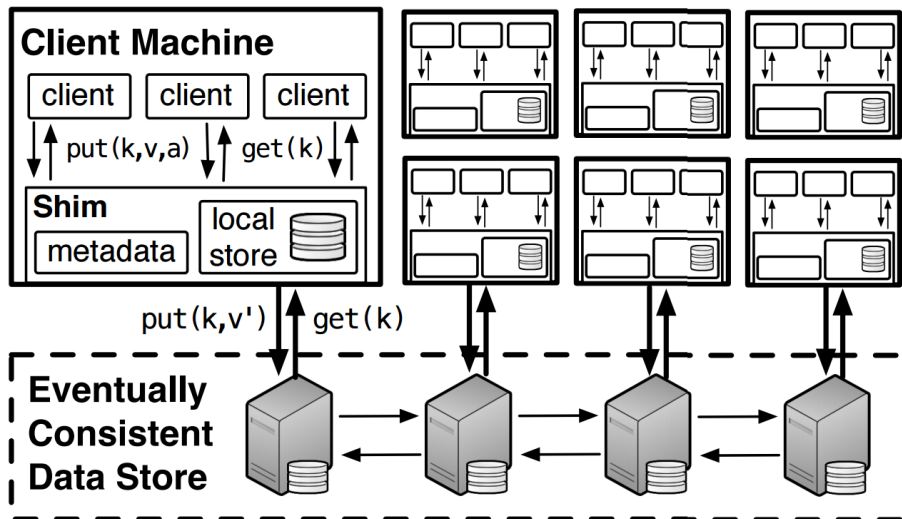
# Writes Follow Reads (WFR)

**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.



# PRAM and Causal Consistency

- ▶ Combinations of previous session consistency guarantees
  - PRAM = MR + MW + RYW
  - Causal Consistency = PRAM + WFR
- ▶ All consistency level up to causal consistency can be guaranteed with **high availability**
- ▶ Example: Bolt-on causal consistency



# Bounded Staleness

- ▶ Either **time-based**:

**t-Visibility ( $\Delta$ -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**:

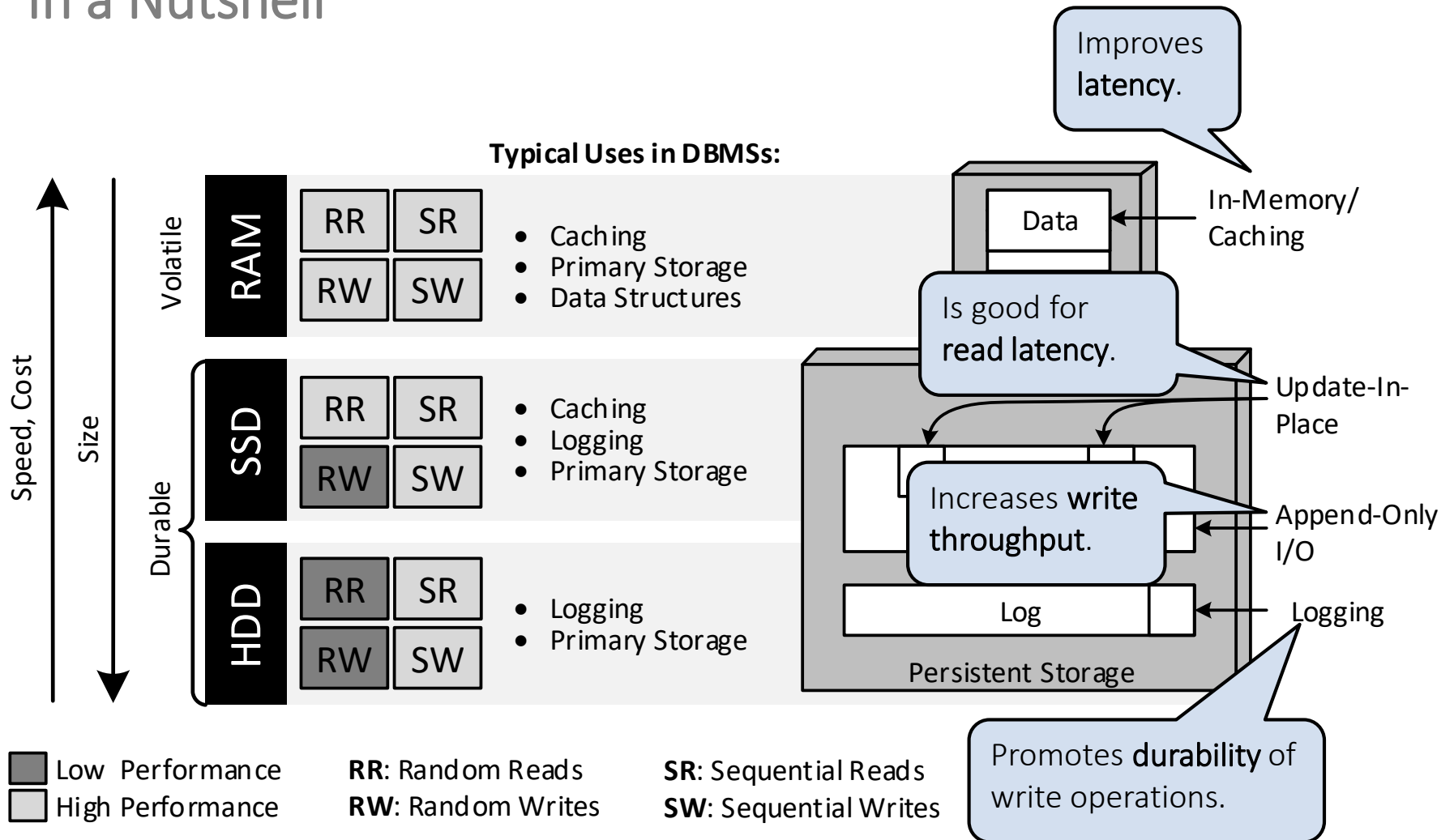
**k-Staleness:** the inconsistency window comprises at most  $k$  versions; that is, lags at most  $k$  versions behind the most recent version.

- ▶ Both are *not* achievable with high availability



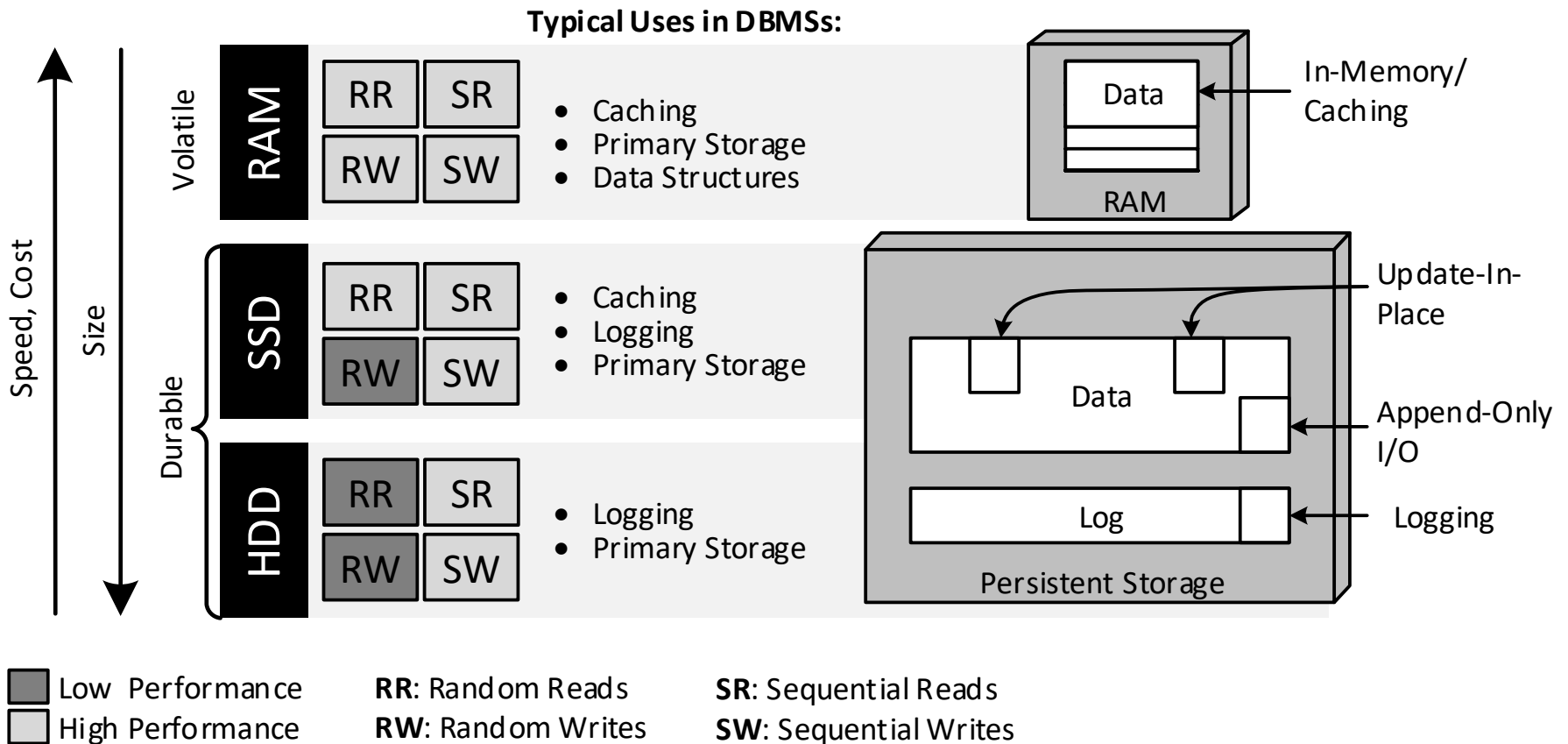
# NoSQL Storage Management

## In a Nutshell



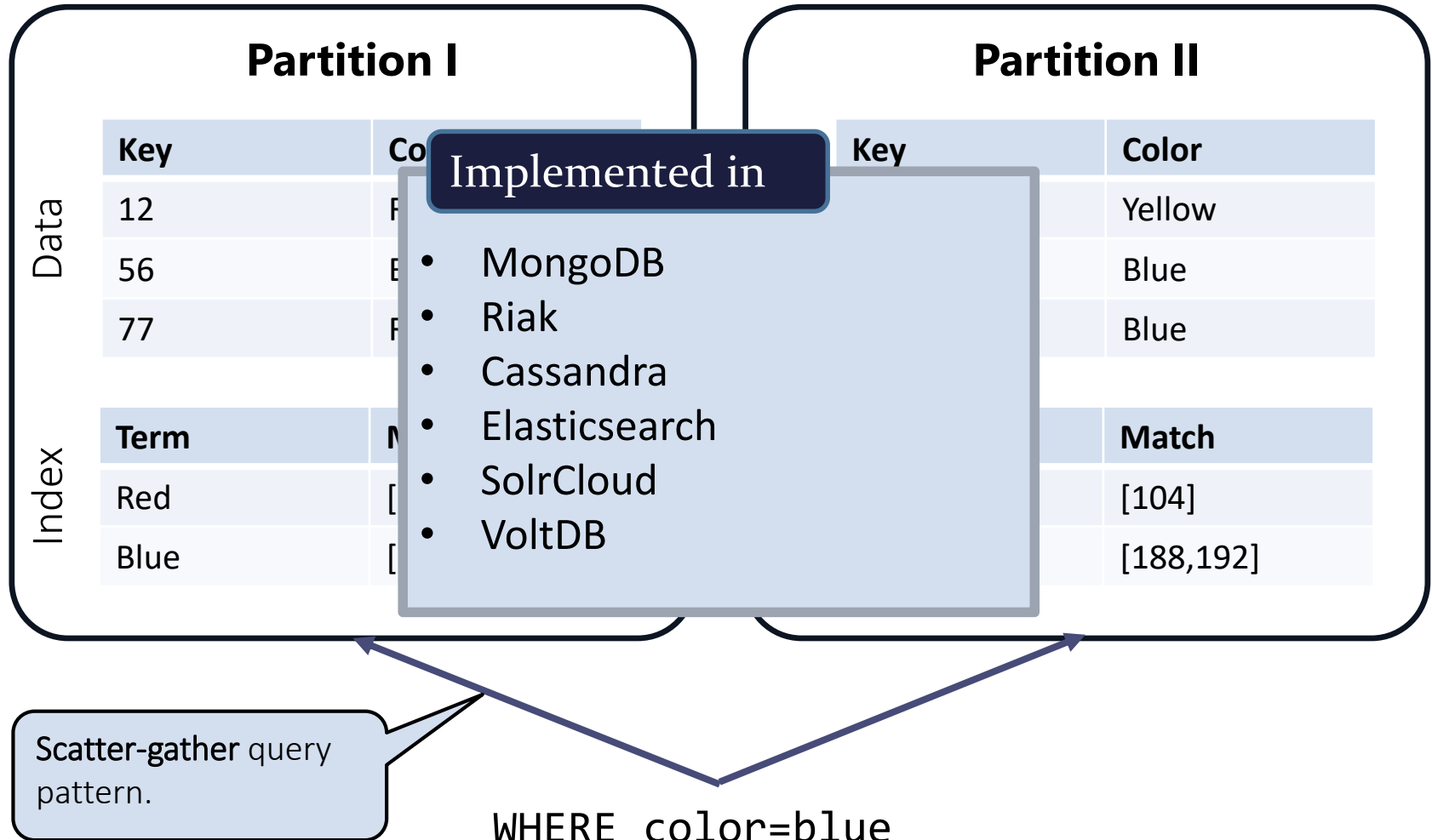
# NoSQL Storage Management

## In a Nutshell



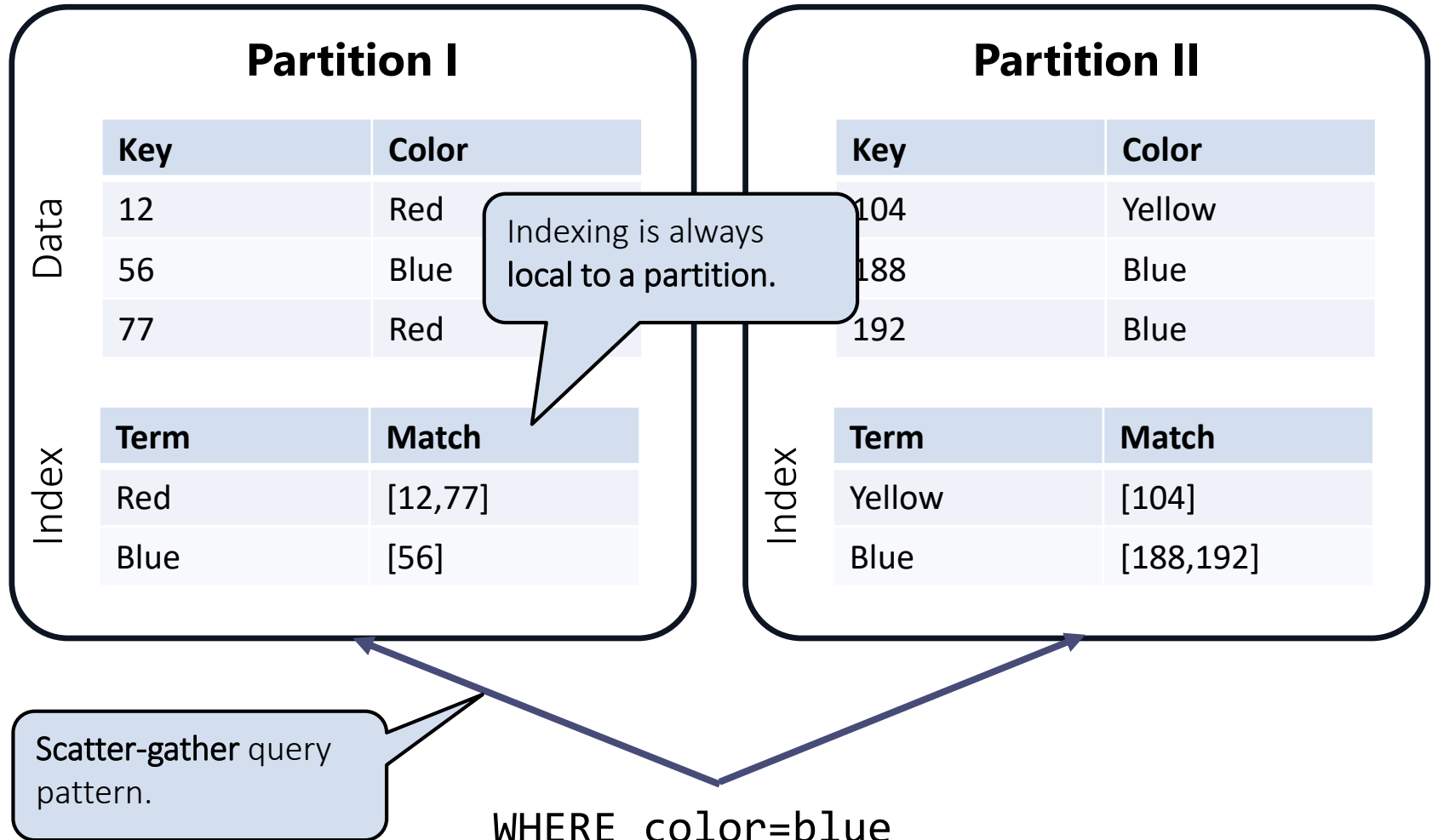
# Local Secondary Indexing

## Partitioning By Document



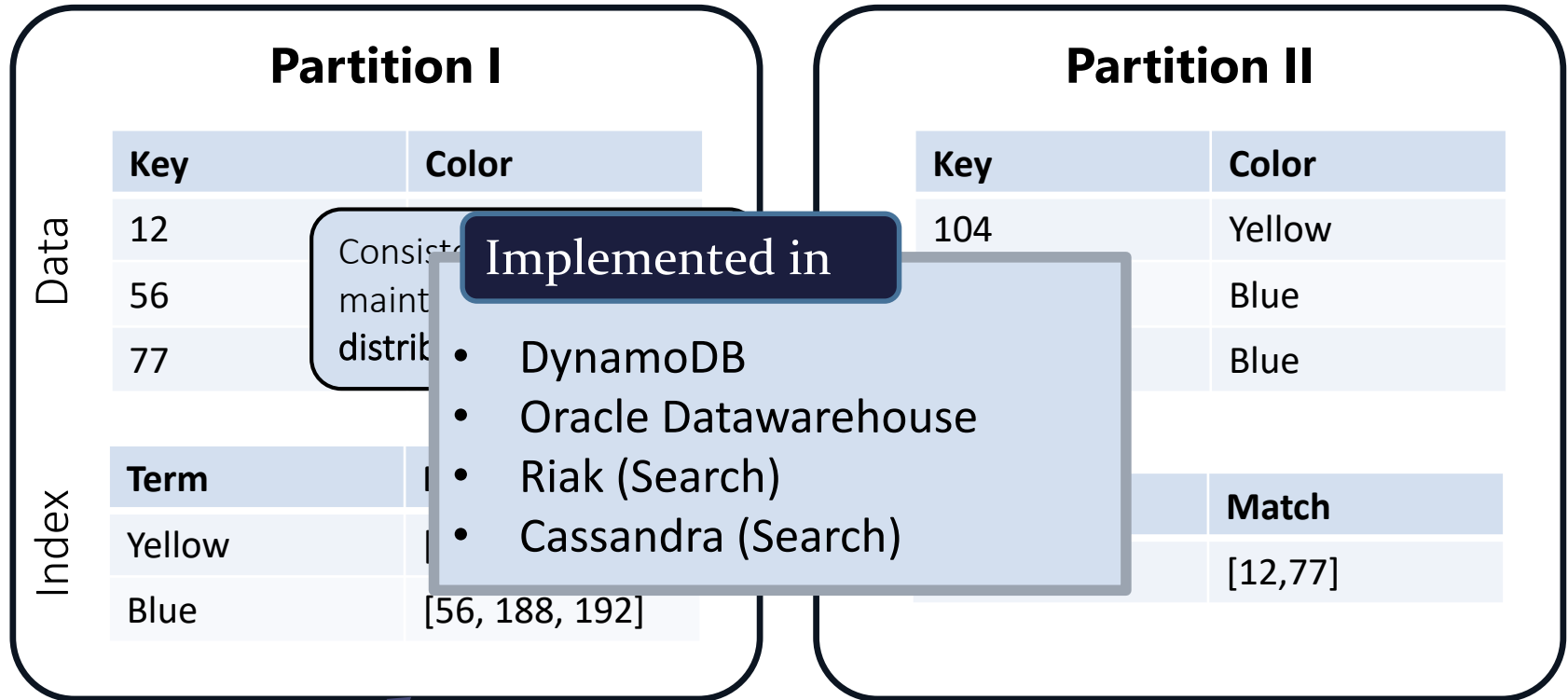
# Local Secondary Indexing

## Partitioning By Document



# Global Secondary Indexing

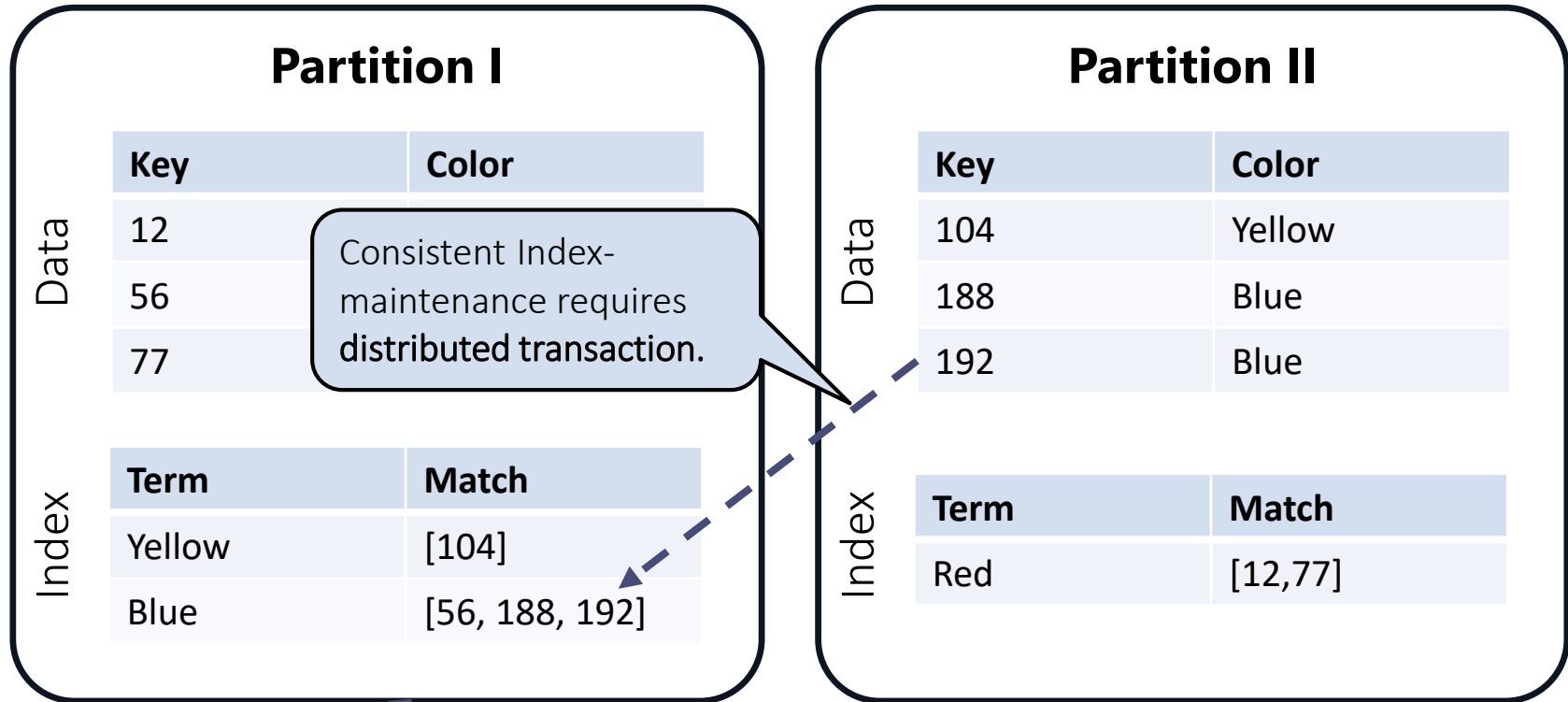
## Partitioning By Term





# Global Secondary Indexing

## Partitioning By Term



Targeted Query

WHERE color=blue

# Wrap-up



- ▶ High data volumes, unstructured sources and new kinds of applications triggered BigData and 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
  - **BASE** (Basically available, soft-state, eventually consistent) paradigm as an alternative to ACID
- ▶ **2-Phase-Commit (2PC)**: popular protocol for atomic commitment; without availability guarantee
- ▶ **Consistency** can be relaxed in various ways

