# Tutorial 18: Going for Speed

| | |
|---|---|
| **09:00** | **Introduction** |
| **09:15** | **Part 1: Efficient Frontend Design** |
| 09:40 | Q&A |
| 09:50 | Coffee Break |
| **10:00** | **Part 2: High-Performance Networking** |
| 10:40 | Q&A |
| 10:50 | Coffee Break |
| **11:00** | **Part 3: Scalable Backend Architectures** |
| 11:40 | Q&A |
| 11:50 | Coffee Break |
| **12:00** | **Part 4: Performance Tracking & Analysis** |
| 12:00 | The Core Web Vitals (Google Guest Speaker!) |
| 12:30 | Measuring Web Performance |
| 12:50 | Q&A |

Up next!

Baqend

# Going for Speed

## Part 3:

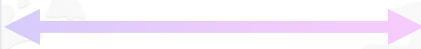## Scalable Backend Architectures

The Web Conference
April 15, 2021
Tutorial

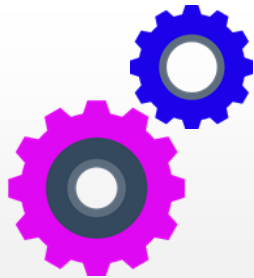Wolfram Wingerath, Benjamin Wollmer, Felix Gessert, Stephan Succo, Norbert Ritter

# The 4 **Challenges** of Web Performance

**② Network Delays**

**① Client**

**Up next!**

**③ Backend Processing**

**?**

**④ Tracking & Analysis**

# How to **Choose** a Database System

Application Layer

Billing Data

Nested Application Data

Session data

Files

Friend network

DB2

Cached data & metrics

mongoDB

Search Index

cassandra

Recommen- dation Engine

Google Cloud Storage

Neo4j the graph database

redis

elasticsearch.

Amazon Elastic MapReduce

# How to Choose a Database System

# How to get from requirements to a concrete system?

# NoSQL !?

- „NoSQL" term coined in 2009
- Interpretation: „**N**ot **O**nly **SQL**"
- Typical properties:
  - Non-relational
  - Open-Source
  - Schema-less (*schema-free*)
  - Optimized for distribution (clusters)
  - Tunable consistency

*NoSQL-Databases.org*:
Current list has over 200
NoSQL systems

# NoSQL: Two Main Motivations

## Scalability



User-generated data,
Request load

## Impedance Mismatch



ID
Customer

Line Item 1: …
Line Item2: …

Payment: Credit Card, …

Orders

Line Items

Payment

Customers

# Schema-Free Data Remodeling

**RDBMS:**



SELECT Name, Age
FROM Customers

Customers

Explicit schema

**NoSQL DB:**



Item[Price] -
Item[Discount]

Implicit schema

# Scale-Up vs. Scale-Out

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



More RAM

More CPU

More HDD

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



Commodity Hardware

Shared-Nothing Architecture

# Paradigm Shift: Open–Source & Commodity

| | |
|---|---|
| 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: Shared Nothing

Shift towards higher distribution & less coordination:



**Shared Memory**
e.g. "Oracle 11g"

**Shared Disk**
e.g. "Oracle RAC"

**Shared Nothing**
e.g. "NoSQL"

# Typical Classification Schemes

▸ Two common criteria:

**Data Model**

→ Key-Value

→ Wide-Column

→ Document

→ Graph

**Consistency/Availability Trade-Off**

→ **AP**: Available & Partition Tolerant

→ **CP**: Consistent & Partition Tolerant

→ **CA**: Not Partition Tolerant

# CAP Theorem



In a distributed system, only 2 out of the following 3 properties are achievable at the same time:

- **Consistency**: all clients have the same view on the data
- **Availability**: every request to a non-failed node most result in correct response
- **Partition tolerance**: the system has to continue working, even under arbitrary network partitions

