# **Real-Time Data Management**

## For Big Data

Wolfram Wingerath, Felix Gessert, Norbert Ritter

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

March 29, EDBT 2018, Vienna

Universität Hamburg

**BaQend**

www.baqend.com

# Who We Are

**Norbert Ritter**
Professor

**Felix Gessert**
CEO

**Wolfram Wingerath**
Developer

**Research:**
- NoSQL & Cloud Databases
- Polyglot Persistence
- Database Benchmarking
- …

**Practice**:
Backend-as-a-Service •
Web Caching •
Real-Time Database •
… •

UH
Universität Hamburg

BaQend
www.baqend.com

# Outline

**Introduction**
Where From? Where To?

**Stream Processing**
Big Data + Low Latency
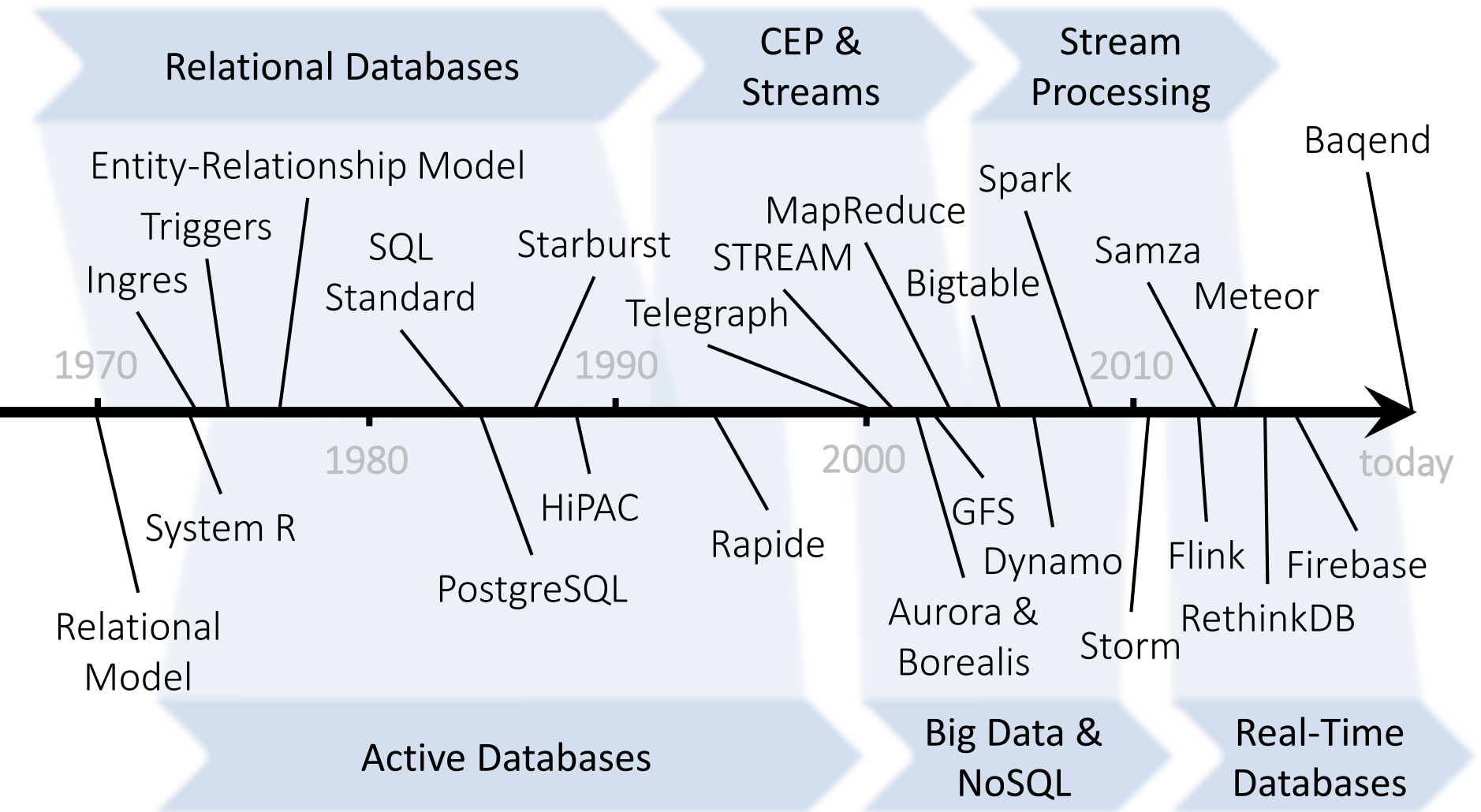
**Real-Time Databases**
Push-Based Collections

**Future Directions**
Current Research & Outlook

- A Short History of Data Management
- **Database Management:**
  - (No)SQL Decision Tree
  - (No)SQL Toolbox
  - Active Database Features
- **Data Stream Management:**
  - General Architecture
  - Stream Operators
  - Approximation & Sampling
  - CEP