# Outline



Foundations: Big Data,  
Scalability, Availability



The 4 Classes of NoSQL  
Databases



NoSQL Examples: concrete  
Architectures, Systems, APIs

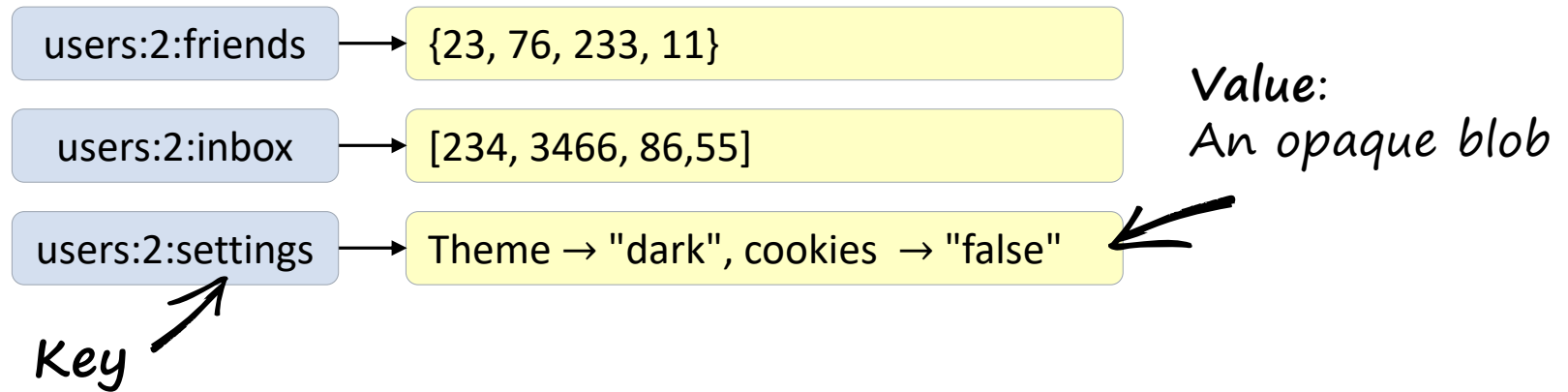


Cloud Databases

- Key-Value stores
- Wide-Column stores
- Document stores
- Graph databases
- Other classes

# Key-Value Stores

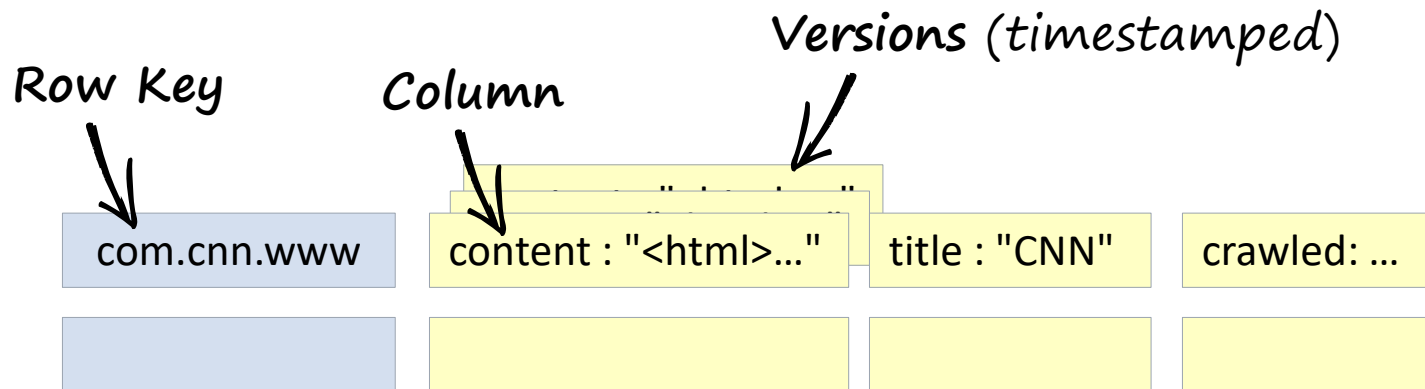
- ▶ **Data model:** (key) -> value
- ▶ **Interface:** CRUD (Create, Read, Update, Delete)



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

# Wide-Column Stores

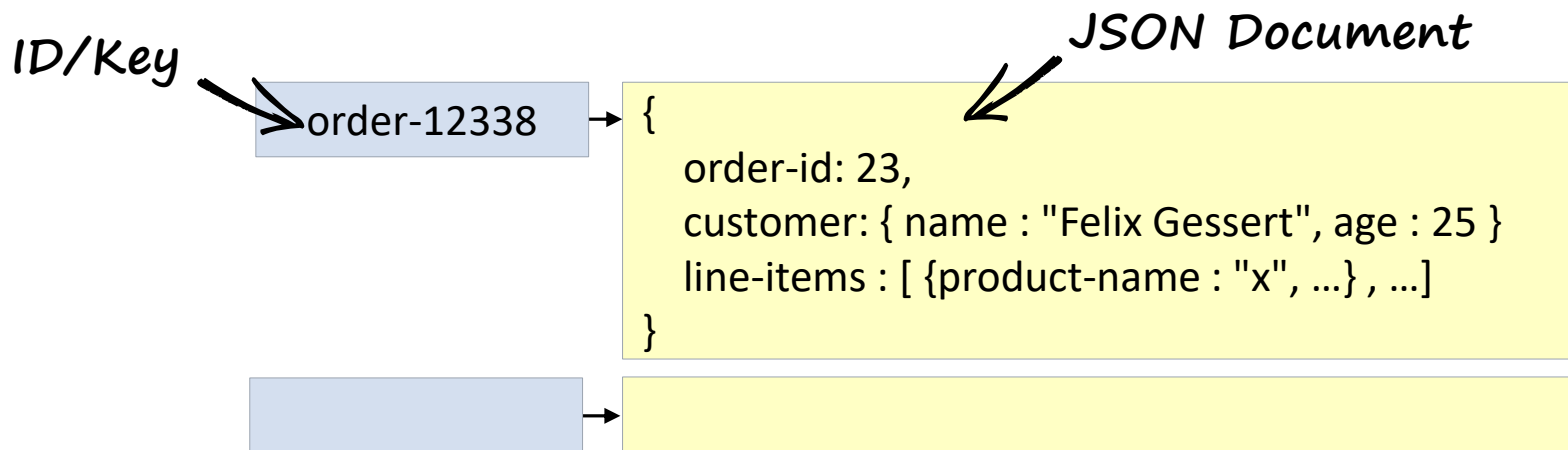
- ▶ **Data model:** (rowkey, column, timestamp) -> value
- ▶ **Interface:** CRUD, Scan



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

# Document Stores

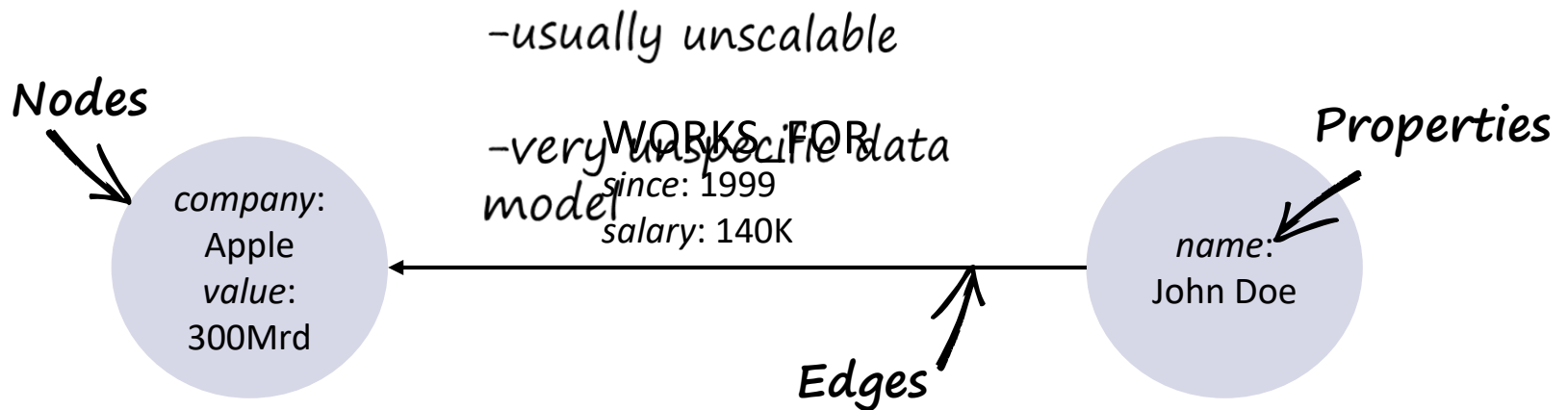
- ▶ **Data model:** (collection, key) -> document
- ▶ **Interface:** CRUD, Querys, Map-Reduce



- ▶ Examples: CouchDB (AP), Amazon SimpleDB (AP), MongoDB (CP)

# Graph Databases

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

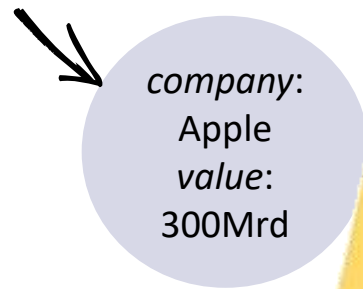


- ▶ Examples: Neo4j (CA), InfiniteGraph (CA)

# Graph Databases

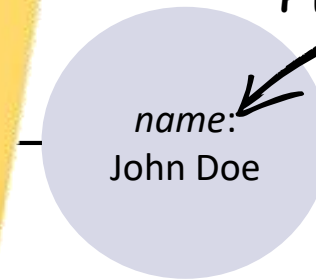
- ▶ **Data model:**  $G = (V, E)$ : Graph-Property Modell
- ▶ **Interface:** Traversal, queries, transactions

*Nodes*



*-usually unscalable*  
*-very unspecific data model*

*Properties*

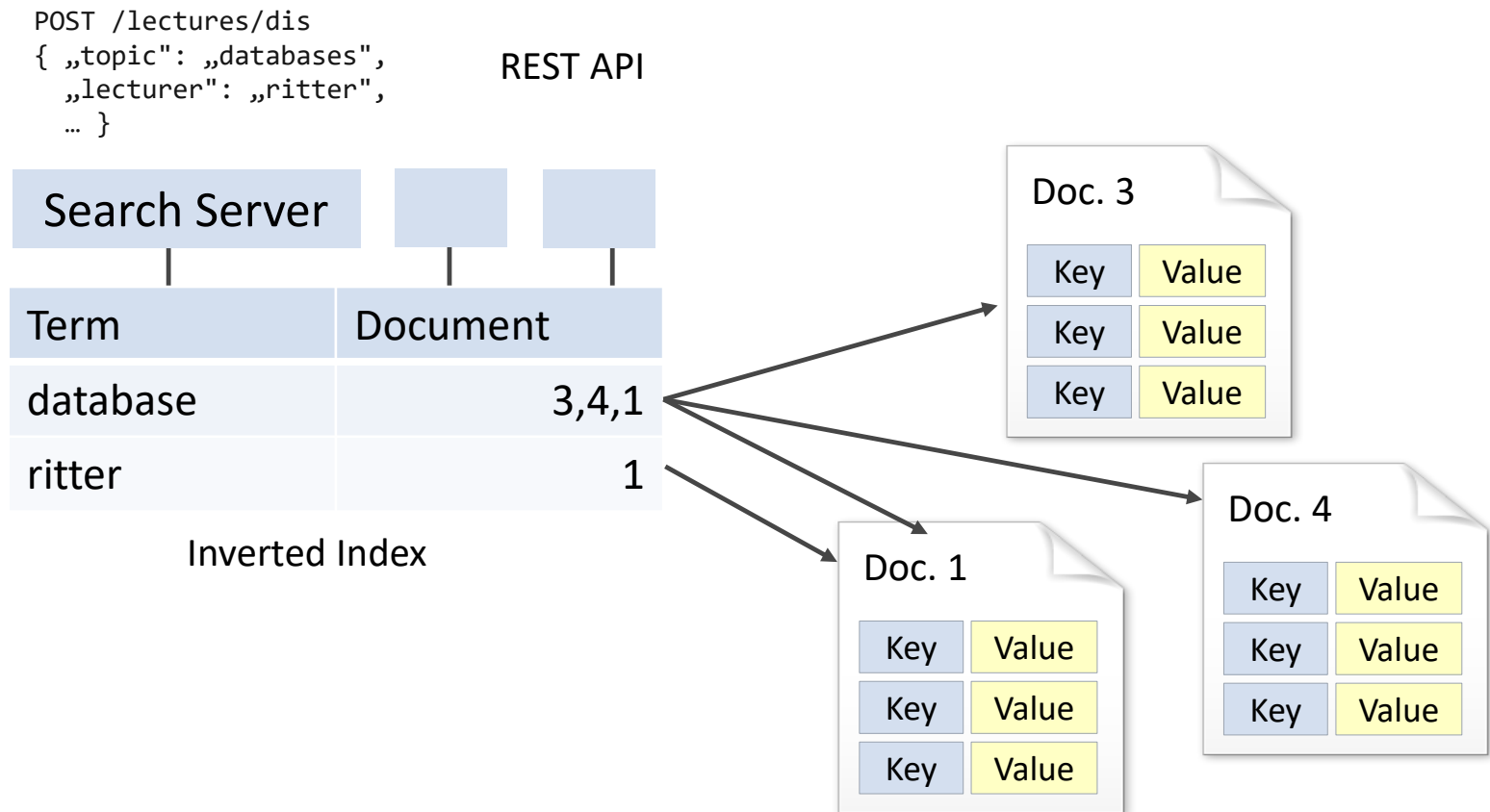


- ▶ **Examples:** Neo4j (CA), InfiniteGraph (CA)



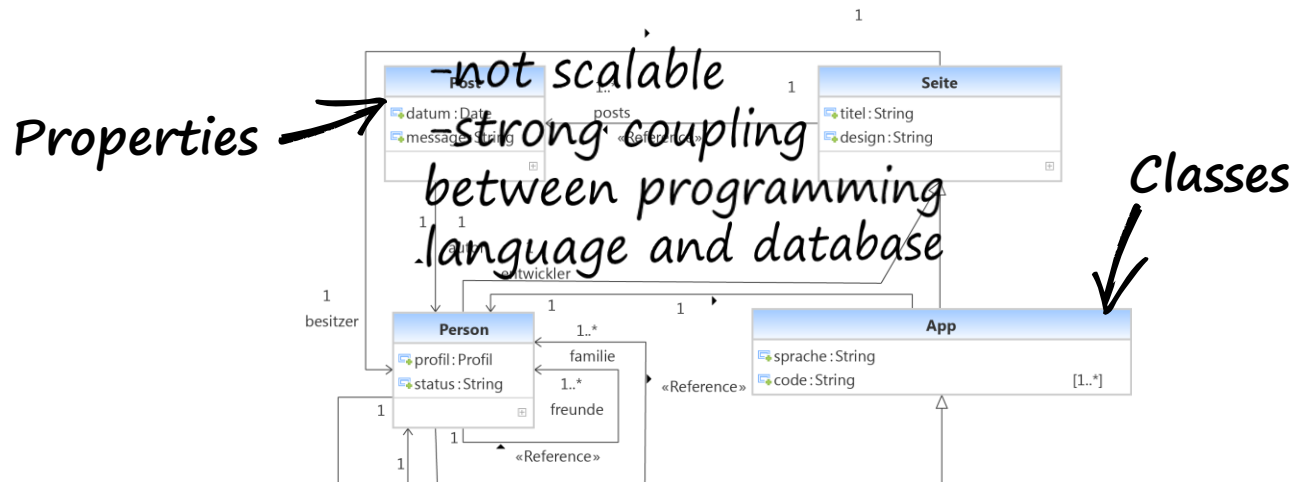
# Search Platforms

- ▶ **Data model:** vectorspace model, docs + metadata
- ▶ Examples: Solr, ElasticSearch



# Object-oriented Databases

- ▶ **Data model:** Classes, objects, relations (references)
- ▶ **Interface:** CRUD, queries, transactions



- ▶ Examples: Versant (CA), db4o (CA), Objectivity (CA)

# Object-oriented Databases

- ▶ **Data model:** Classes, objects, relations (references)
- ▶ **Interface:** CRUD

Properties

-not scalable  
-strong coupling  
between programming  
language and database

Classes

[1..\*]

- ▶ **Examples:** Versant (CA), db4o (CA), Objectivity (CA)

# XML databases, RDF Stores

- ▶ **Data model:** XML, RDF
- ▶ **Interface:** CRUD, queries (XPath, XQuerys, SPARQL), transactions (some)
- ▶ **Examples:** MarkLogic (CA), AllegroGraph (CA)
  - not scalable
  - not widely used
  - specialized data model

# XML databases, RDF Stores

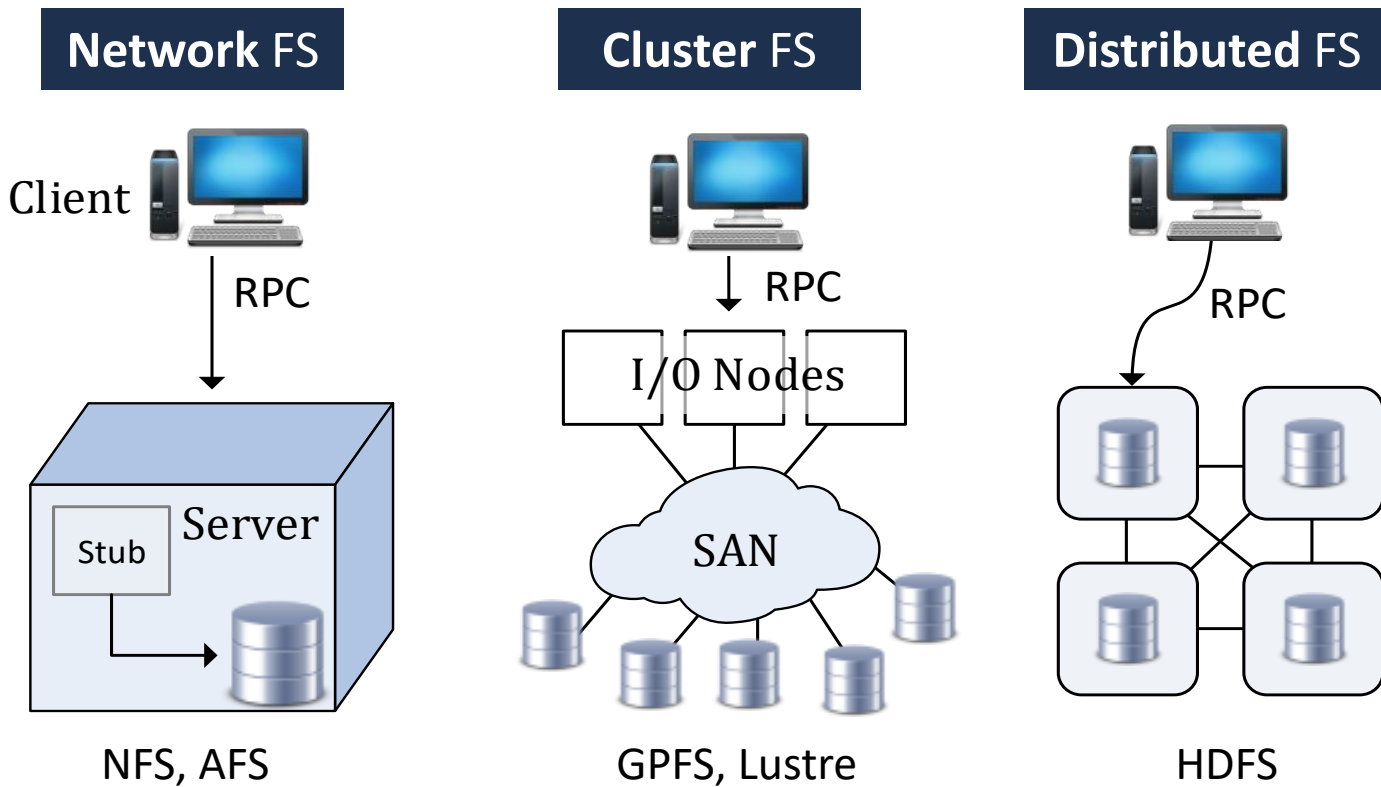
- ▶ Data model: XML, RDF
- ▶ Interface: CRUD, XQuery, XQL, XQL, XQuery, SPARQL), transactions (s
- ▶ Examples: Ma