Consistency

Partition Tolerance

Availability

Impossible

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: Intuitive Explanation

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

Block response until
ACK arrives
→ *Consistency*

Response before
successful replication
→ *Availability*

Value = $V_1$

Replication

Value = $V_0$

**N₁**

**N₂**

*Network partition*

# NoSQL Triangle

Data models | Relational
Key-Value
Wide-Column
Document-Oriented

Every client can always
read and write

A

**CA**
Oracle, MySQL, ...

**AP**
Dynamo, Redis, Riak, Voldemort
Cassandra
SimpleDB

C ──────────────── P

**CP**
Postgres, MySQL Cluster, Oracle RAC
BigTable, HBase, Accumulo, Azure Tables
MongoDB, RethinkDB, DocumentsDB

All clients share the
same view on the data

All nodes continue
working under network
partitions

# PACELC: What About Normal Operation?

▸ **Idea**: Classify systems by behavior during network partitions *and* normal operation



**PA/EL** - Dynamo-Style
Cassandra, Riak, etc.

**PA/EC** - MongoDB

**PC/EC** – Always Consistent
HBase, BigTable and ACID systems

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

# ACID vs. BASE

## ACID
,,Gold standard''
for RDBMSs

- **A**tomicity
- **C**onsistency
- **I**solation
- **D**urability

## BASE
Model of many
NoSQL systems

- **B**asically **A**vailable
- **S**oft State
- **E**ventually Consistent

# NoSQL Landscape

**Document**

Amazon DynamoDB

H·BASE

HYPERTABLE

Google Datastore

Cassandra

mongoDB

**Wide Column**

RethinkDB

**Key-Value**

redis

CouchDB relax

riak

amazon web services | S3

Couchbase

RAVENDB

**Graph**

Project Voldemort

Neo4j the graph database

InfiniteGraph Powered by Objectivity

AEROSPIKE

# Requirements vs. Techniques

**Functional**

Scan Queries

ACID Transactions

Conditional or Atomic Writes

Joins

Sorting

Filter Queries

Full-text Search

Aggregation and Analytics

**Techniques**

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

**Replication**
Commit/Consensus Protocol
Synchronous
Asynchronous
Primary Copy
Update Anywhere

**Storage Management**
Logging
Update-in-Place
Caching
In-Memory Storage
Append-Only Storage

**Query Processing**
Global Secondary Indexing
Local Secondary Indexing
Query Planning
Analytics Framework
Materialized Views

**Non-Functional**

Data Scalability

Write Scalability

Read Scalability

Elasticity

Consistency

Write Latency

Read Latency

Write Throughput

Read Availability

Write Availability

Durability

# Requirements vs. Techniques



Functional

*enable* → Techniques ← *enable*

Non-Functional

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

**application requirements**

**Pivotal NoSQL techniques**

**operational requirements**

**Replication**
- Commit/Consensus Protocol
- Synchronous
- Asynchronous
- Primary Copy
- Update Anywhere

**Storage Management**
- Logging
- Update-in-Place
- Caching
- In-Memory
- Append-Only

**Query**
- Global Secondary Indexing
- Local Secondary Indexing
- Query Planning
- Analytics Framework
- Materialized Views