# A Short History of Data Management
## Hot Topics Through The Ages

CONCEPTS & REQUIREMENTS

# The NoSQL Toolbox

# 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.
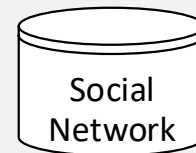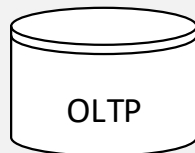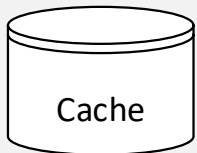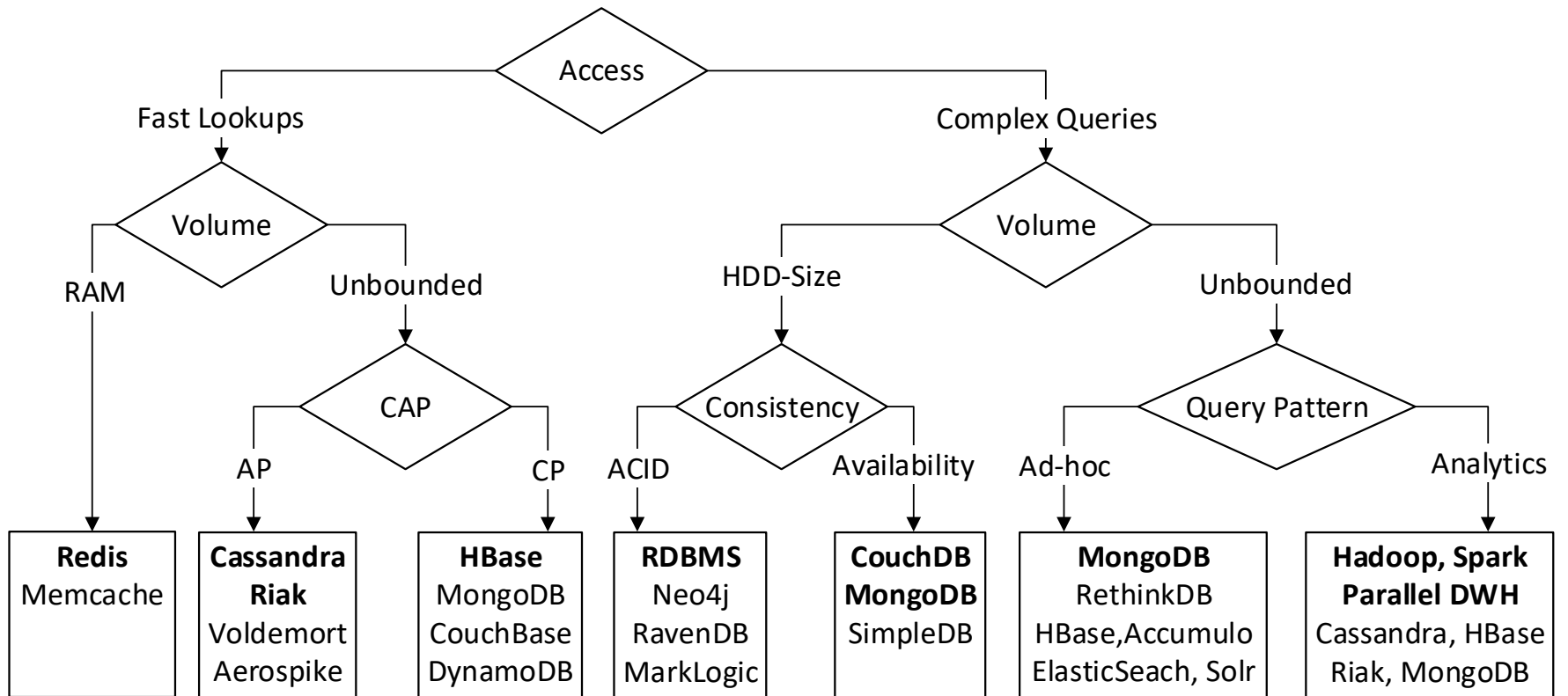
## 1  Introduction