*-not scalable  
-not widely used  
-specialized data  
model*

aph (CA)

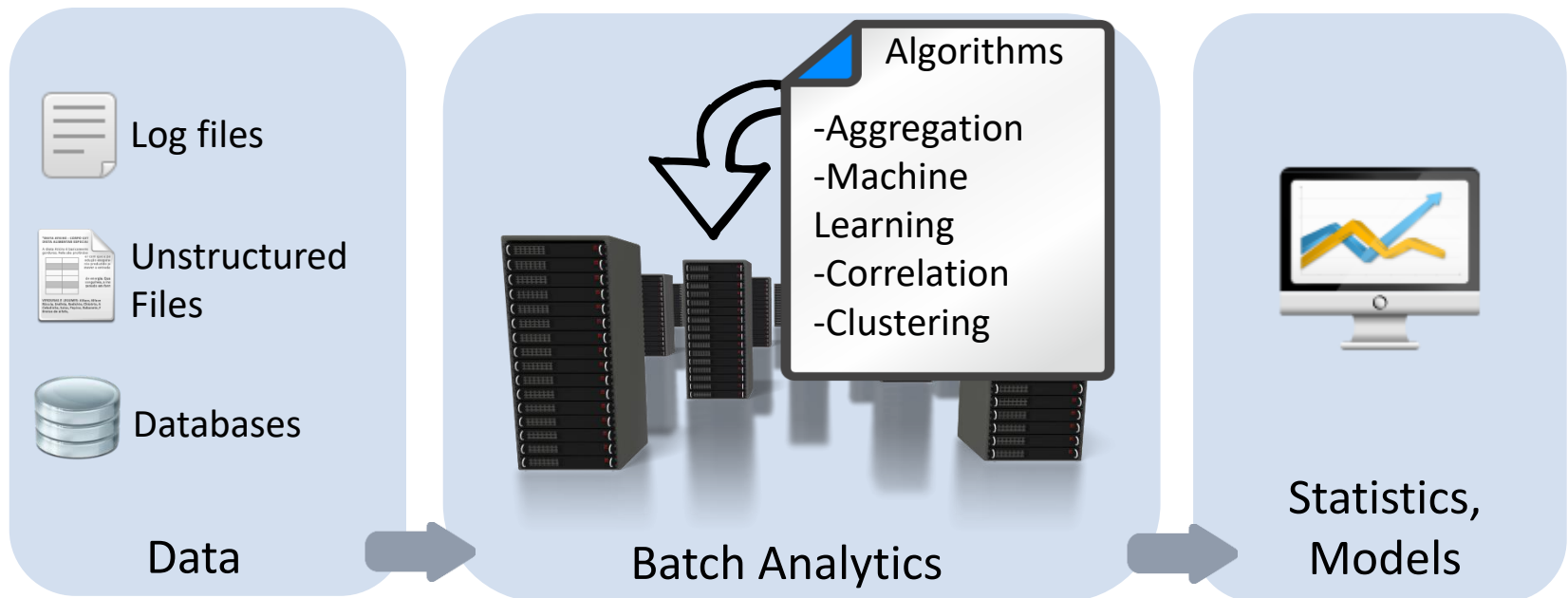
# Distributed File System

- ▶ Data model: files + folders



# Big Data Frameworks

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



# Wrap-up



- ▶ 4 core NoSQL classes
  - **Key-Value** Stores: store opaque key-value pairs
  - **Document** Stores: store nested, rich, schema-free documents
  - **Wide-Column** Stores: extensible table data model
  - **Graph Databases**: graph-property-model (vertices and edges)
- ▶ Other NoSQL-related systems: Object-oriented databases, Search platforms, XML databases, Big Data Frameworks, Distributed File Systems



# Outline



Foundations: Big Data,  
Scalability, Availability



The 4 Classes of NoSQL  
Databases



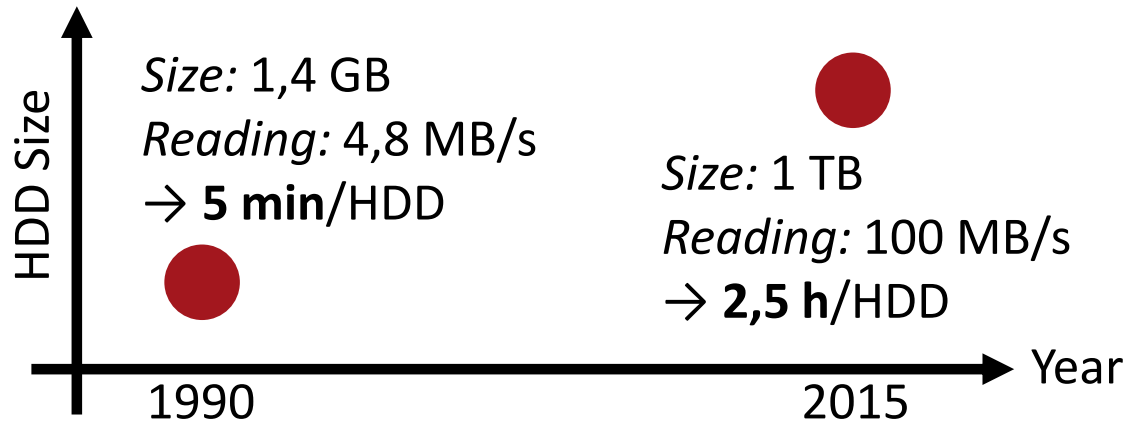
NoSQL Examples: concrete  
Architectures, Systems, APIs



Cloud Databases

- MapReduce (Hadoop)
- Dynamo (Riak)
- BigTable (HBase)
- MongoDB
- Others

# Hadoop Distributed FS (CP)

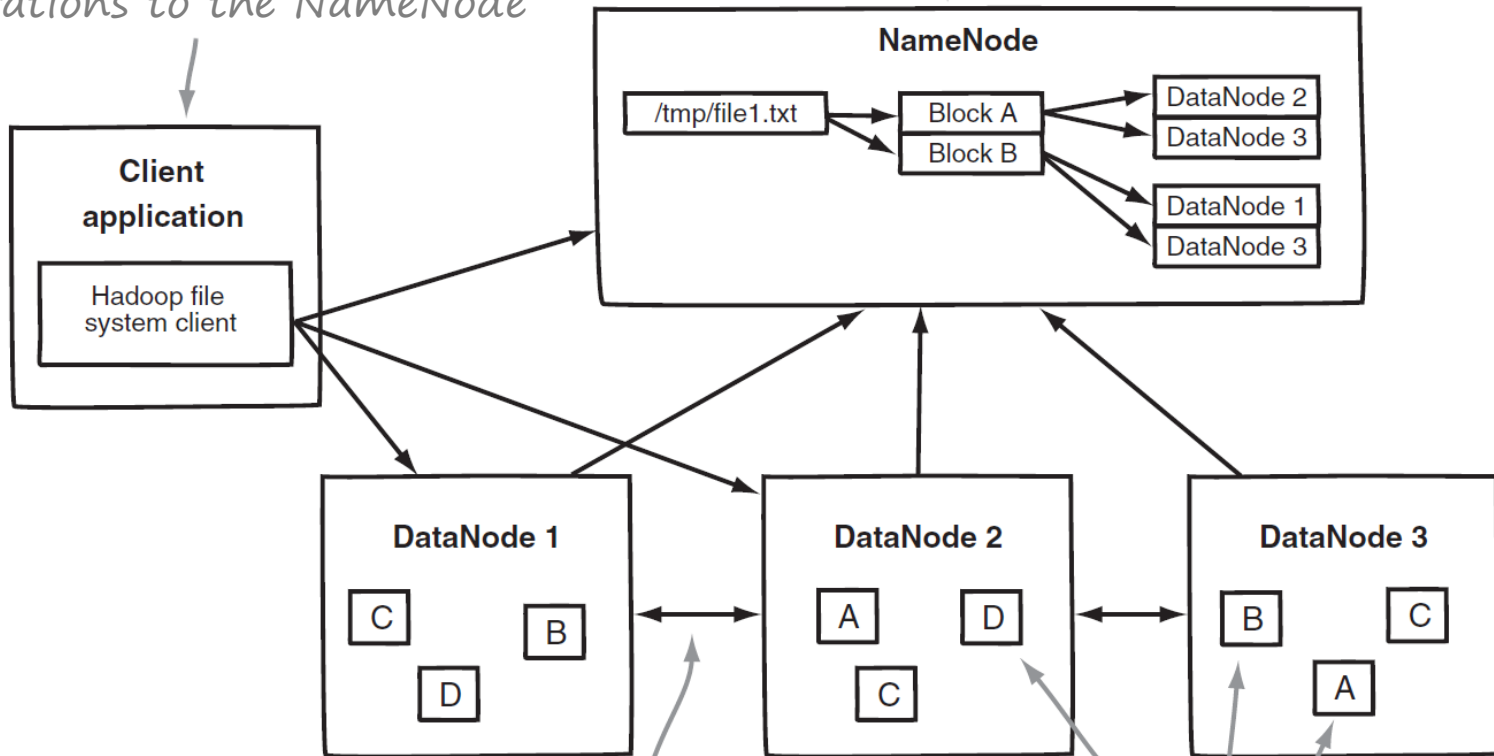


HDFS
Model:
File System
License:
Apache 2
Written in:
Java

- ▶ Modelled after: Googles GFS (2003)
- ▶ **Master-Slave** Replication
  - **Namenode:** Metadata (files + block locations)
  - **Datanodes:** Save file blocks (usually 64 MB)
- ▶ **Design goal:** Maximum Throughput and data locality for Map-Reduce

*Sends data operations to DataNodes and metadata operations to the NameNode*

*Holds filesystem data and block locations in RAM*



*DataNodes communicate to perform 3-way replication*

*Files are split into blocks and scattered over DataNodes*



# Hadoop

- ▶ For many synonymous to *Big Data Analytics*
- ▶ Large Ecosystem
- ▶ Creator: Doug Cutting (Lucene)
- ▶ Distributors: Cloudera, MapR, HortonWorks
- ▶ Gartner Prognosis: By 2015 65% of all complex analytic applications will be based on Hadoop
- ▶ Users: Facebook, Ebay, Amazon, IBM, Apple, Microsoft, NSA

## Hadoop

Model:

Batch-Analytics  
Framework

License:

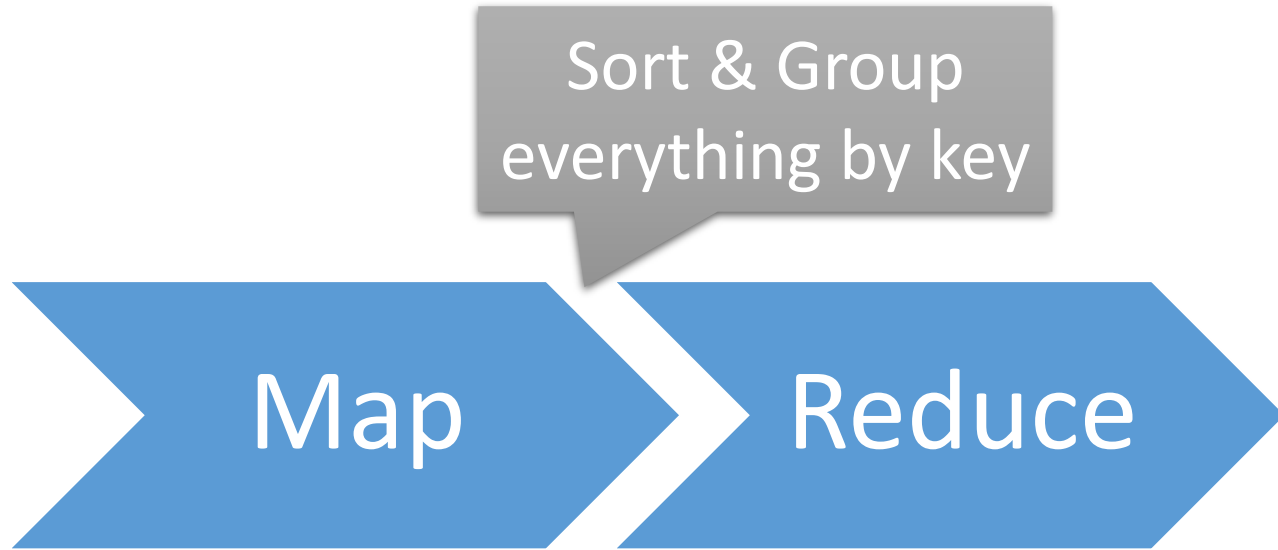
Apache 2

Written in:

Java



# MapReduce



Sort & Group  
everything by key

Map

Reduce

Convert Each Input to a key-  
value pair

Process the list of  
values for each key

# MapReduce: Word Count

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

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

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

1: ich bin ich

(ich, 1)

(bin, 1)

(ich, 1)

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

```
  sum = 0
```

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

```
    sum = sum + v
```

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

(ich, [1,1])

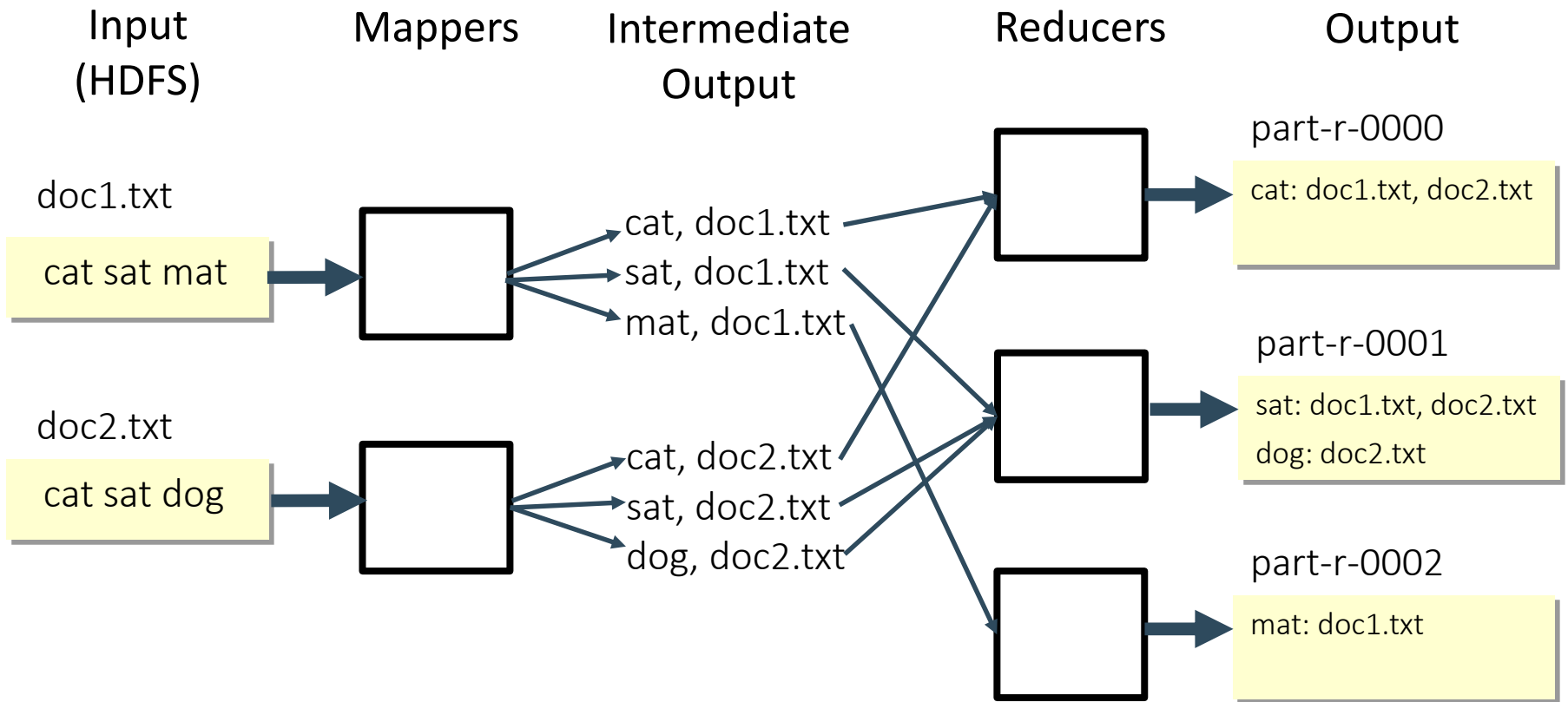
(bin, [1])

(ich, 2)

(bin, 1)

# MapReduce: Example

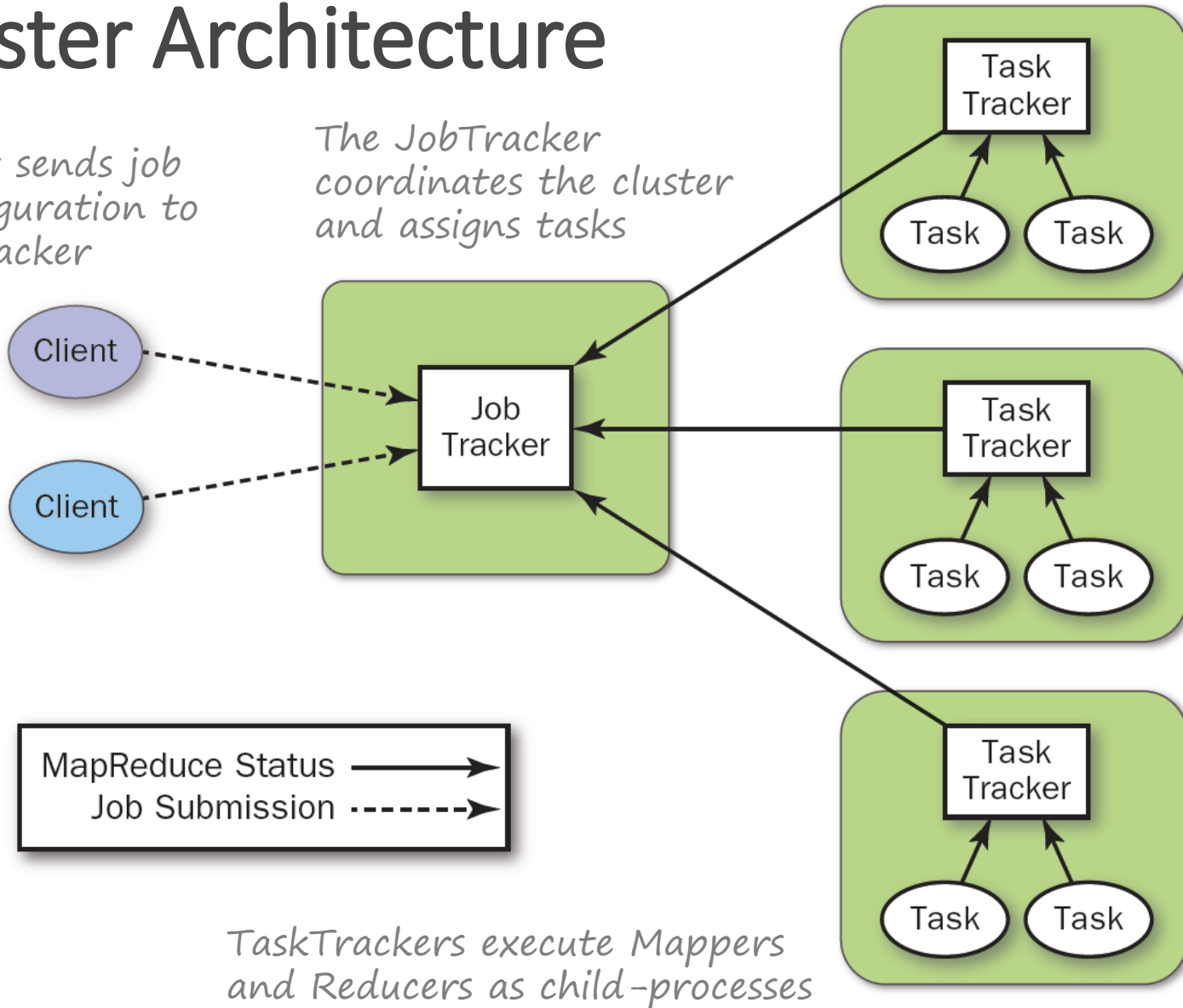
## Constructing a reverse-index



# Cluster Architecture

The client sends job and configuration to the JobTracker

The JobTracker coordinates the cluster and assigns tasks





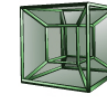
# Summary: Hadoop Ecosystem



- ▶ **Hadoop:** Ecosystem for Big Data Analytics
- ▶ **Hadoop Distributed File System:** scalable, shared-nothing file system for throughput-oriented workloads
- ▶ **Map-Reduce:** Paradigm for performing scalable distributed batch analysis
- ▶ Other Hadoop projects:
  - **Hive:** SQL(-dialect) compiled to YARN jobs (Facebook)
  - **Pig:** workflow-oriented scripting language (Yahoo)
  - **Mahout:** Machine-Learning algorithm library in Map-Reduce
  - **Flume:** Log-Collection and processing framework
  - **Whirr:** Hadoop provisioning for cloud environments
  - **Giraph:** Graph processing à la Google Pregel
  - **Drill, Presto, Impala:** SQL Engines