Scan Queries

ACID Transactions

Conditional or Atomic Writes

Joins

Sorting

Queries

Search

Aggregation and Analytics

Data Scalability

Write Scalability

Read Scalability

Elasticity

Consistency

Write Latency

Read Latency

Write Throughput

Read Throughput

Write Availability

Durability

**Functional**

Scan Queries

ACID Transactions

Conditional or Atomic Writes

Joins

Sorting

**Techniques**

Sharding

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

**Non-Functional**

Data Scalability

Write Scalability

Read Scalability

Elasticity

# Sharding : Scaling Storage & 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

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

## Range-based Sharding

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

Implemented in

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

## Entity-Group Sharding

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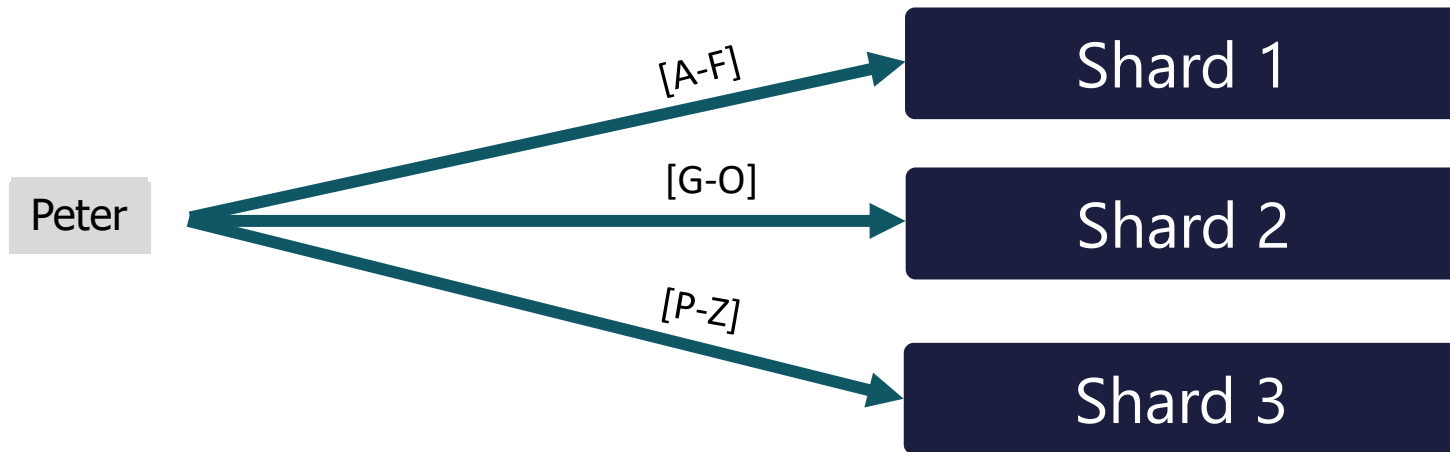- Pro: Enables *ACID Transactions*
- Contra: Partitioning not easily changable

Implemented in

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

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.

Functional | Techniques | Non-Functional

ACID Transactions

Conditional or Atomic Writes

**Replication**
Commit/Consensus Protocol
Synchronous
Asynchronous
Primary Copy
Update Anywhere

Read Scalability

Consistency

Write Latency

Read Latency

Read Availability

Write Availability

# Replication : Read Scalability & Fault Tolerance

▸ Stores *N* copies of each data item



▸ **Consistency model**: synchronous vs asynchronous

▸ **Coordination**: Multi-Master, Master-Slave

Özsu, M.T., Valduriez, P.: Principles of distributed database systems. Springer Science & Business Media (2011)

# Replication : When

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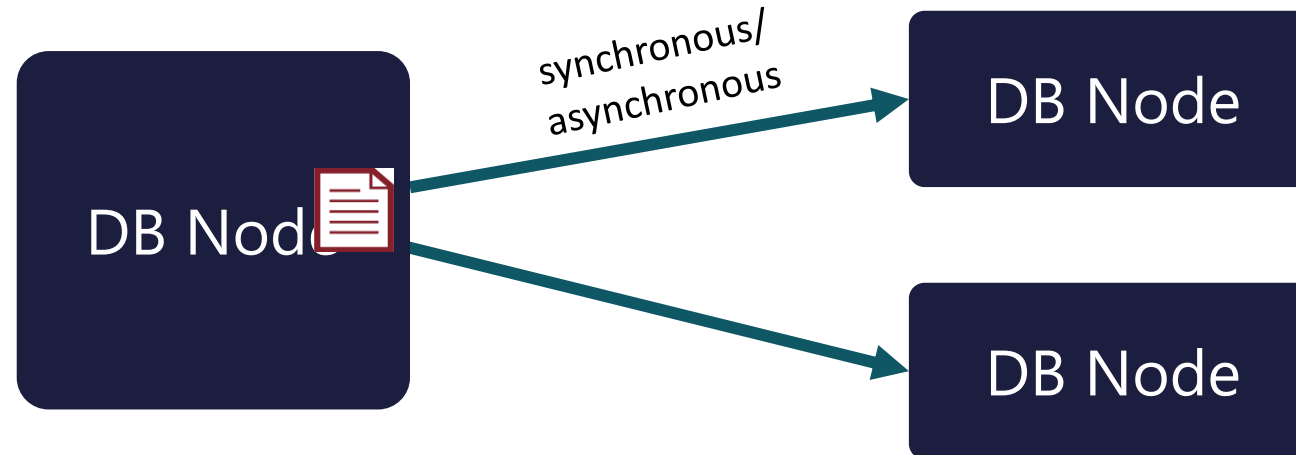- Writes are acknowledged immdediately
- Performed through *log shipping* or *update propagation*
- **Pro**: Fast writes, no coordination needed
- **Contra**: Replica data potentially stale (*inconsistent*)