# NoSQL landscape

## Document



HYPERTABLE



Google Datastore



Cassandra

## Wide Column

## Key-Value



redis



CouchBase

## Graph



Project Voldemort



AEROSPIKE

# Popularity

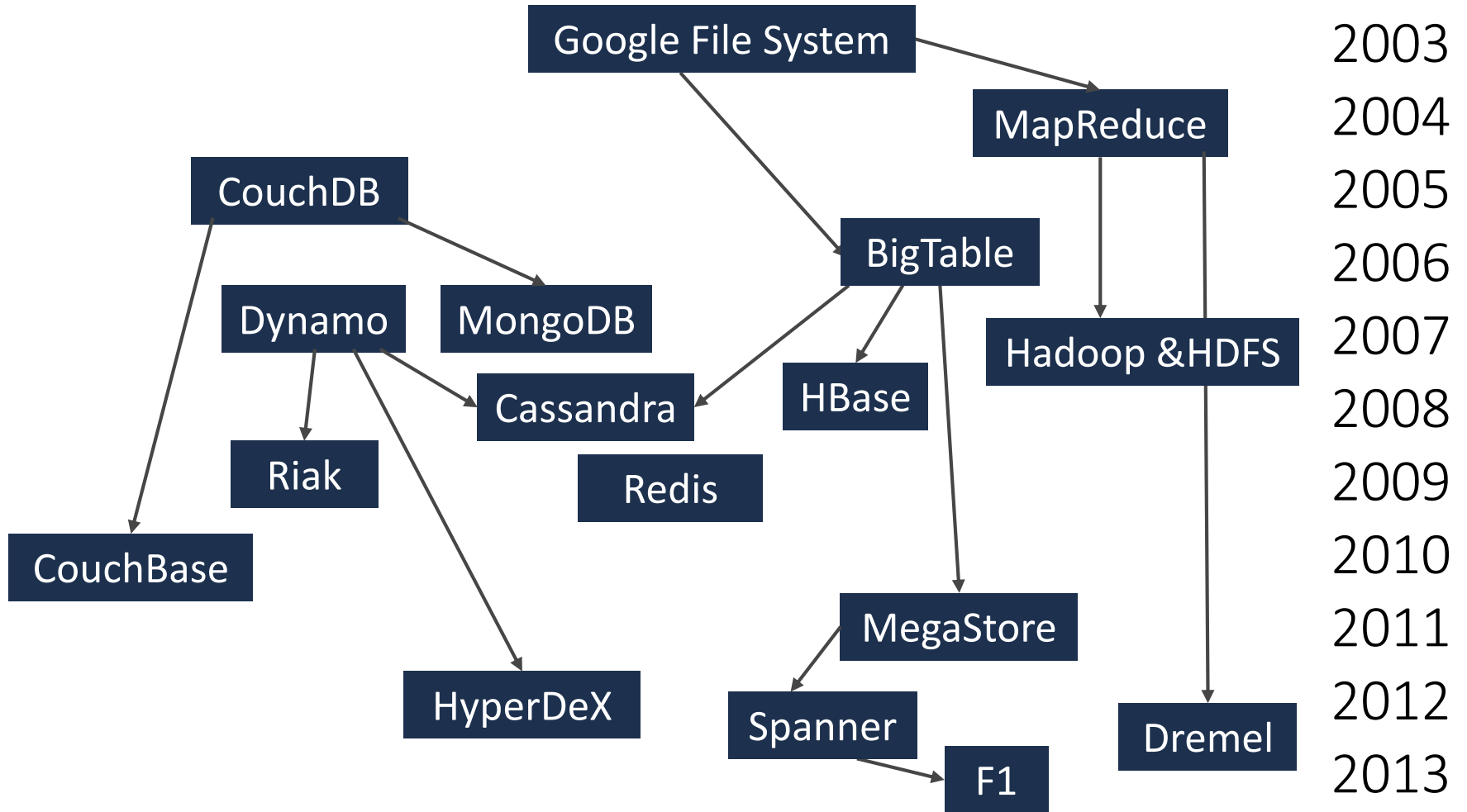
<http://db-engines.com/de/ranking>

Rang	DBMS	Modell	Punkte
1.	Oracle	Relational DBMS	1514,90
2.	MySQL	Relational DBMS	1334,94
3.	Microsoft SQL Server	Relational DBMS	1286,22
4.	PostgreSQL	Relational DBMS	199,39
5.	DB2	Relational DBMS	177,04
6.	Microsoft Access	Relational DBMS	149,66
<b>7.</b>	<b>MongoDB</b>	<b>Document Store</b>	<b>137,49</b>
8.	Sybase	Relational DBMS	88,41
9.	SQLite	Relational DBMS	87,81
10.	Teradata	Relational DBMS	51,11
11.	Solr	Suchmaschine	46,43
<b>12.</b>	<b>Cassandra</b>	<b>Wide Column Store</b>	<b>37,64</b>
<b>13.</b>	<b>Redis</b>	<b>Key-Value Store</b>	<b>34,22</b>

Rang	DBMS	Modell	Punkte
<b>14.</b>	<b>Memcached</b>	<b>Key-Value Store</b>	<b>30,73</b>
<b>15.</b>	<b>HBase</b>	<b>Wide Column Store</b>	<b>25,78</b>
16.	Informix	Relational DBMS	24,73
17.	Hive	Relational DBMS	22,16
<b>18.</b>	<b>CouchDB</b>	<b>Document Store</b>	<b>15,93</b>
19.	Firebird	Relational DBMS	14,55
20.	Netezza	Relational DBMS	11,44
21.	dBASE	Relational DBMS	10,44
22.	Elasticsearch	Suchmaschine	9,51
23.	Sphinx	Suchmaschine	9,02
<b>24.</b>	<b>Riak</b>	<b>Key-Value Store</b>	<b>8,99</b>
<b>25.</b>	<b>Neo4j</b>	<b>Graph DBMS</b>	<b>8,83</b>

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

# 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, Accumolo



- ▶ **Dynamo** (2007, Amazon)

- Available, Partition tolerant
- **Key-Value** interface
- Eventually Consistent, always writable, fault-tolerant
- Riak, Cassandra, Voldemort, DynamoDB



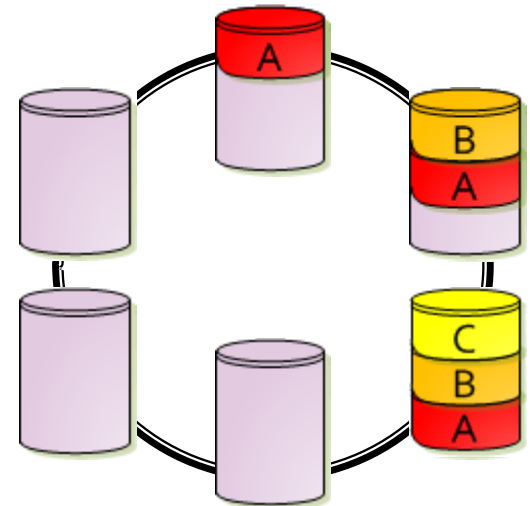
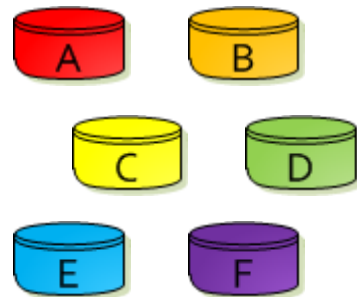
Chang, Fay, et al. "Bigtable: A distributed storage system for structured data."



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

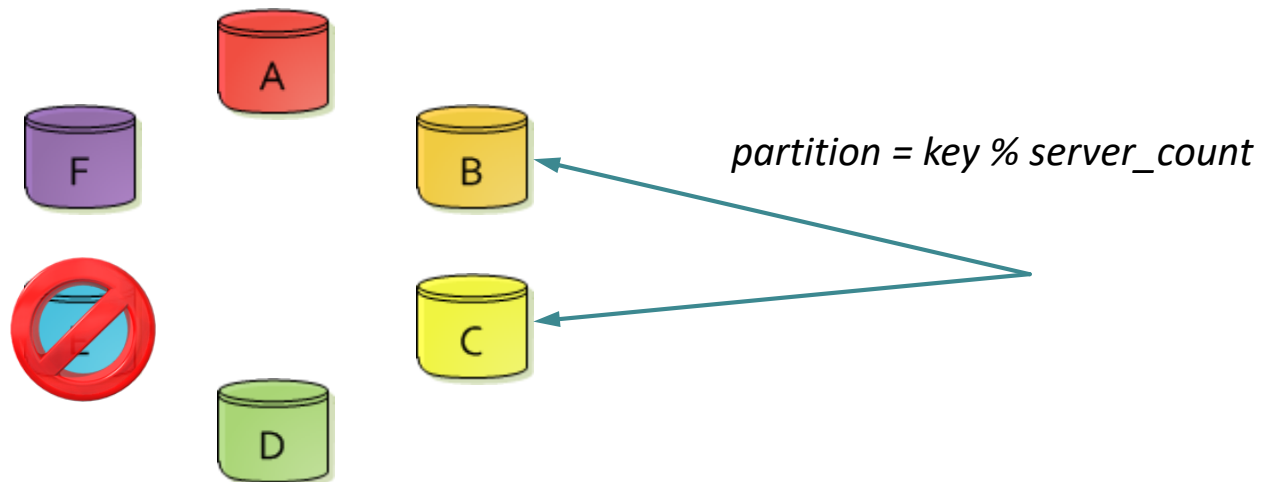
# Dynamo (AP)

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



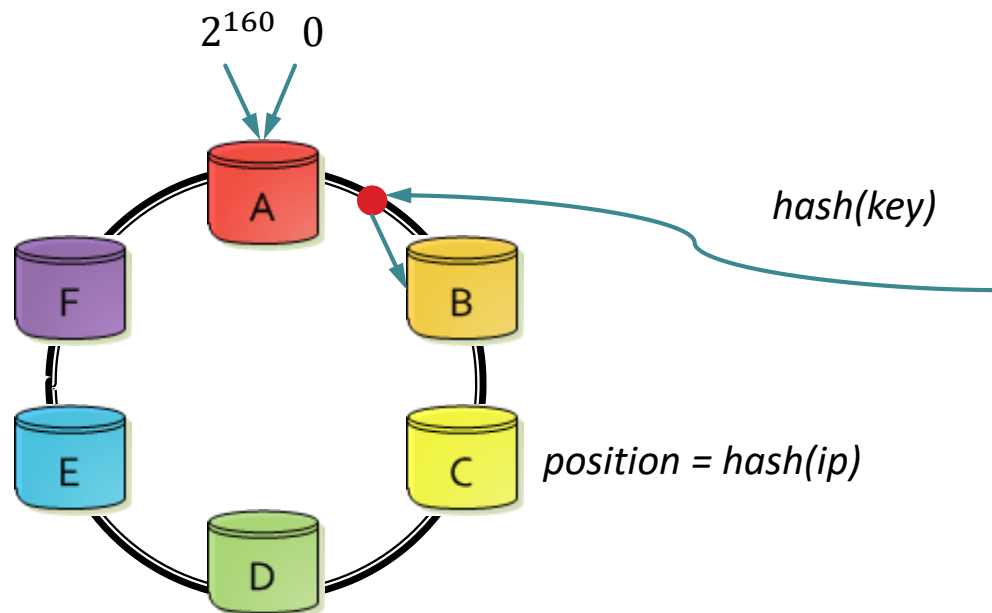
# Consistent Hashing

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



# Consistent Hashing

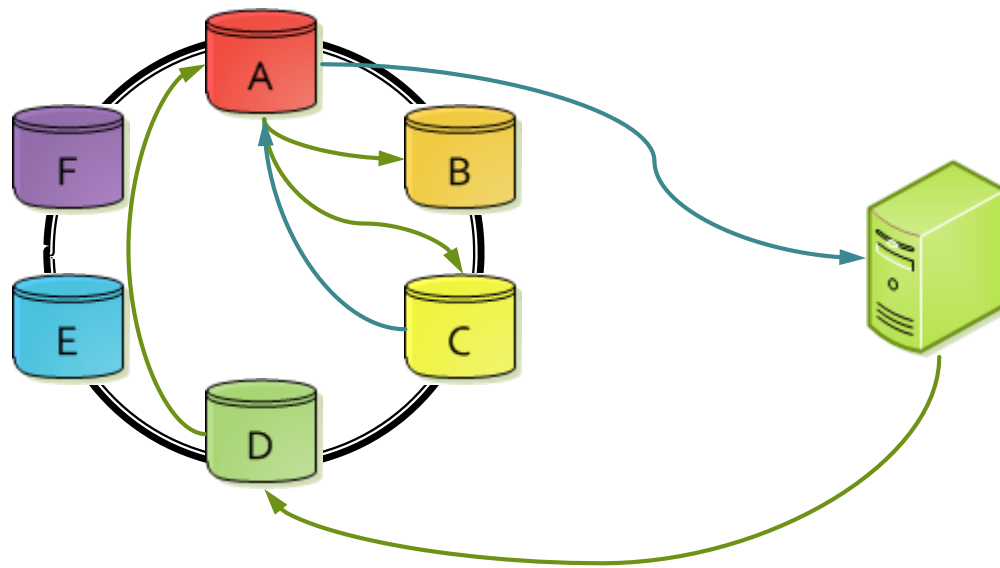
- ▶ Solution: **Consistent Hashing** – mapping of data to nodes is stable under topology changes





# Reading and Writing

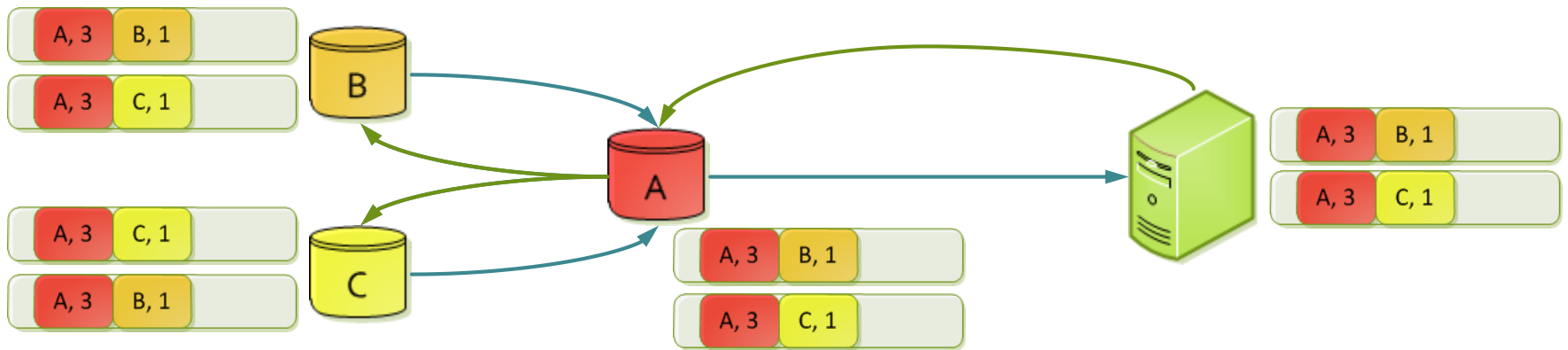
- ▶ An arbitrary node acts as a coordinator
- ▶ **N**: number of replicas
- ▶ **R**: number of nodes that need to confirm a read
- ▶ **W**: number of nodes that need to confirm a write



N=3  
R=2  
W=1

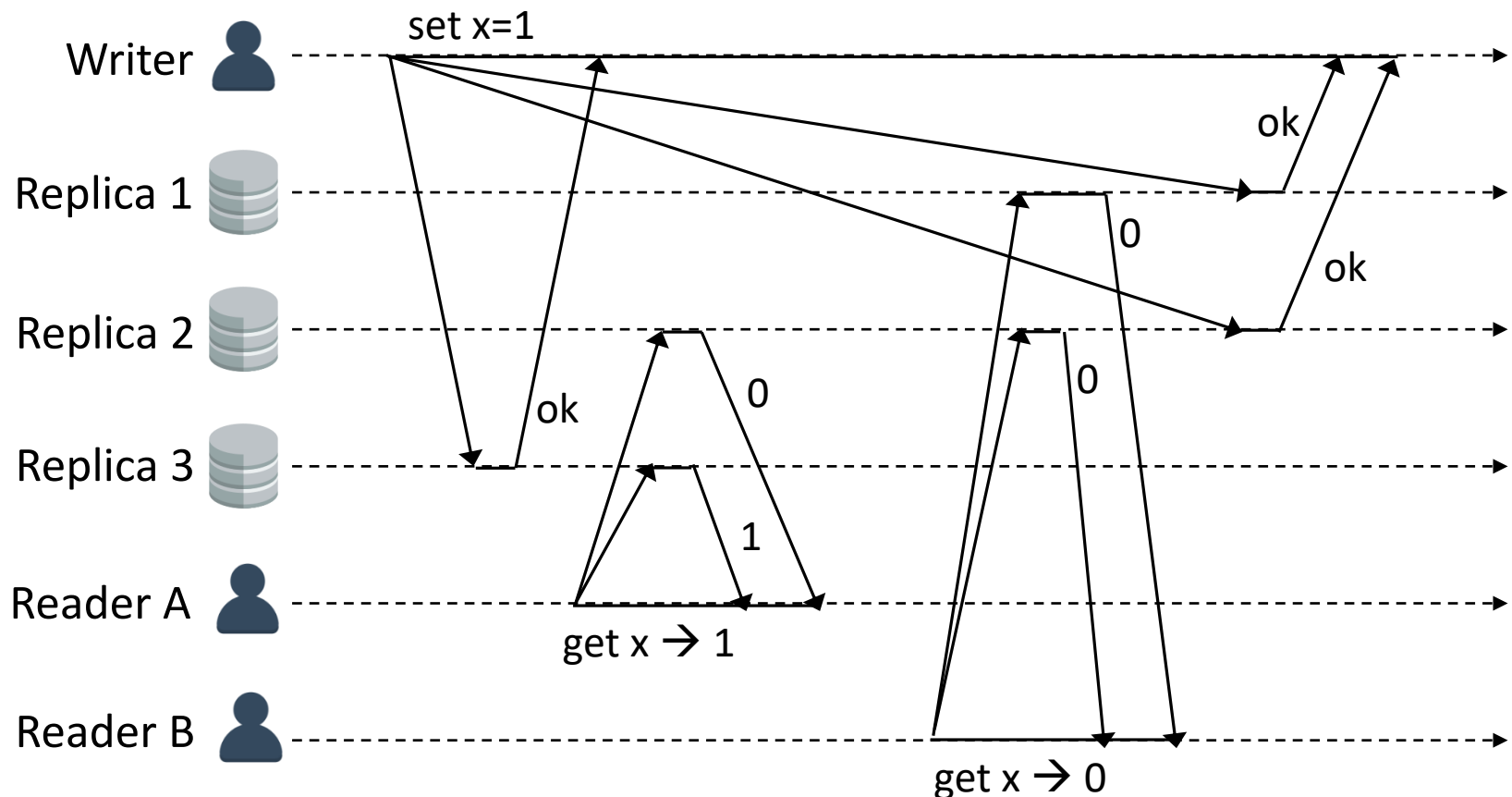
# Versioning and Consistency

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



# $R + W > N$ does not imply linearizability

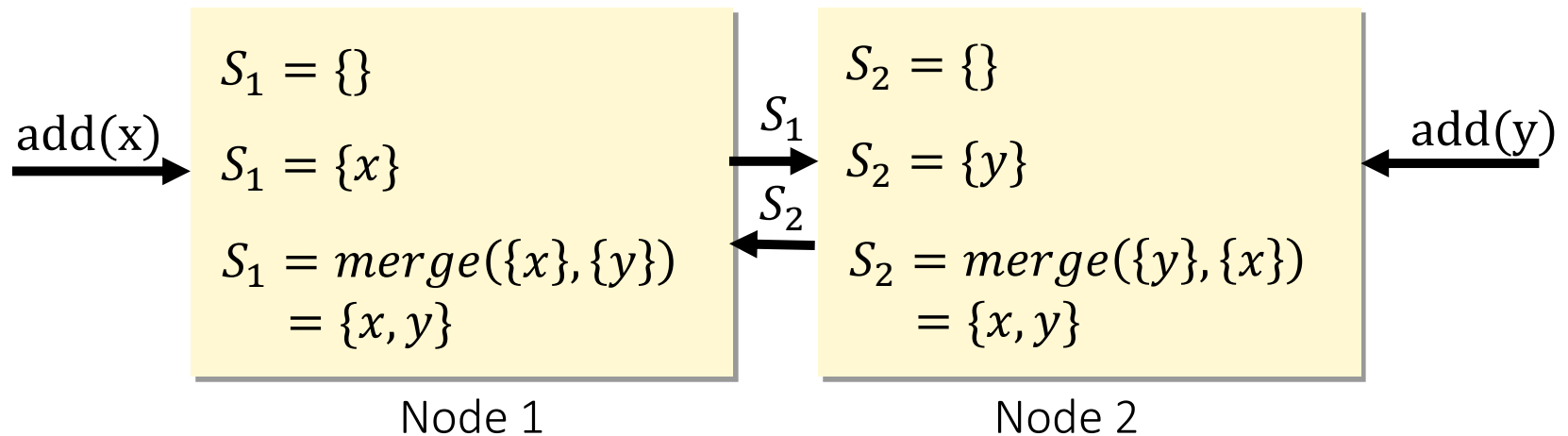
- ▶ Consider the following execution:



# CRDTs

## Convergent/Commutative Replicated Data Types

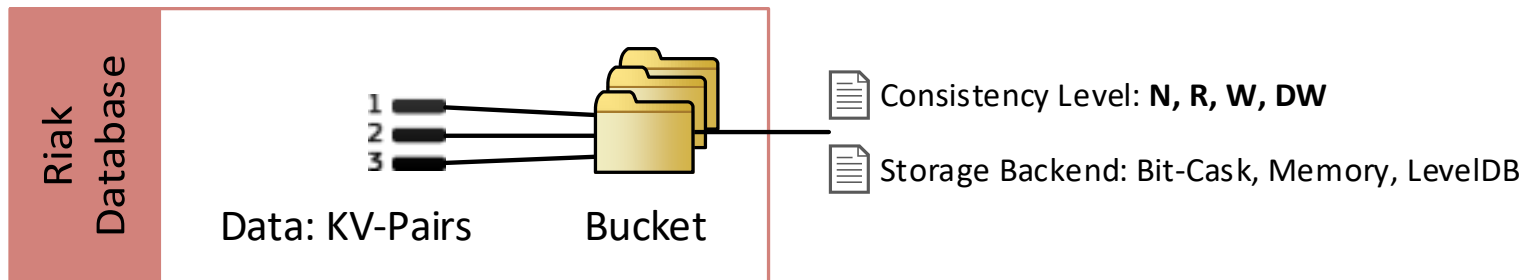
- ▶ **Goal:** avoid manual conflict-resolution
- ▶ **Approach:**
  - **State-based** – commutative, idempotent merge function
  - **Operation-based** – broadcasts of commutative updates
- ▶ Example: State-based Grow-only-Set (G-Set)



# 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
  - **REST-API**

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



# Summary: Dynamo and Riak



- ▶ **Consistent Hashing:** hash-based distribution with stability under topology changes (e.g. machine failures)
- ▶ Parameters: **N** (Replicas), **R** (Read Acks), **W** (Write Acks)
  - $N=3, R=W=1$  → fast, potentially inconsistent
  - $N=3, R=3, W=1$  → slower reads, most recent object version contained
- ▶ Available and Partition-Tolerant
- ▶ **Vector Clocks:** concurrent modification can be detected, inconsistencies are healed through the application
- ▶ **API:** Create, Read, Update, Delete (CRUD) on key-value pairs
- ▶ **Riak:** Open-Source Implementation of the Dynamo paper



# 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 on one machine
- ▶ Redis Cluster adds sharding

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

# 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

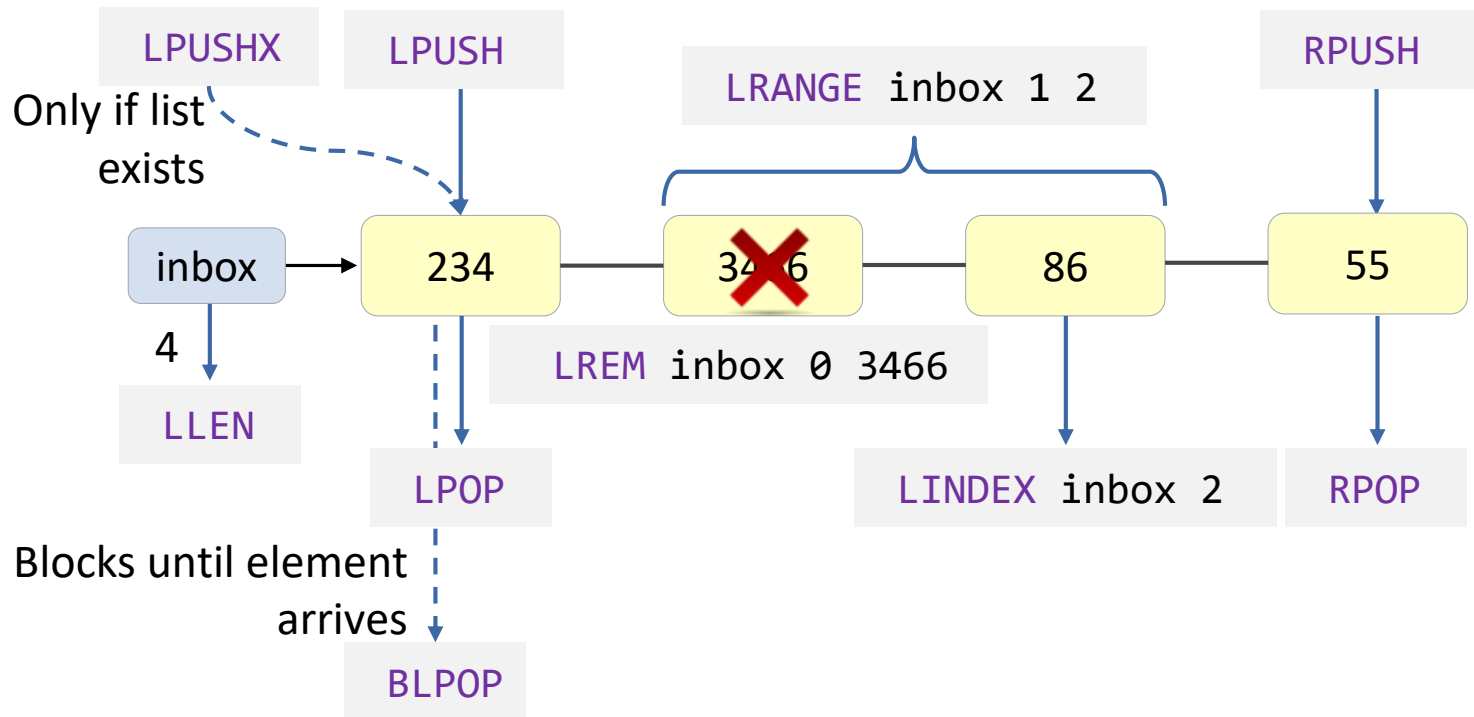
users:2:notif

"{event: 'comment posted', time : ..."



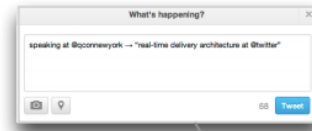
# Example Redis Data Structure: lists

## ▶ (Linked) Lists:



# Example Redis Use-Case: Twitter

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

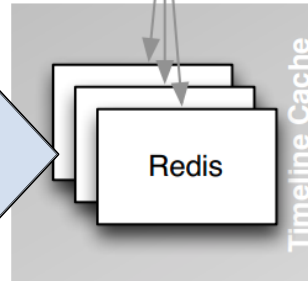
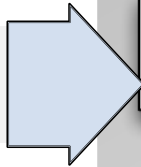


Write API

Fanout

>150 million users  
~300k timeline queries/s

`RPUSHX user_id tweet`







Tweet ID	User ID	Bits	
Tweet ID	User ID	Bits	Tweet ID
Tweet ID	User ID	Bits	
Tweet ID	User ID	Bits	
Tweet ID	User ID	Bits	Tweet ID
Tweet ID	User ID	Bits	



# Classification: Redis

## Techniques

 <b>Sharding</b>	Range-Sharding	Hash-Sharding	Entity-Group Sharding	Consistent Hashing	Shared Disk
 <b>Replication</b>	Transaction Protocol	Sync. Replication	<b>Async. Replication</b>	<b>Primary Copy</b>	Update Anywhere
 <b>Storage Management</b>	<b>Logging</b>	Update-in-Place	<b>Caching</b>	<b>In-Memory</b>	Append-Only Storage
 <b>Query Processing</b>	Global Index	Local Index	Query Planning	Analytics	Materialized Views

# Google BigTable (CP)

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

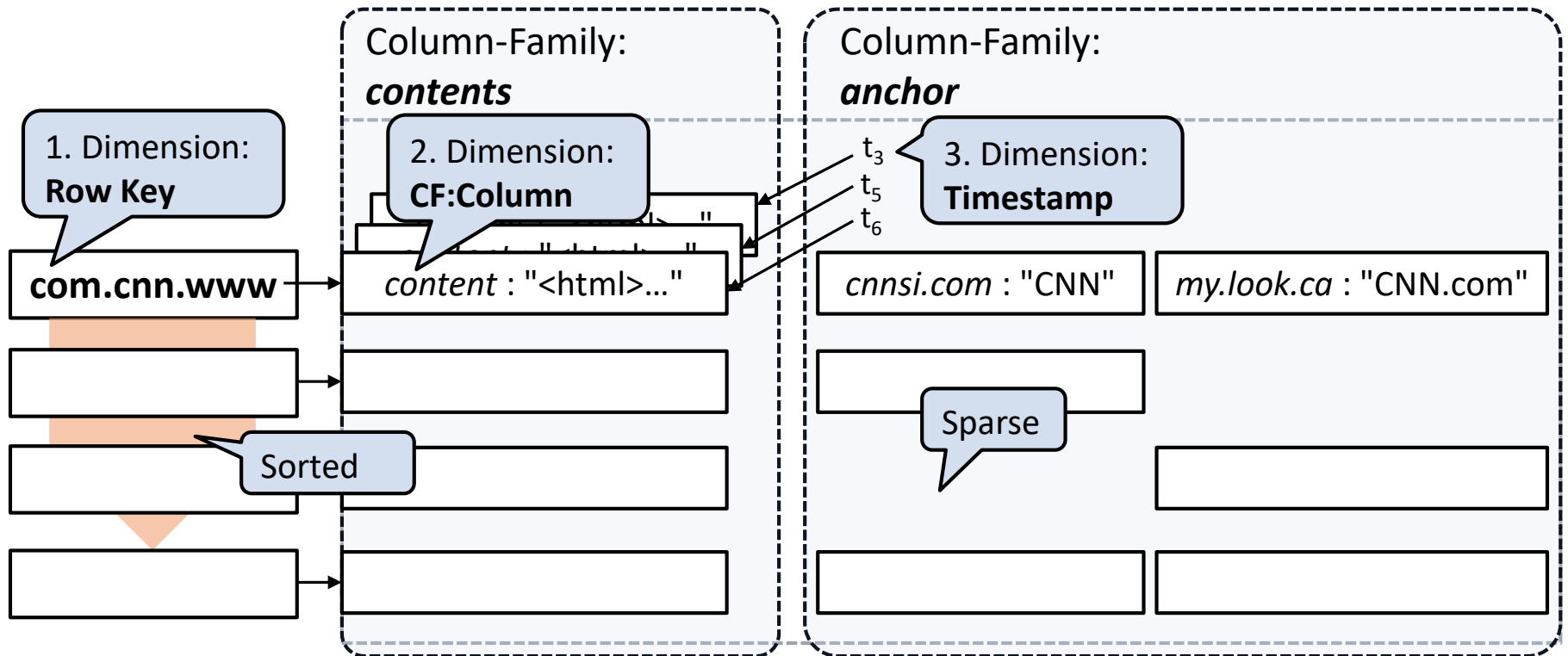
A Bigtable is a sparse, distributed, persistent multidimensional sorted map.

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

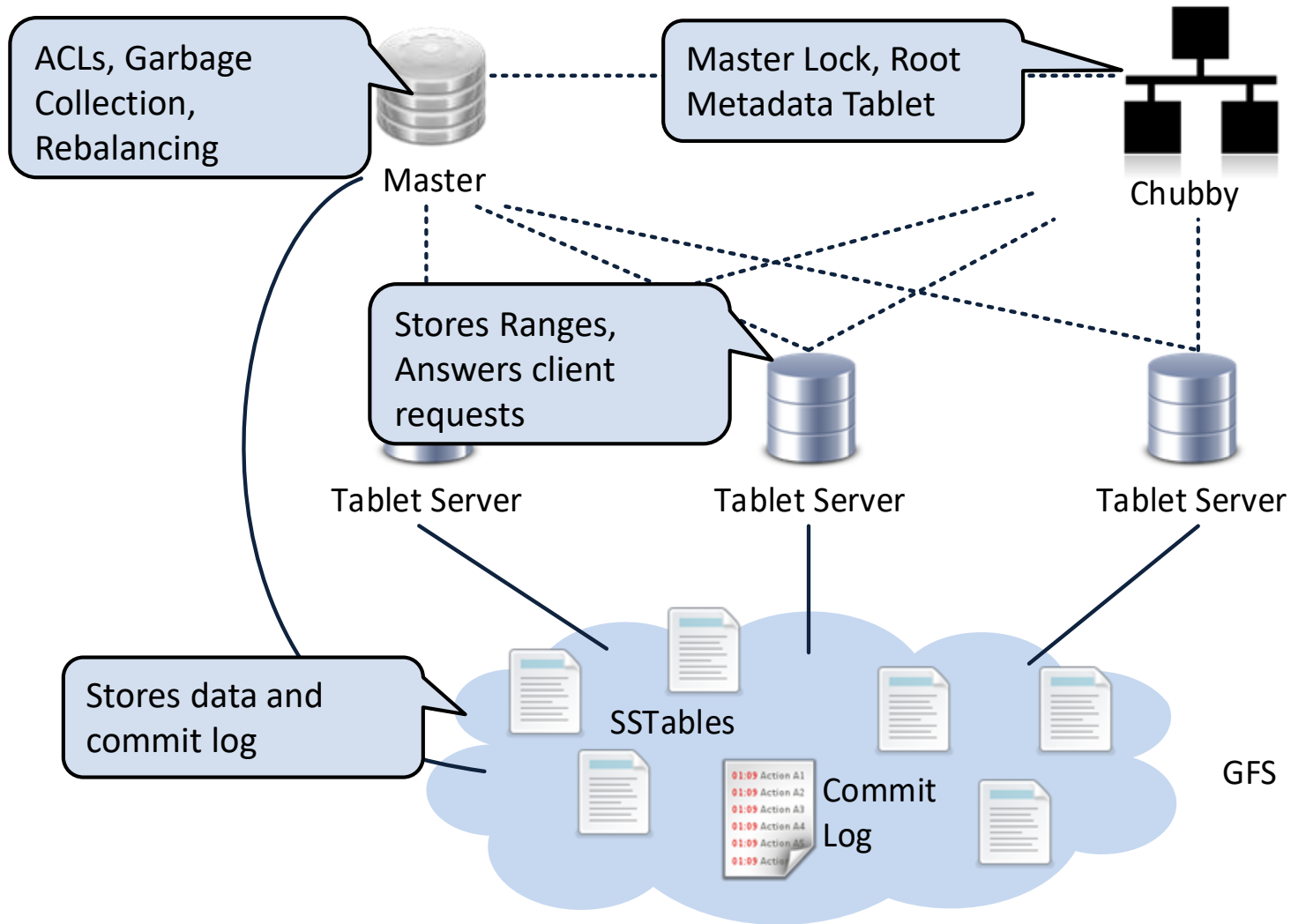


# Wide-Column Data Modelling

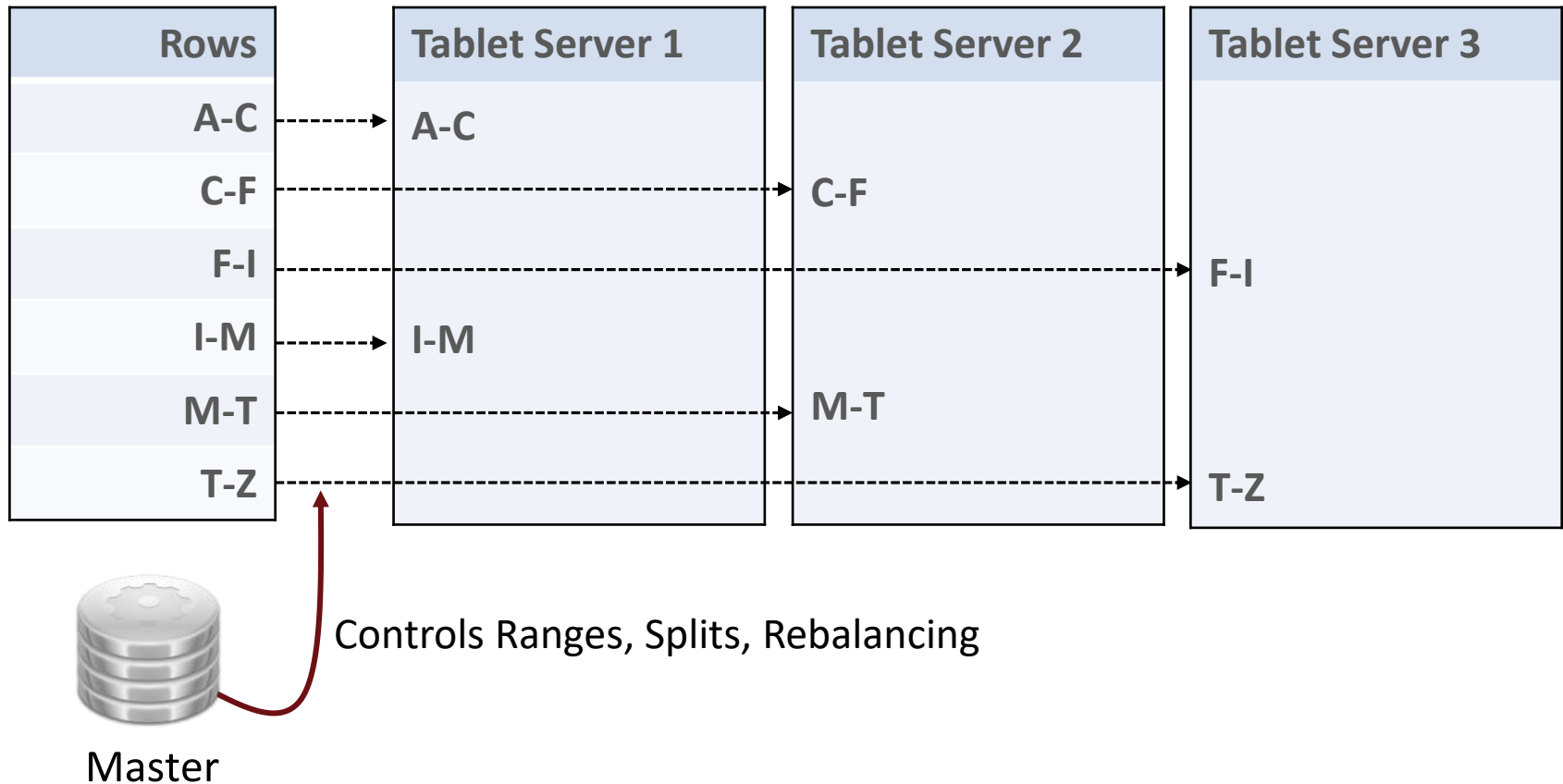
- ▶ Storage of crawled web-sites („Webtable“):



# Architecture

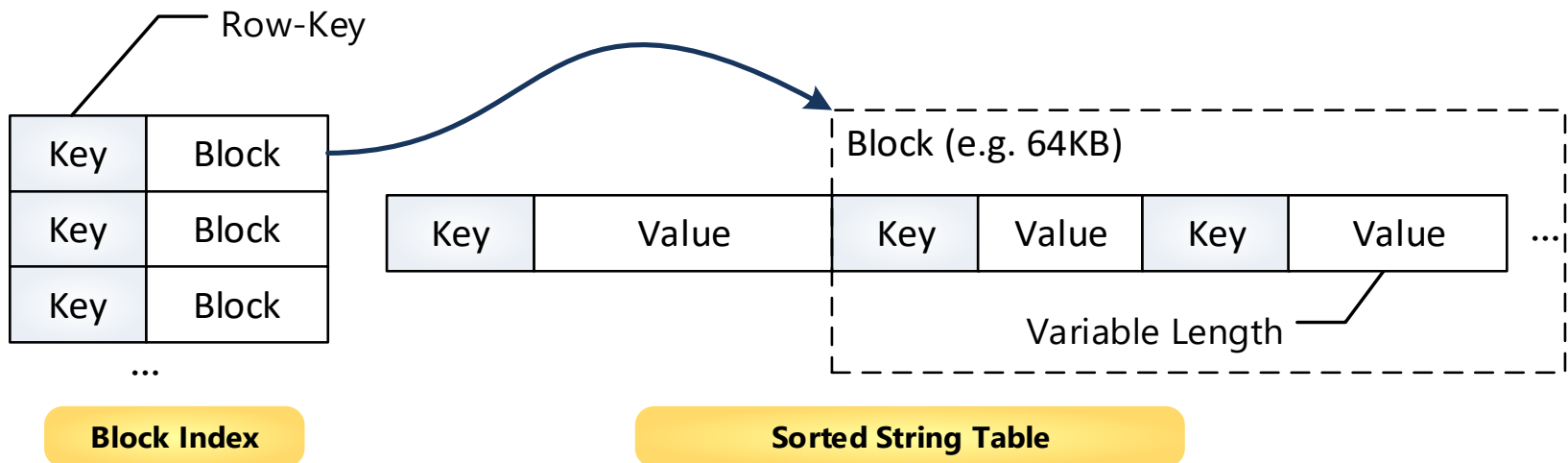


# Range-based Sharding



# 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





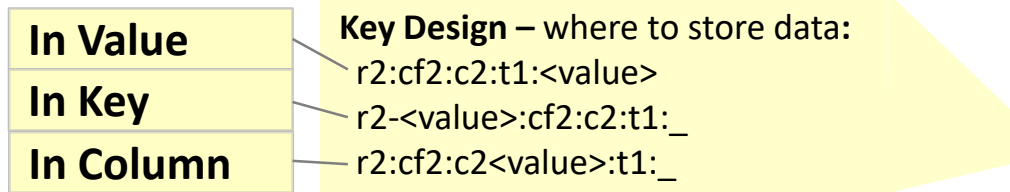
# Apache HBase (CP)

- ▶ 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

# HBase Storage

## ▶ Logical to physical mapping:



Key	cf1:c1	cf1:c2	cf2:c1	cf2:c2
r1	■		■	
r2		■		■
r3				■
r4		■		■
r5	■		■	

```

r1:cf2:c1:t1:<value>
r2:cf2:c2:t1:<value>
r3:cf2:c2:t2:<value>
r3:cf2:c2:t1:<value>
r5:cf2:c1:t1:<value>
  