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

**Implemented in**

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

**Implemented in**

BigTable, HBase, Accumulo, CouchBase, MongoDB, RethinkDB

Charron-Bost, B., Pedone, F., Schiper, A. (eds.): Replication: Theory and Practice, Lecture Notes in Computer Science, vol. 5959. Springer (2010)

# 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
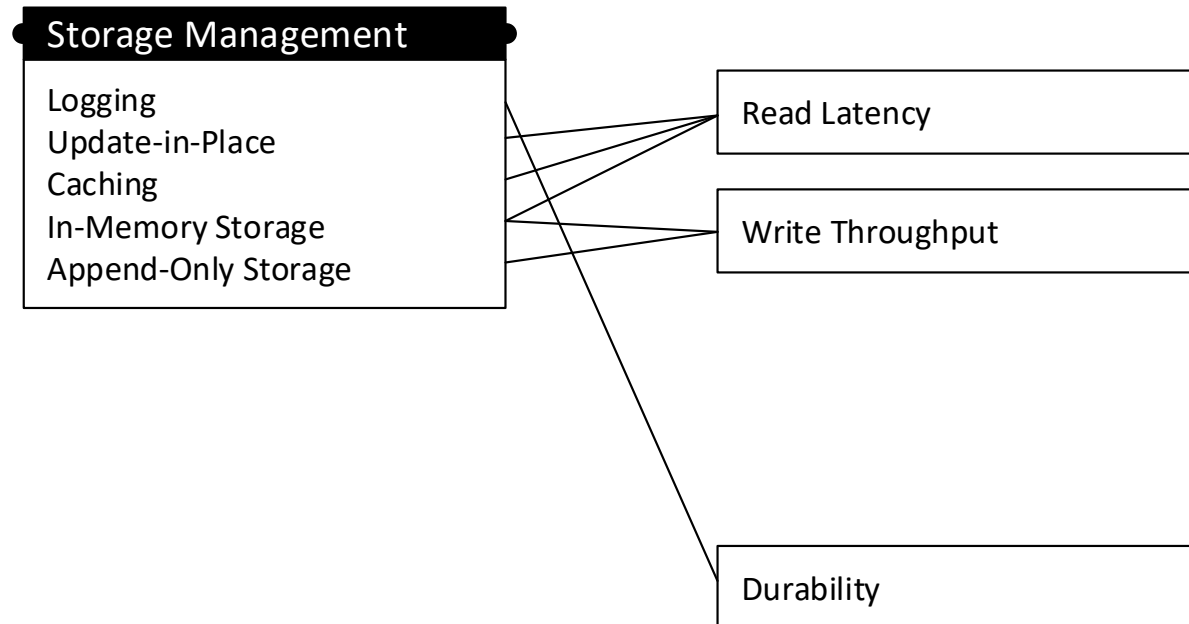
## Multi-Master (*Update anywhere*)

◦ 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
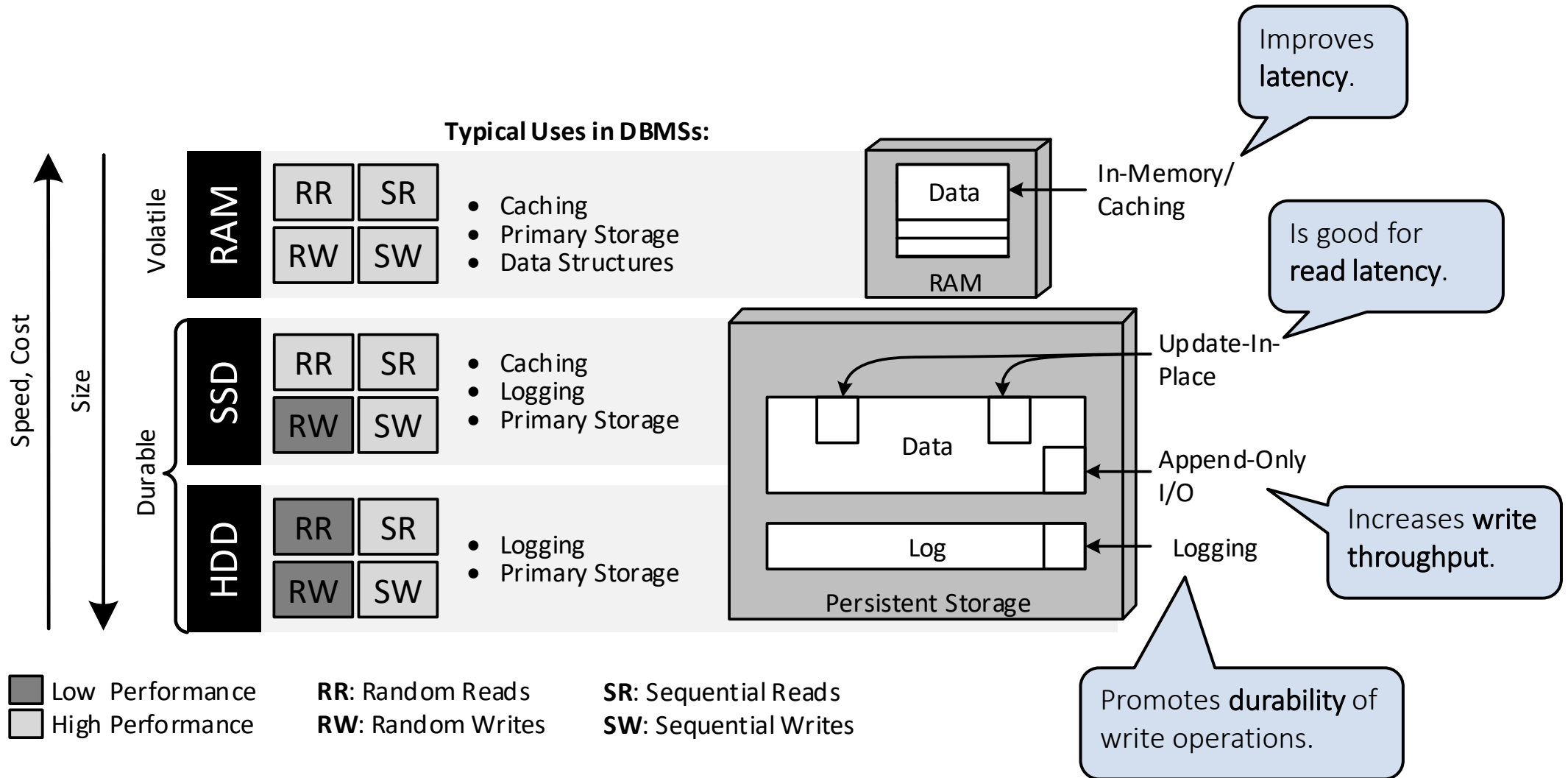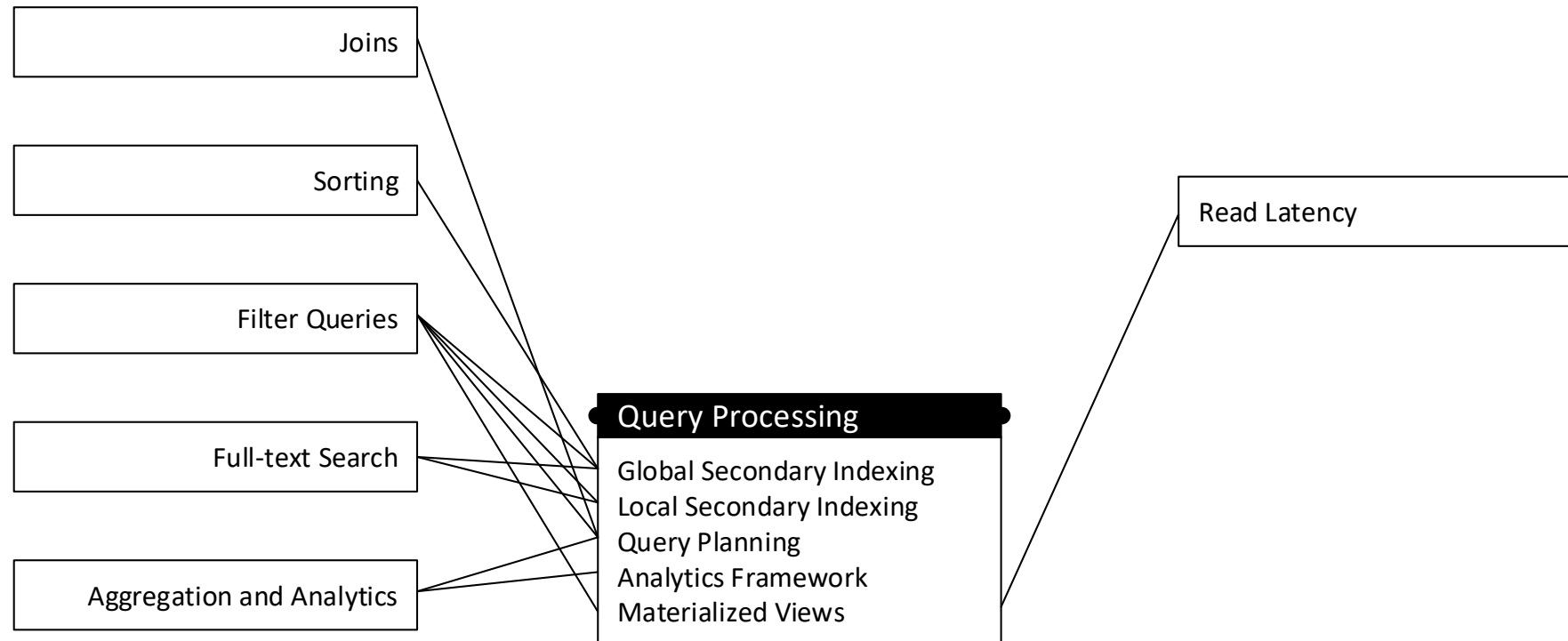
Charron-Bost, B., Pedone, F., Schiper, A. (eds.): Replication: Theory and Practice, Lecture Notes in Computer Science, vol. 5959. Springer (2010)