```

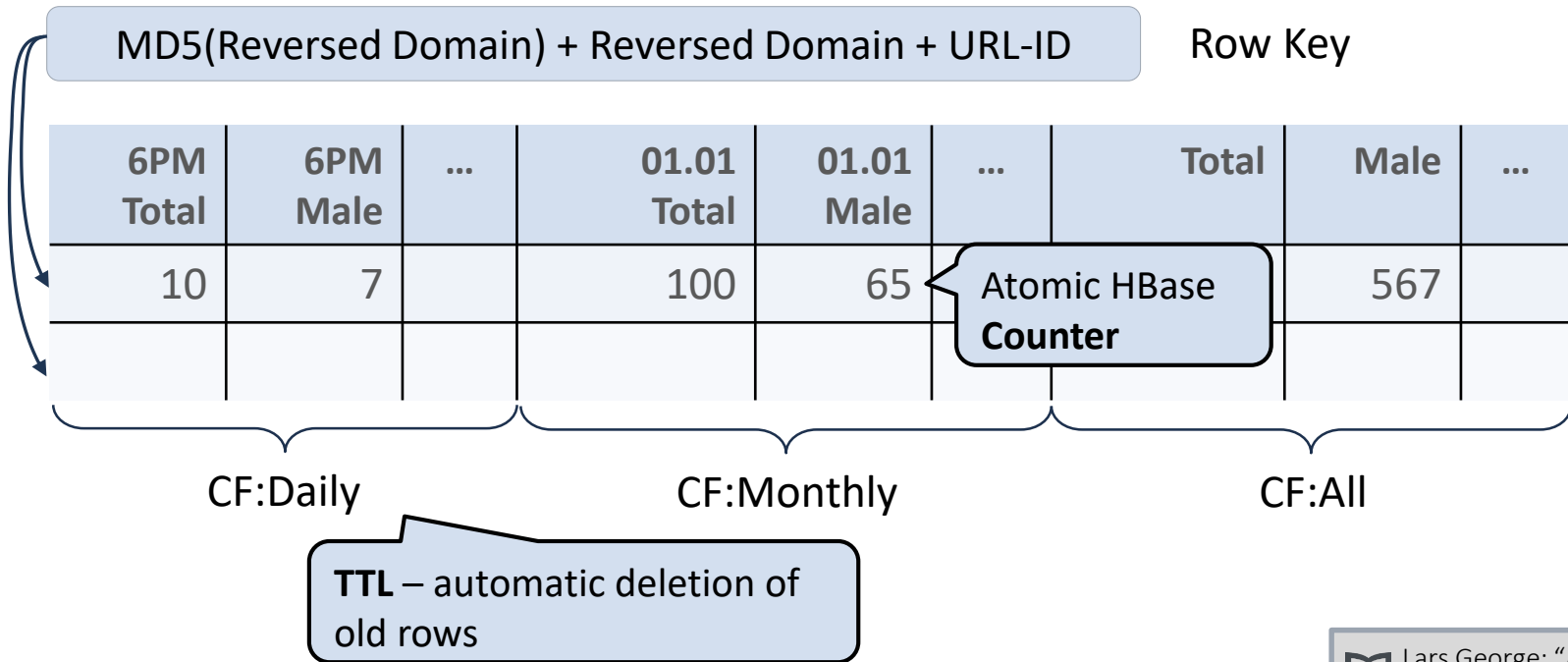
**HFile cf2**

```

r1:cf1:c1:t1:<value>
r2:cf1:c2:t1:<value>
r3:cf1:c2:t1:<value>
r3:cf1:c1:t2:<value>
r5:cf1:c1:t1:<value>
  
```

**HFile cf1**

# Example: Facebook Insights







# Summary: BigTable, HBase



- ▶ Data model:  $(rowkey, cf: column, timestamp) \rightarrow value$
- ▶ **API**: CRUD + Scan(*start-key*, *end-key*)
- ▶ Uses distributed file system (GFS/HDFS)
- ▶ Storage structure: **Memtable** (in-memory data structure) + **SSTable** (persistent; append-only-IO)
- ▶ **Schema design**: only primary key access  $\rightarrow$  implicit schema (key design) needs to be carefully planned
- ▶ **HBase**: very literal open-source implementation BigTable
- ▶ **Cassandra**: combination of Dynamo and BigTable ideas

# Classification: HBase

## Techniques

 <b>Sharding</b>	Range-Sharding	Hash-Sharding	Entity-Group Sharding	Consistent Hashing	Shared Disk
 <b>Replication</b>	Transaction Protocol	Sync. Replication	Async. Replication	Primary Copy	Update Anywhere
 <b>Storage Management</b>	Logging	Update-in-Place	Caching	In-Memory	Append-Only Storage
 <b>Query Processing</b>	Global Index	Local Index	Query Planning	Analytics	Materialized Views







# 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
- ▶ **Secondary indices:** hidden table with mapping → queries with simple equality condition

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

# Classification: Cassandra

## Techniques

 <b>Sharding</b>	Range-Sharding	Hash-Sharding	Entity-Group Sharding	Consistent Hashing	Shared Disk
 <b>Replication</b>	Transaction Protocol	Sync. Replication	Async. Replication	Primary Copy	Update Anywhere
 <b>Storage Management</b>	Logging	Update-in-Place	Caching	In-Memory	Append-Only Storage
 <b>Query Processing</b>	Global Index	Local Index	Query Planning	Analytics	Materialized Views

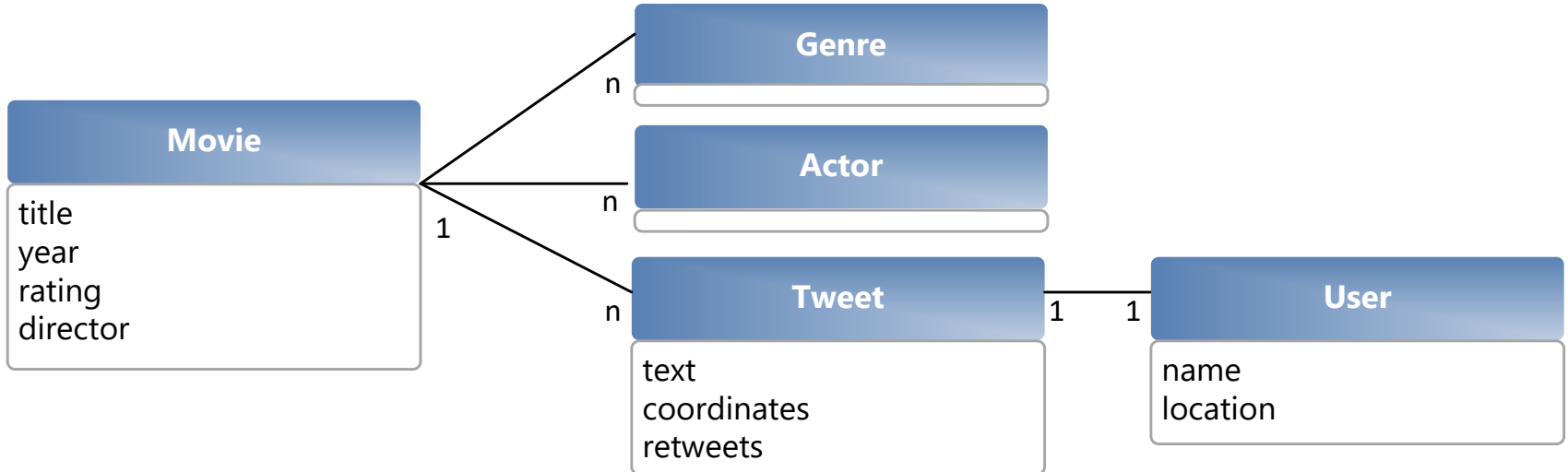
# MongoDB (CP)

- ▶ From humongous  $\cong$  gigantic
- ▶ Tunable consistency
- ▶ Schema-free document database
- ▶ Allows complex queries and indexing
- ▶ **Sharding** (either range- or hash-based)
- ▶ **Replication** (either synchronous or asynchronous)
- ▶ Storage Management:
  - **Write-ahead logging** for redos (*journaling*)
  - **Memory-mapped** storage files, buffer management handled by operating system (paging)

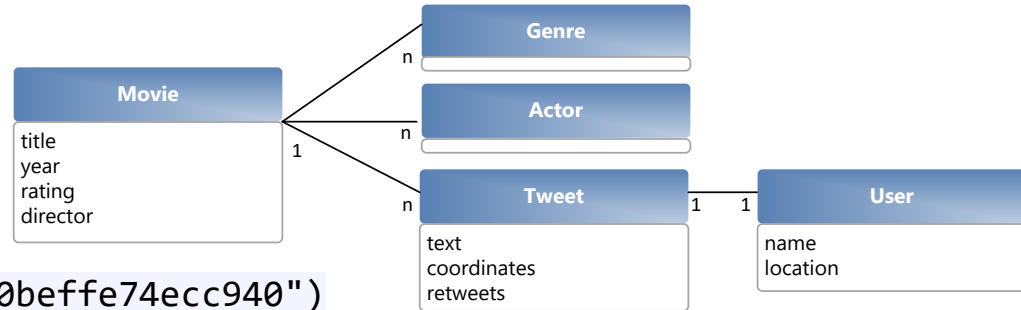
MongoDB
Model:
Document
License:
GNU AGPL 3.0
Written in:
C++



# Data Modelling



# Data Modelling



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

**Movie Document**

**Denormalisation** instead of joins

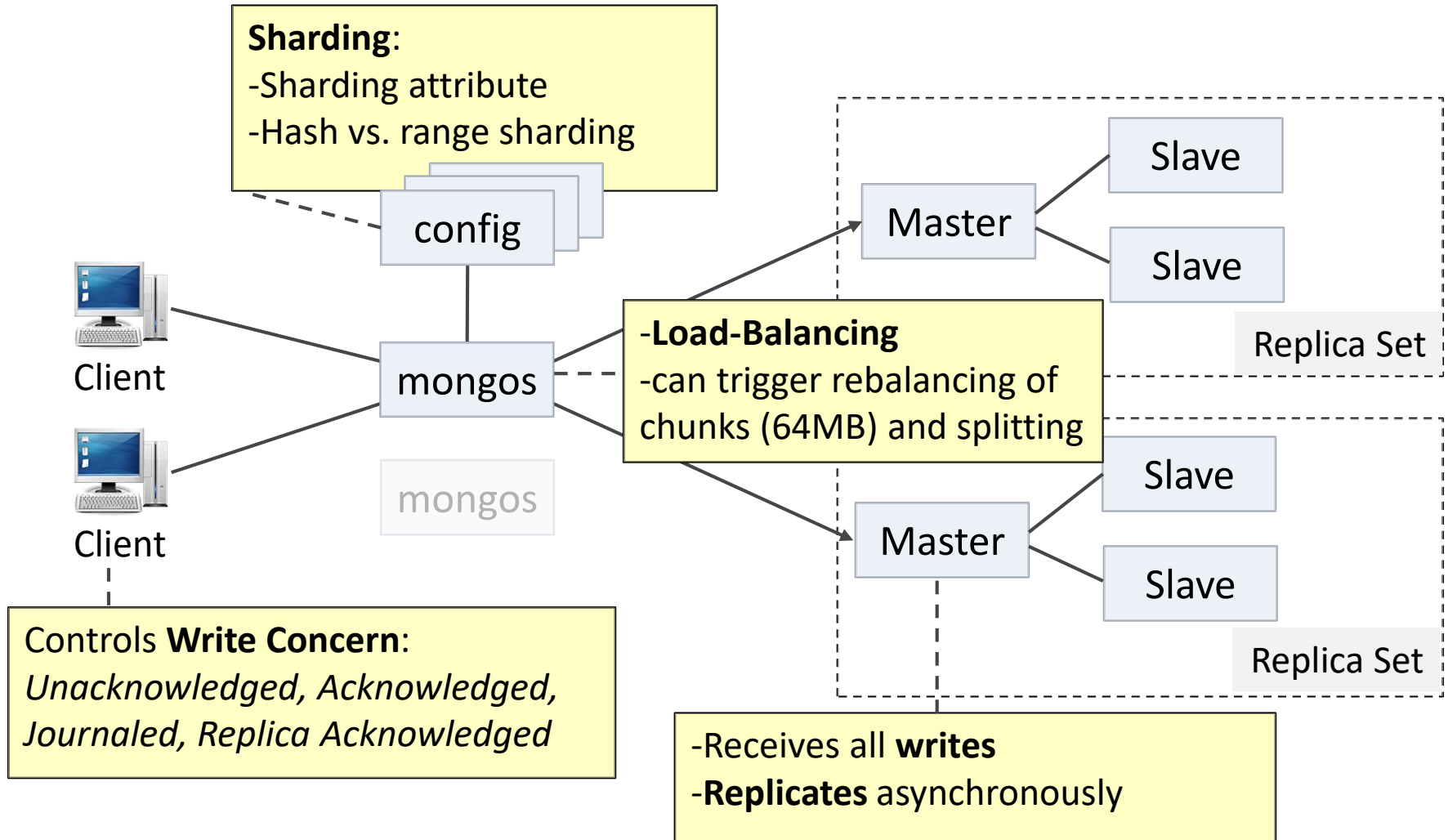
**Nesting** replaces 1:n and 1:1 relations

**Schemafreeness:**  
Attributes per document

**Unit of atomicity:**  
document





**Principles**

# Sharding und Replication



# Classification: MongoDB

## Techniques

 <b>Sharding</b>	Range-Sharding	Hash-Sharding	Entity-Group Sharding	Consistent Hashing	Shared Disk
 <b>Replication</b>	Transaction Protocol	Sync. Replication	Async. Replication	Primary Copy	Update Anywhere
 <b>Storage Management</b>	Logging	Update-in-Place	Caching	In-Memory	Append-Only Storage
 <b>Query Processing</b>	Global Index	Local Index	Query Planning	Analytics	Materialized Views

# Outline



Foundations: Big Data,  
Scalability, Availability



The 4 Classes of NoSQL  
Databases



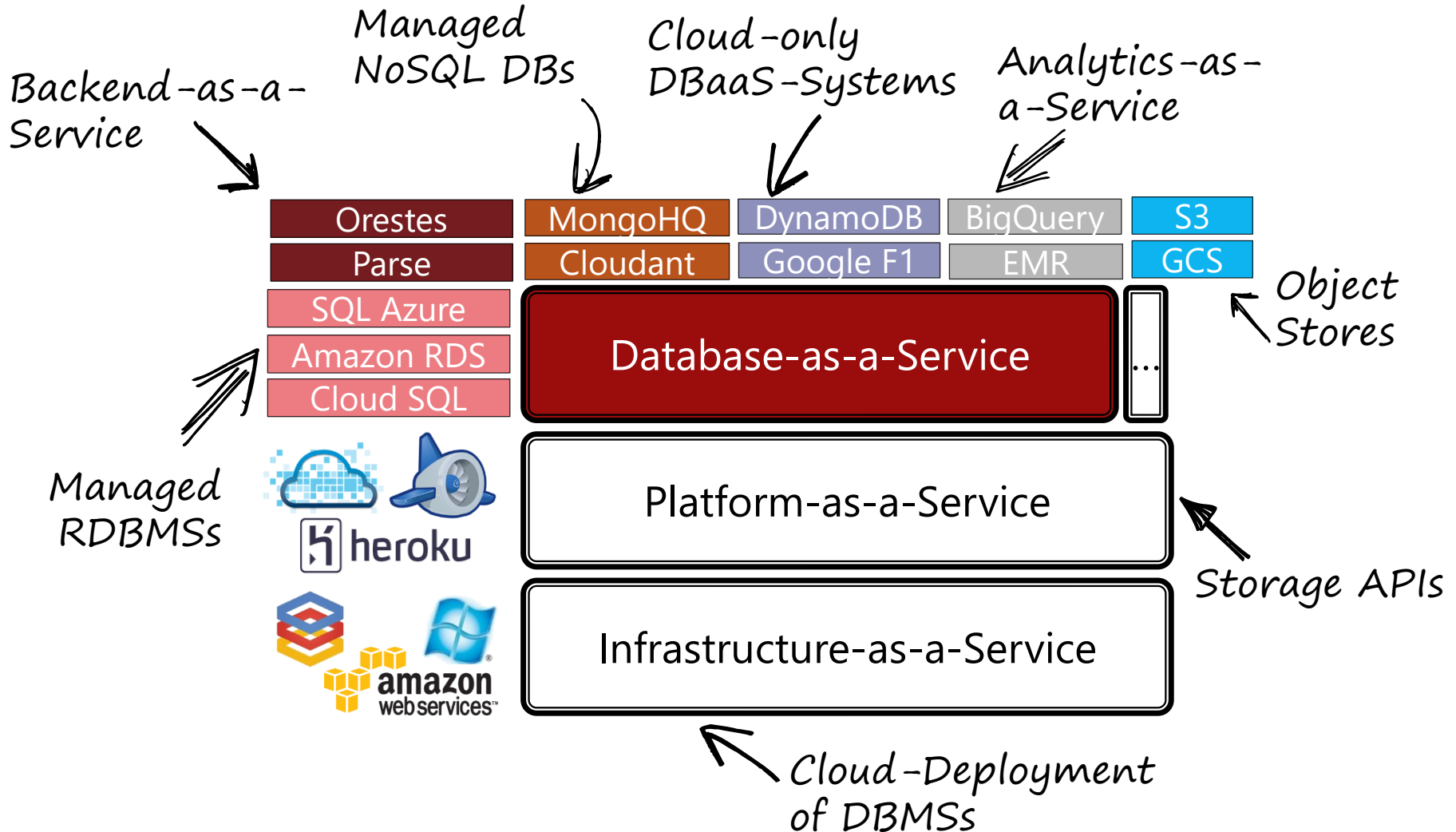
NoSQL Examples: concrete  
Architectures, Systems, APIs



Cloud Databases

- Database-as-a-Service
- Backend-as-a-Service

# Cloud Databases



# Database-as-a-Service

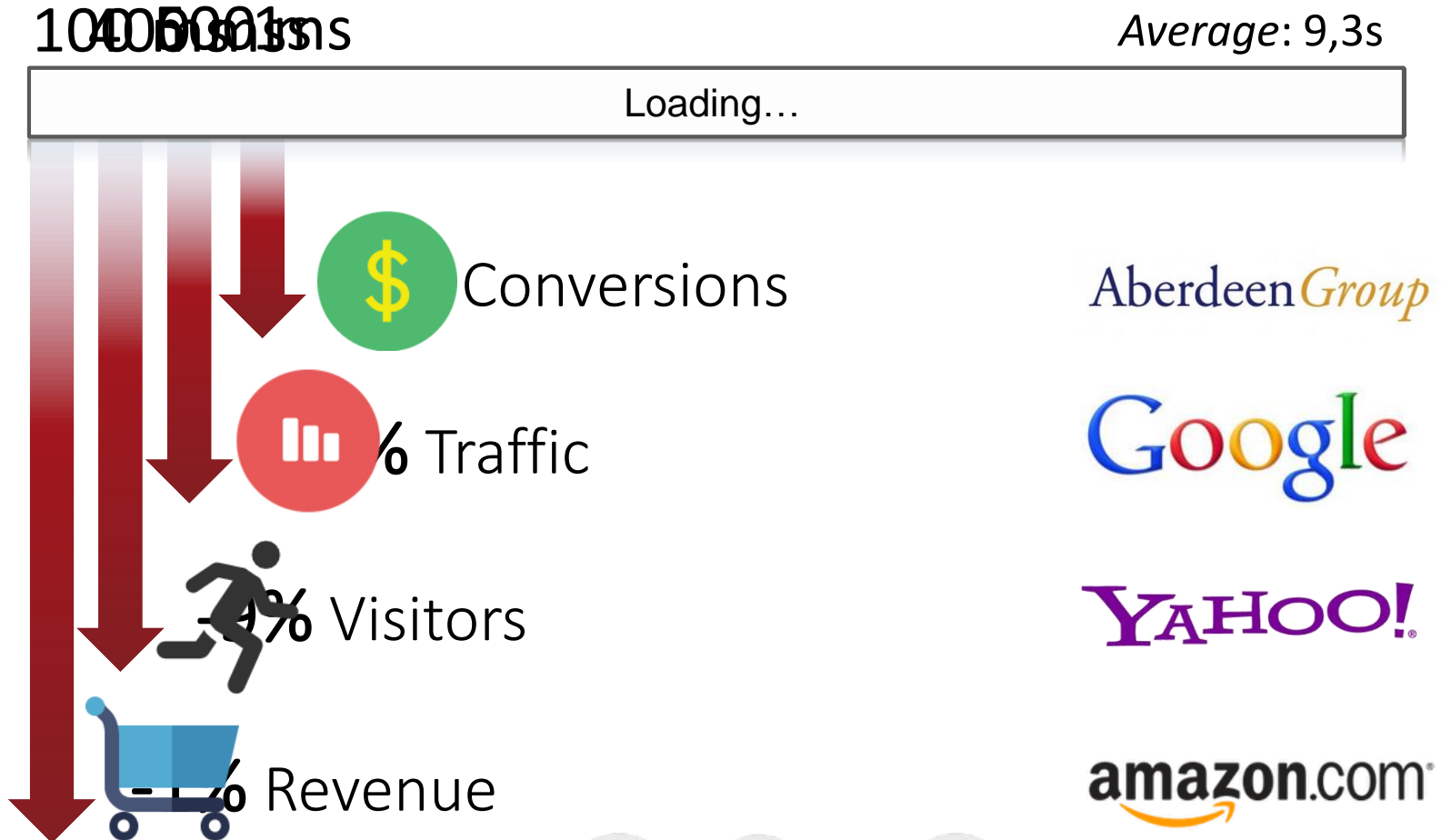
- ▶ Cloud databases with a **pay-per-use** pricing model
- ▶ **Managed Database Service:** Existing DBMS deployed and managed in the cloud
  - Managed NoSQL System (e.g. MongoHQ, Redis2Go)
  - Managed RDBMS (e.g. Amazon Relational Database Service)
- ▶ **Proprietary Database Service:** special DBMS built for cloud environments (e.g. Amazon DynamoDB)
- ▶ **Object Stores:** cloud-based file storage (e.g. Amazon S3)
- ▶ **Backend-as-a-Service:** Database + Implementation of standard app concerns (e.g. user management, push)

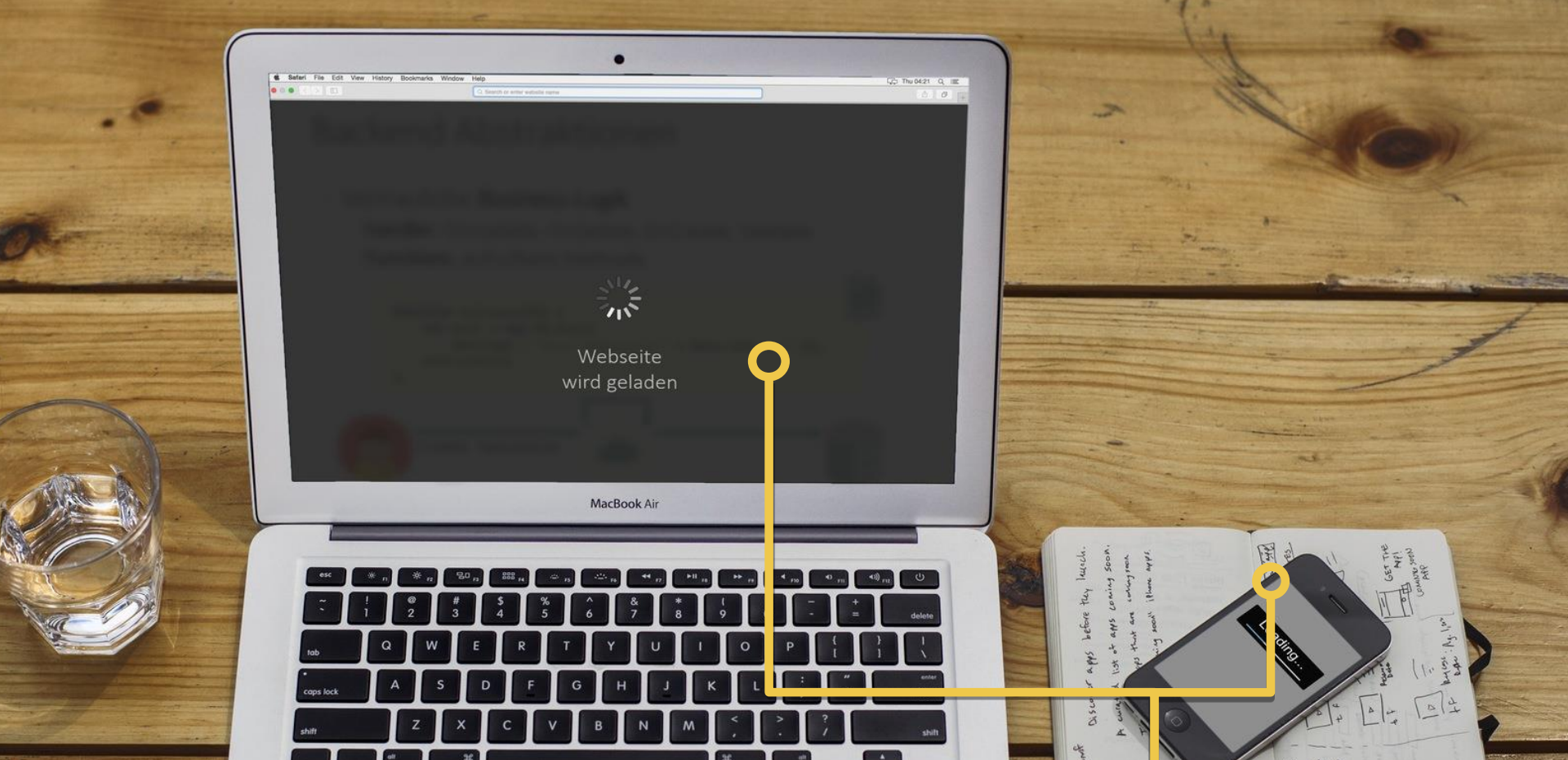


Presentation  
is loading



# The Latency Problem



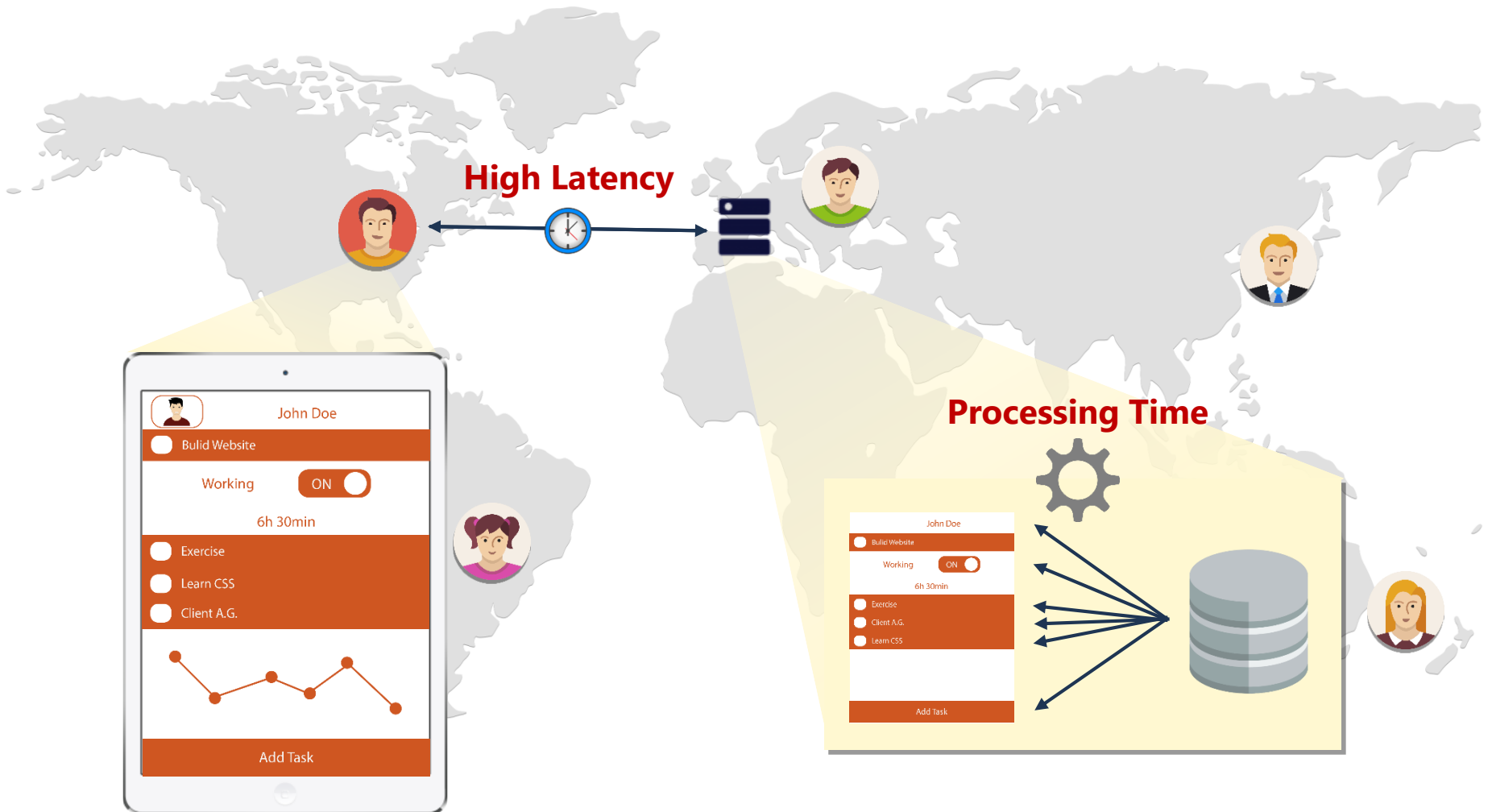


If perceived speed is such an import factor

...what causes slow page load times?