Functional                    Techniques                    Non-Functional

**Storage Management**

Logging
Update-in-Place
Caching
In-Memory Storage
Append-Only Storage

Read Latency

Write Throughput

Durability

Functional                          Techniques                          Non-Functional

| Joins |

| Sorting |

| Filter Queries |

| Full-text Search |

| Aggregation and Analytics |

| Read Latency |

**Query Processing**

Global Secondary Indexing
Local Secondary Indexing
Query Planning
Analytics Framework
Materialized Views

# Local Secondary Indexing

## Partition I

**Data**

| Key | Color |
|-----|-------|
| 12 | Red |
| 56 | Blue |
| 77 | Red |

**Index**

| Term | Match |
|------|-------|
| Red | [12,77] |
| Blue | [56] |

## Partition II

**Data**

| Key | Color |
|-----|-------|
| 104 | Yellow |
| 188 | Blue |
| 192 | Blue |

**Index**

| Term | Match |
|------|-------|
| Yellow | [104] |
| Blue | [188,192] |

Indexing is always local to a partition.

## Implemented in

- MongoDB
- Riak
- Cassandra
- Elasticsearch
- SolrCloud
- VoltDB

Scatter-gather query pattern.

```
WHERE color=blue
```

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

# Global Secondary Indexing

## Partition I

**Data**

| Key | Color |
|-----|-------|
| 12  |       |
| 56  |       |
| 77  |       |

Consistent Index-maintenance requires distributed transaction.

**Index**

| Term   | Match           |
|--------|-----------------|
| Yellow | [104]           |
| Blue   | [56, 188, 192]  |

## Partition II

**Data**

| Key | Color  |
|-----|--------|
| 104 | Yellow |
| 188 | Blue   |
| 192 | Blue   |

**Index**

| Term | Match    |
|------|----------|
| Red  | [12,77]  |

Implemented in

- DynamoDB
- Oracle Datawarehouse
- Riak (Search)
- Cassandra (Search)

Targeted Query

WHERE color=blue

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

# **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

▸ **Materialized Views**: similar to global indexing

# NoSQL **Decision** Tree



- **Access**
  - Fast Lookups
    - **Volume**
      - RAM
        - **Redis**
          Memcache
      - Unbounded
        - **CAP**
          - AP
            - **Cassandra Riak**
              Voldemort
              Aerospike
          - CP
            - **HBase**
              MongoDB
              CouchBase
              DynamoDB
  - Complex Queries
    - **Volume**
      - HDD-Size
        - **Consistency**
          - ACID
            - **RDBMS**
              Neo4j
              RavenDB
              MarkLogic
          - Availability
            - **CouchDB MongoDB**
              SimpleDB
      - Unbounded
        - **Query Pattern**
          - Ad-hoc
            - **MongoDB**
              RethinkDB
              HBase,Accumulo
              ElasticSeach, Solr
          - Analytics
            - **Hadoop, Spark Parallel DWH**
              Cassandra, HBase
              Riak, MongoDB

Example Applications

| Cache | Shopping-basket | Order History | OLTP | Website | Social Network | Big Data |

# What About **Push**-Based Systems?

# A Data Processing Pipeline



Persistence/
Streaming

Processing

Serving

Application

# Scale-Out Made Feasible

Data processing frameworks **hide complexities of scaling**, e.g.:

- **Deployment -** code distribution, starting/stopping work
- **Monitoring -** health checks, application stats
- **Scheduling -** assigning work, rebalancing
- **Fault-tolerance** - restarting workers, rescheduling failed work

Running in cluster

Running on single node

Scaling out

# Big Data Processing Models

**stream**                    **micro-batch**                    **batch**

# Stream Processing System Comparison

| | Storm | Trident | Samza | Spark Streaming | Flink (streaming) |
|---|---|---|---|---|---|
| **Strictest Guarantee** | at-least-once | exactly-once | at-least-once | exactly-once | exactly-once |
| **Achievable Latency** | ≪100 ms | <100 ms | <100 ms | <1 second | <100 ms |
| **State Management** | ◯ (small state) | ◯ (small state) | ✓ | ✓ | ✓ |
| **Processing Model** | one-at-a-time | micro-batch | one-at-a-time | micro-batch | one-at-a-time |
| **Backpressure** | ✓ | ✓ | no (buffering) | ✓ | ✓ |
| **Ordering** | ✗ | between batches | within partitions | between batches | within partitions |
| **Elasticity** | ✓ | ✓ | ✗ | ✓ | ✗ |

Real-Time Databases: Combining Push & Pull

# Push vs. Pull: Trade-Offs in Data Management

**Database Management**

**Real-Time Databases**

**Datastream Management**

query semantics

delivery semantics

static collections

dynamic collections

dynamic data streams

Pull-Based

Push-Based

# Real-Time Databases: Always Up-to-Date

circular shapes

{ }

**Real-Time Queries** for query result maintenance:
→ efficient
→ fast

# Real-Time Database Comparison

| | METEOR | RethinkDB | Parse | Firebase |
|---|:---:|:---:|:---:|:---:|
| | Poll-and-Diff | Change Log Tailing | | Unknown |
| **Write Scalability** | ✓ | ✗ | ✗ | ✗ | ✗ |
| **Read Scalability** | ✗ | ✓ | ✓ | ✓ | ? (100k connections) |
| Composite Filters (AND/OR) | ✓ | ✓ | ✓ | ✓ | ◯ (AND In Firestore) |
| Sorted Queries | ✓ | ✓ | ✓ | ✗ | ◯ (single attribute) |
| Limit | ✓ | ✓ | ✓ | ✗ | ✓ |
| Offset | ✓ | ✓ | ✗ | ✗ | ◯ (value-based) |
| Self-Maintaining Queries | ✓ | ✓ | ✗ | ✗ | ✗ |
| Event Stream Queries | ✓ | ✓ | ✓ | ✓ | ✓ |

W. Wingerath, F. Gessert, N. Ritter: InvaliDB: *Scalable Push-Based Real-Time Queries on Top of Pull-Based Databases (Extended)*, VLDB 2020

**Wrapup**: A Short History of Data Management

For videos & books,
visit slides.baqend.com!