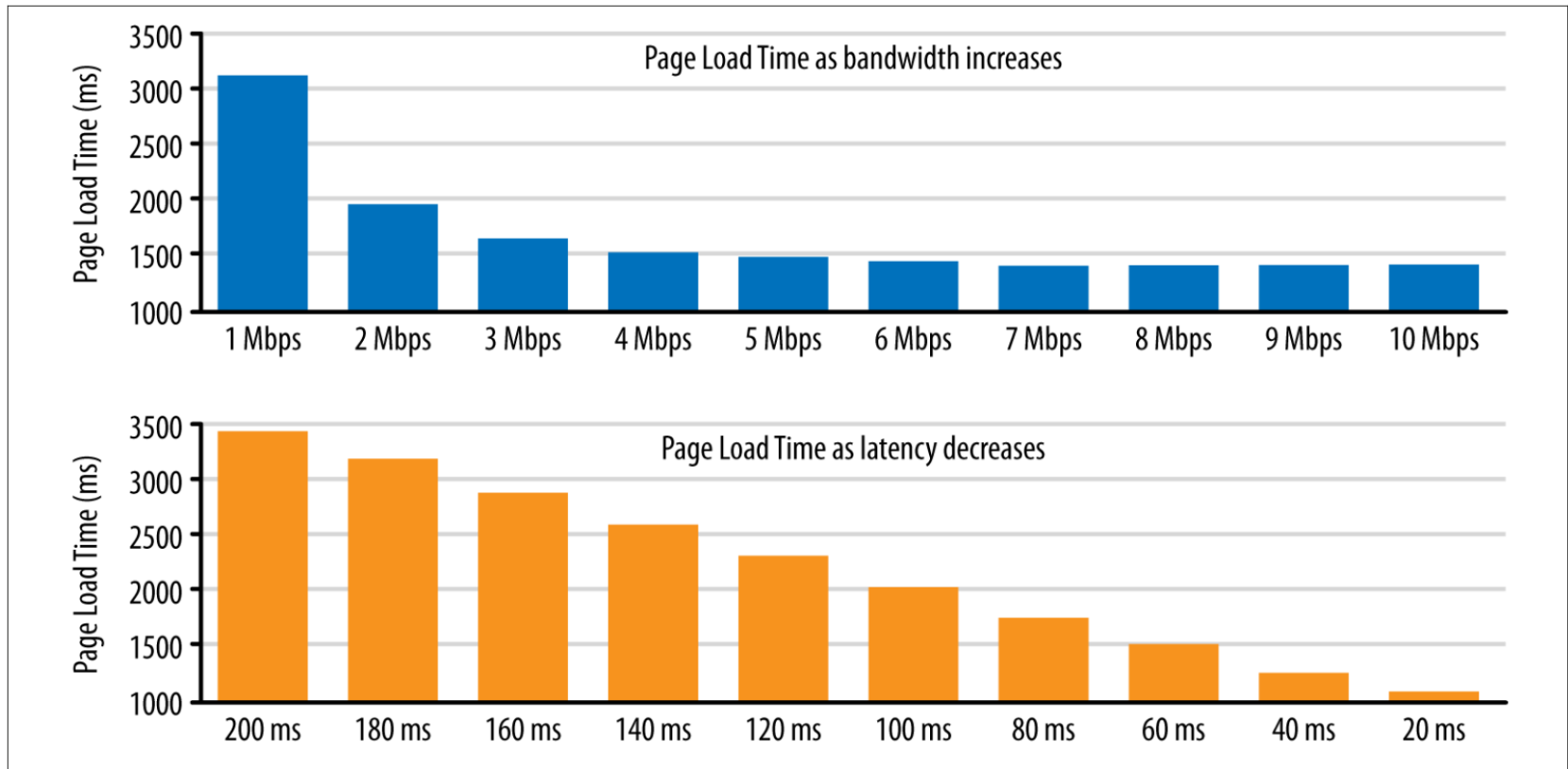
# State of the art

Two bottlenecks: latency und processing



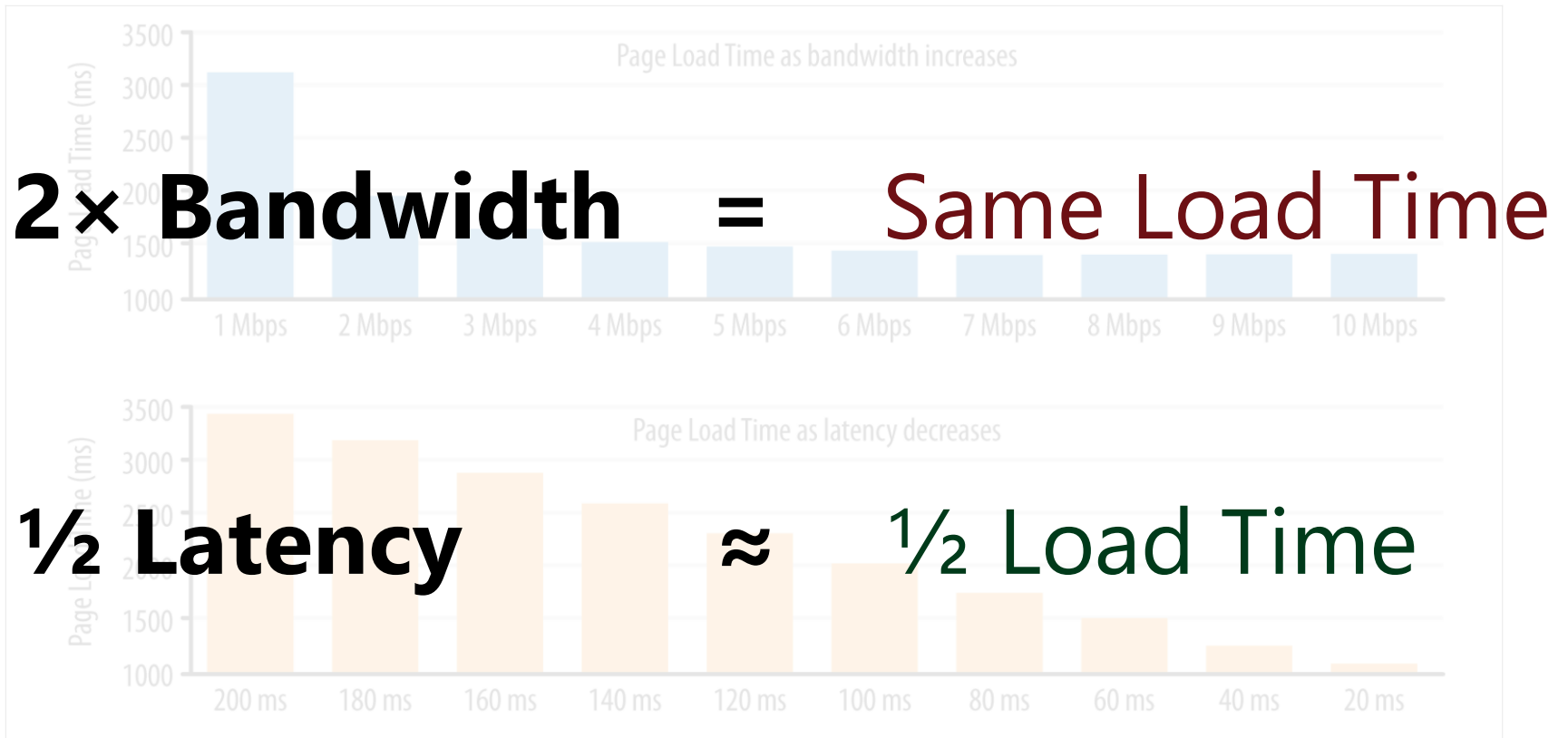
# Network Latency

The underlying problem of high page load times



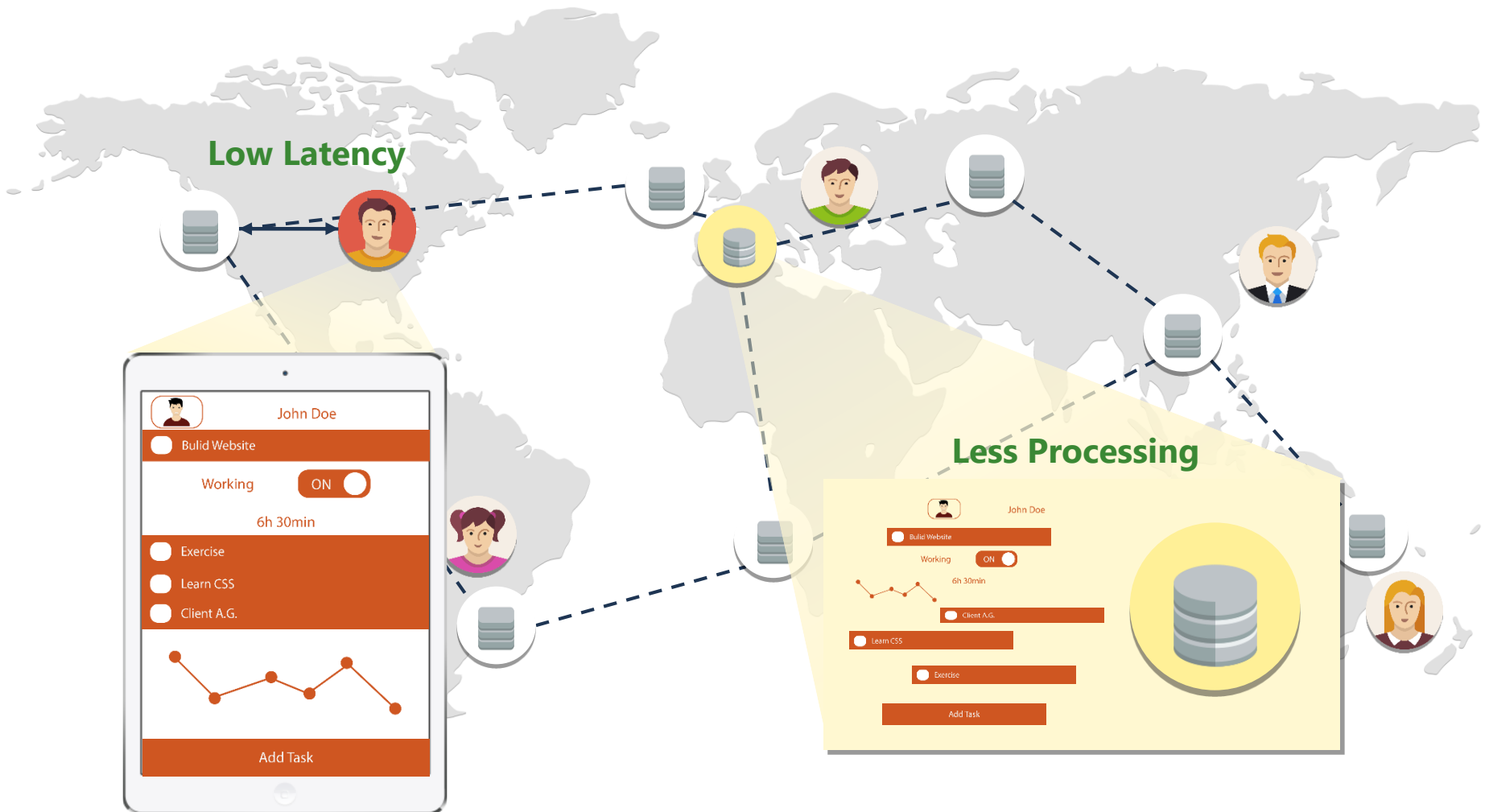
# Network Latency

The underlying problem of high page load times



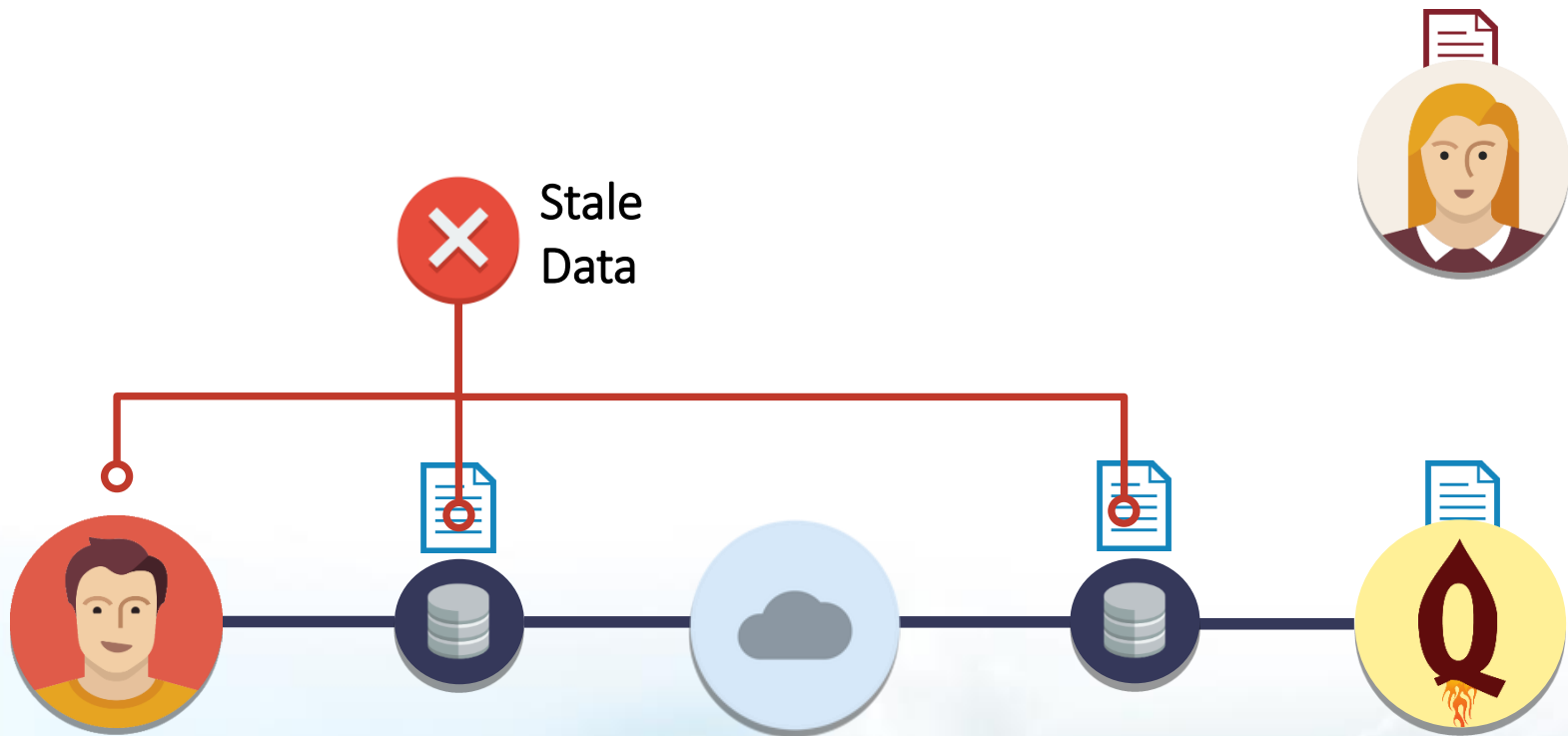
# The low-latency vision

Data is served by ubiquitous web-caches



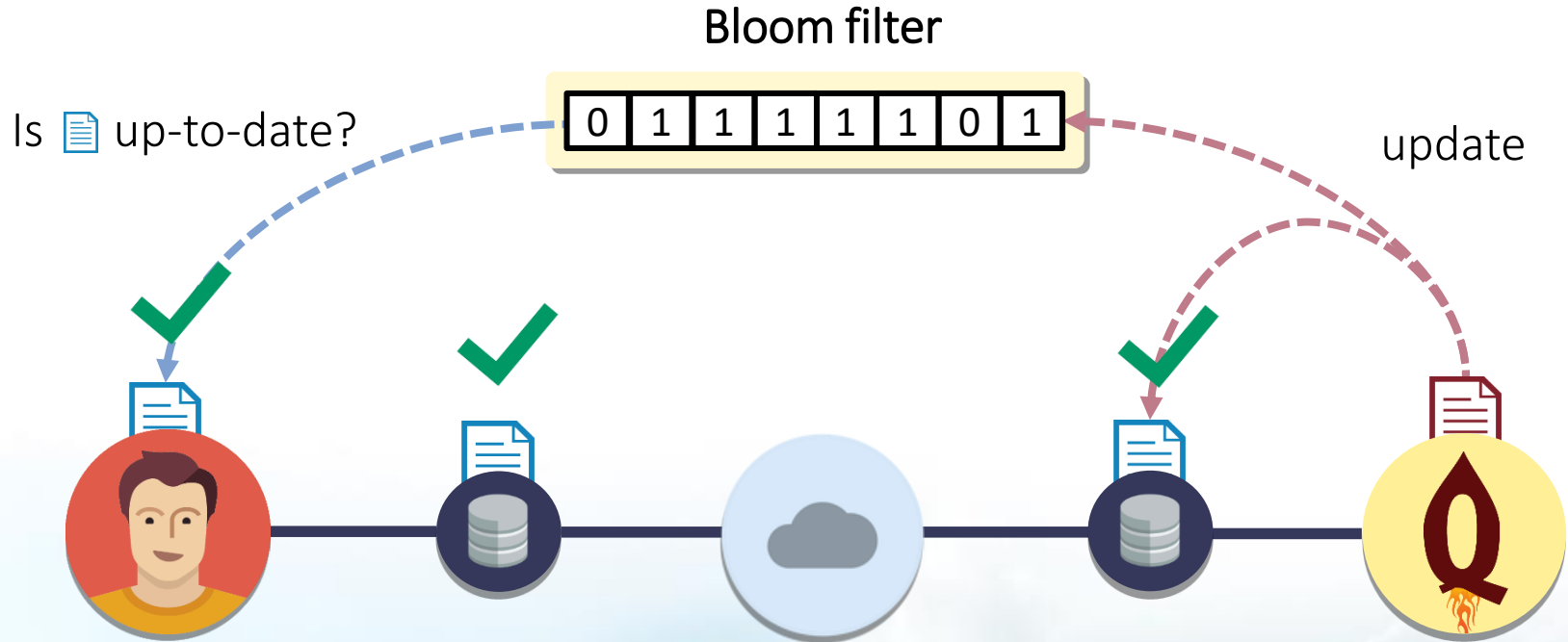
# The Problem with today's caching

Changes invalidate cached data



# Our Research

Keep Data up-to-date through Cache Sketches





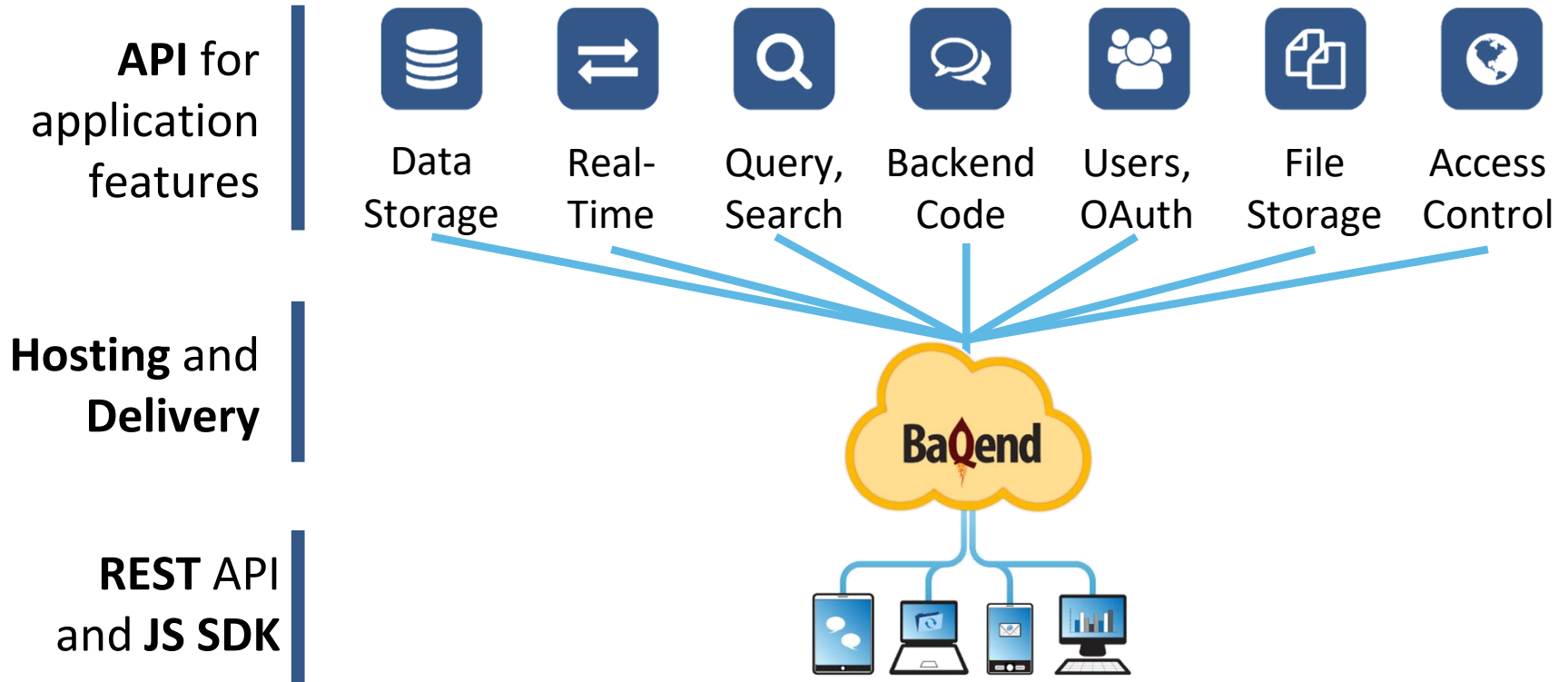
A space shuttle is shown launching vertically against a blue sky with scattered white clouds. The shuttle is white with a large orange external tank and two white boosters. A massive plume of white smoke and fire trails behind it. To the left, the launch pad's service structure is visible. The entire scene is framed by a semi-transparent white horizontal band across the middle.

# BaQend

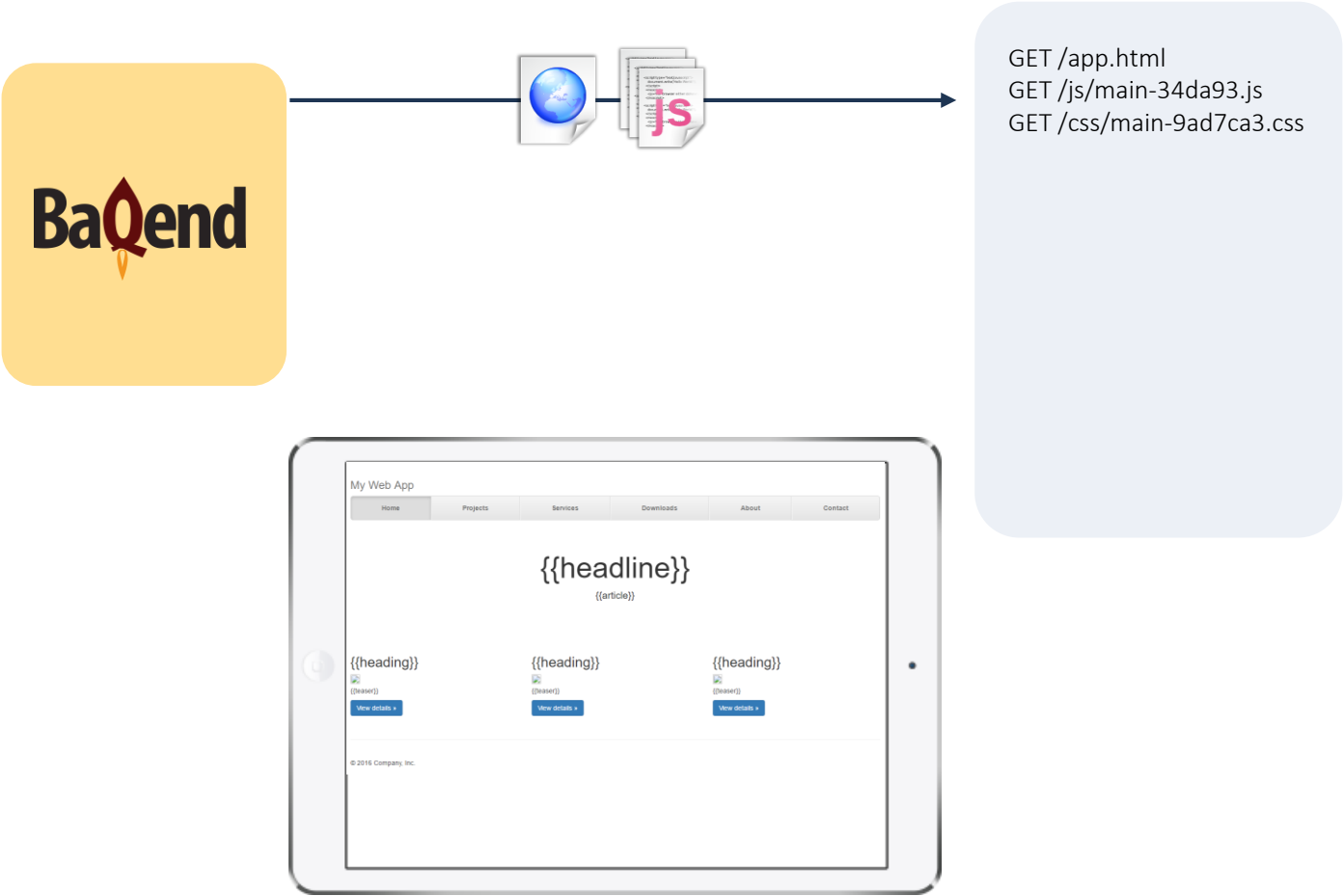
Build a faster **web**.

# Backend-as-a-Service

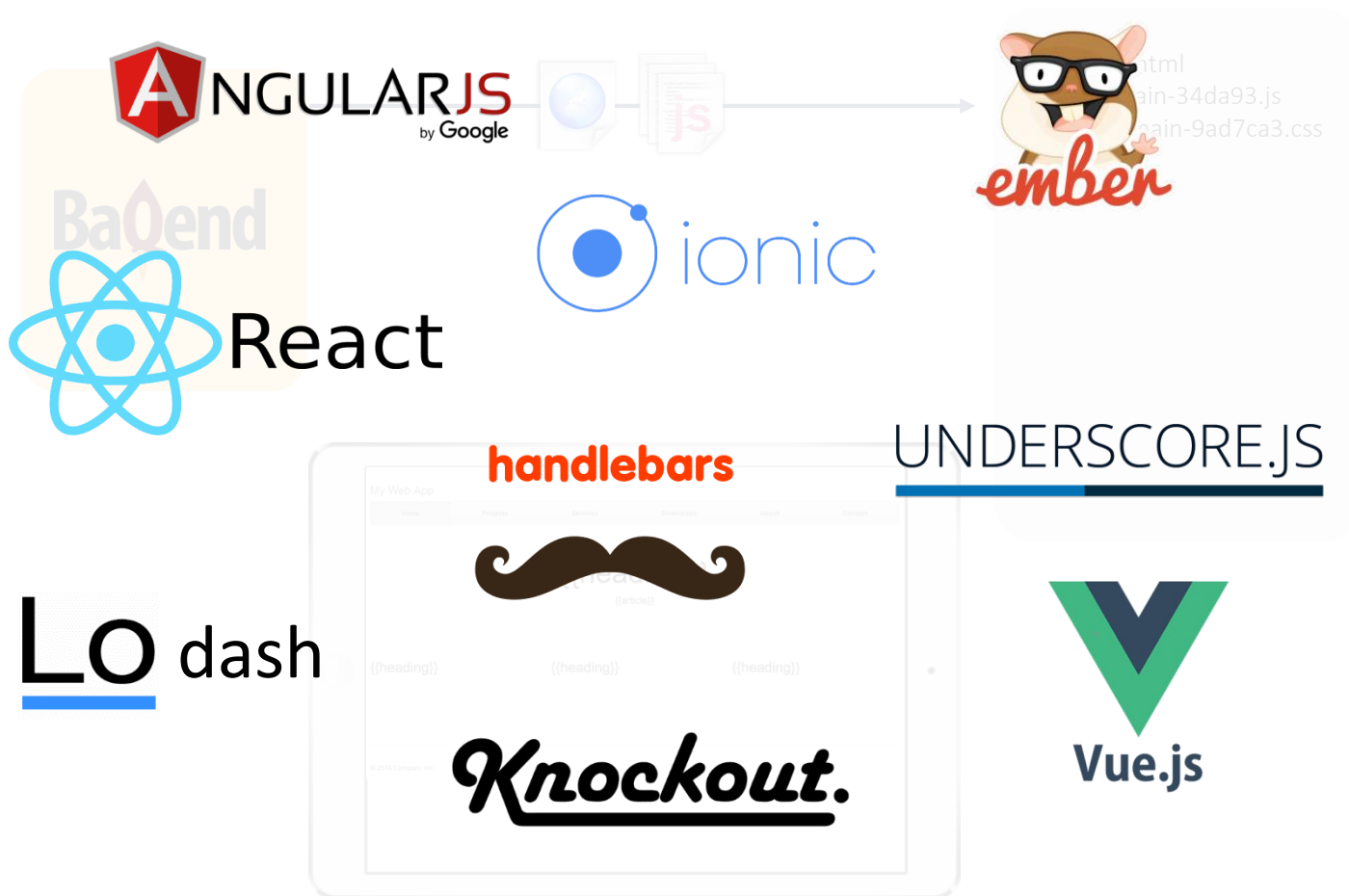
## Feature Sets



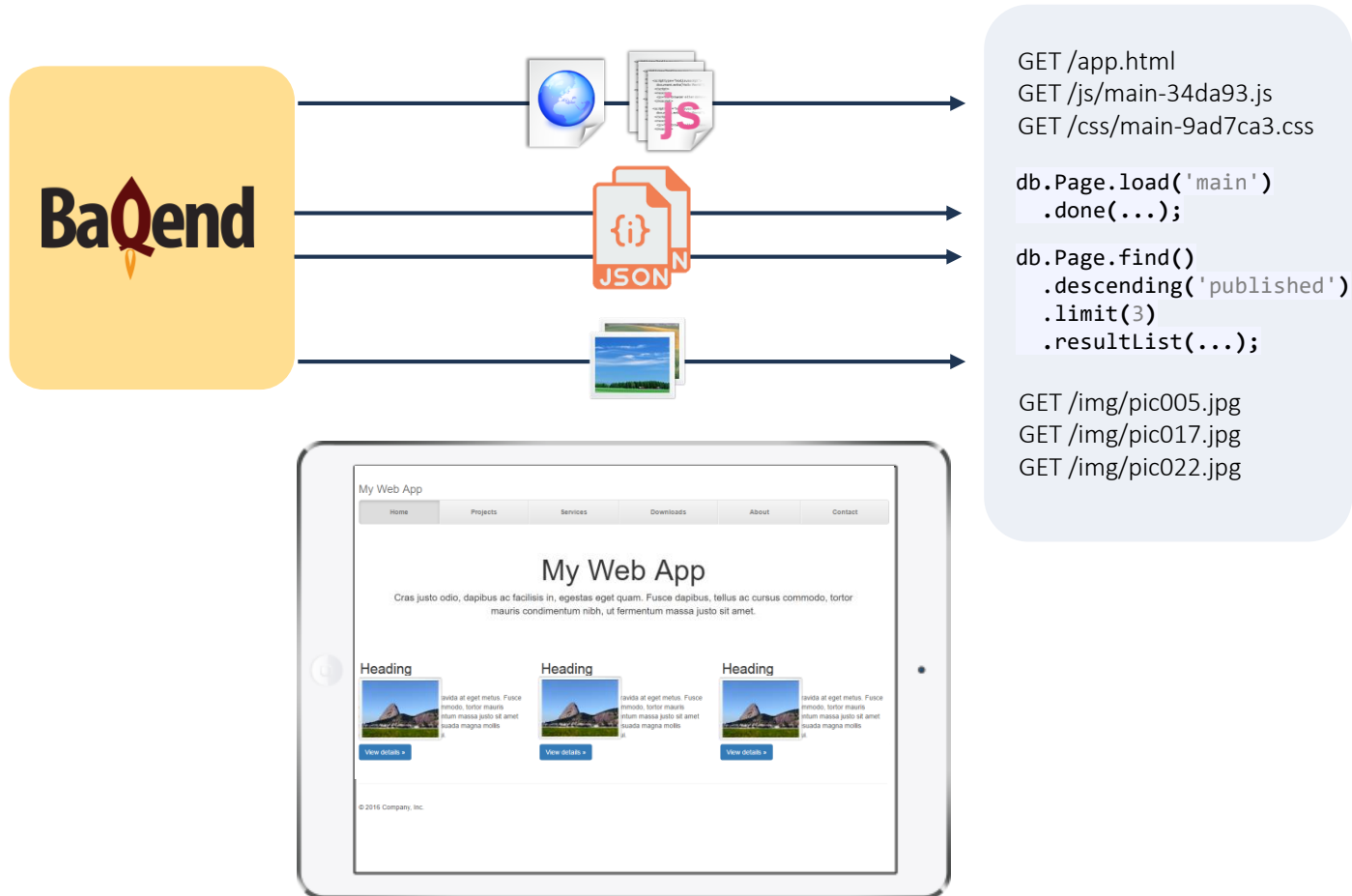
# Frontend



Compatible with:

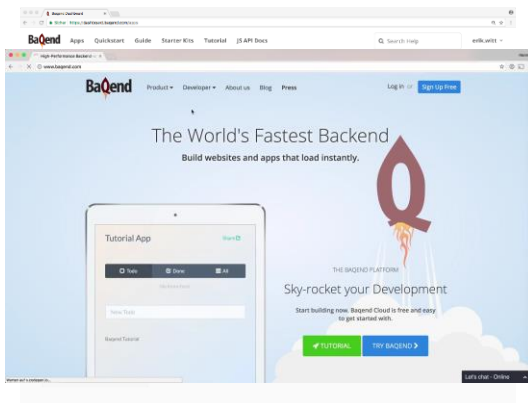


# Frontend



# Development On Baqend

## Dashboard



Create Schema, configure, browse data, etc.

## CLI

```
hannes@bq1:~$ baqend --help
Usage: baqend [command] [options] <args...>

Commands:

  login [options]           Logs you in and lo
  register                  Registers an accou
  open [app]                Opens the url to y
  dashboard                 Opens the url to t
  deploy [options] [app]   Deploys your baqen
  logout [options]         Removes your store
  typings [options] <app> Generates addition
  start [name] [dir]       clones the starter

Type in one of the above commands followed by

Options:
```

Develop, deploy and test frontend und backend Code

## REST & SDK

```
script.js
1 function leaveMessage(name, message) {
2   //Create new message object
3   var msg = new DB.Message();
4   //Set the properties
5   msg.name = name;
6   msg.message = message;
7   msg.date = new Date();
8   //Insert it to the database
9   msg.insert().then(showMessages);
10 }
11
12 function showMessages() {
13   DB.Message.find()
14     .descending("date")
15     .limit(30)
16     .resultList()
```

Website logic: load site, get data, login, etc.

# Orestes & Baqend

Learn more

If you are interested in topics combining web/mobile with scalable data management:



Bachelor and Master thesis topics  
<http://tiny.cc/orestes> (frequently updated)



Hiwi/student positions

Tasks: building real applications using cutting edge technology  
(e.g. ES6, Angular, React, MongoDB, Redis, Node.js,...)

Contact me directly or at [fg@baqend.com](mailto:fg@baqend.com)

# Summary



- ▶ Variety of different NoSQL systems:
  - ▶ **HDFS and Hadoop**: Map-Reduce platform for batch analytics
  - ▶ **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
- ▶ Cloud Databases
  - **Database-as-a-Service**: managed (NoSQL) database provided as a pay-per-use service
  - **Orestes and Baqend**: Backend-as-a-Service research project and startup with the goal of solving the web's latency problem