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
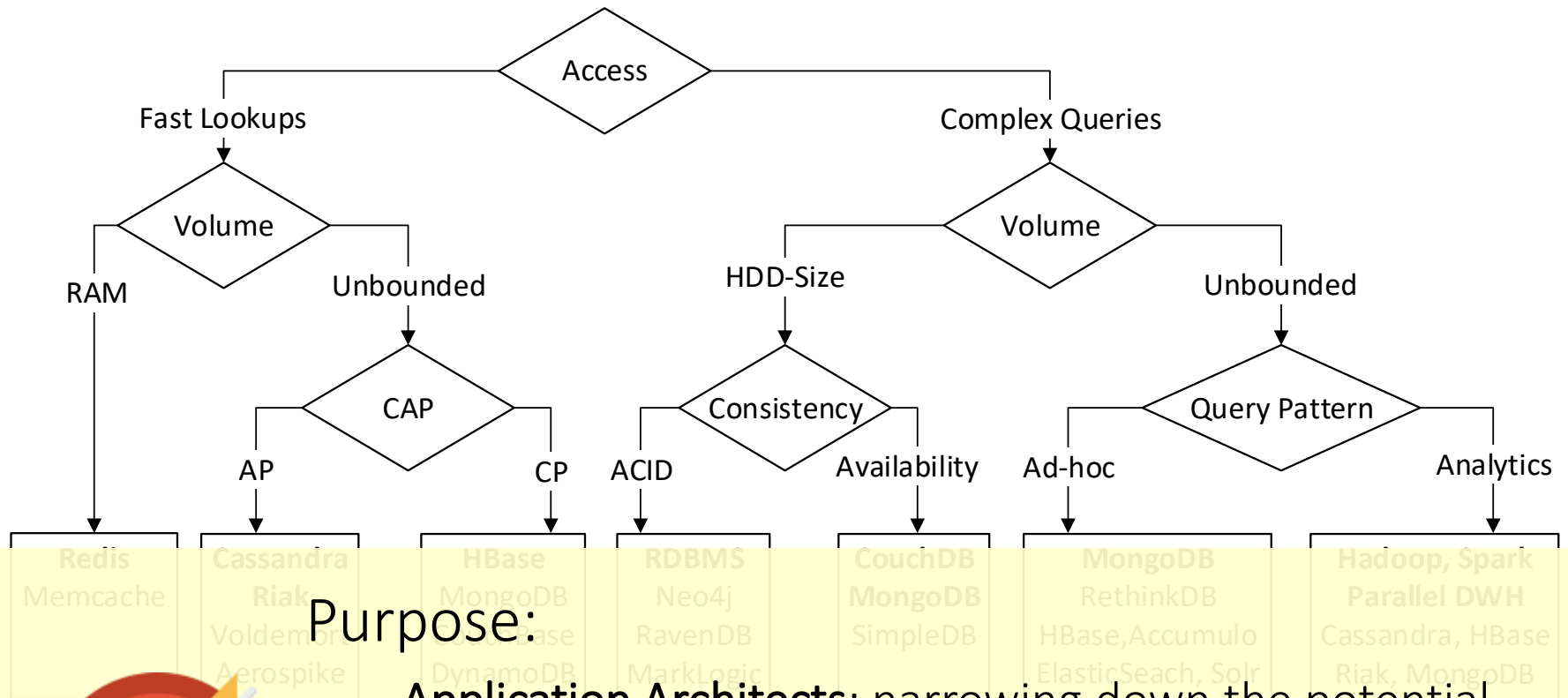
# (No)SQL Decision Tree

**Access**

Fast Lookups — Complex Queries

**Fast Lookups → Volume**
- RAM
- Unbounded

**Complex Queries → Volume**
- HDD-Size
- Unbounded

**Unbounded (Fast Lookups) → CAP**
- AP
- CP

**HDD-Size → Consistency**
- ACID
- Availability

**Unbounded (Complex Queries) → Query Pattern**
- Ad-hoc
- Analytics

| **Redis**<br>Memcache | **Cassandra Riak**<br>Voldemort<br>Aerospike | **HBase**<br>MongoDB<br>CouchBase<br>DynamoDB | **RDBMS**<br>Neo4j<br>RavenDB<br>MarkLogic | **CouchDB MongoDB**<br>SimpleDB | **MongoDB**<br>RethinkDB<br>HBase,Accumulo<br>ElasticSeach, Solr | **Hadoop, Spark Parallel DWH**<br>Cassandra, HBase<br>Riak, MongoDB |
|---|---|---|---|---|---|---|

## Example Applications

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

# (No)SQL Decision Tree



**Access**

Fast Lookups — Complex Queries

**Volume** (Fast Lookups)
- RAM
- Unbounded

**Volume** (Complex Queries)
- HDD-Size
- Unbounded

**CAP**
- AP
- CP

**Consistency**
- ACID
- Availability

**Query Pattern**
- Ad-hoc
- Analytics

| Redis Memcache | Cassandra Riak Voldemort Aerospike | HBase MongoDB CouchBase DynamoDB | RDBMS Neo4j RavenDB MarkLogic | CouchDB MongoDB SimpleDB | MongoDB RethinkDB HBase,Accumulo ElasticSeach, Solr | Hadoop, Spark Parallel DWH Cassandra, HBase Riak, MongoDB |

Purpose:

**Application Architects**: narrowing down the potential system candidates based on requirements

**Database Vendors/Researchers**: clear communication and design of system trade-offs

Example Applications

| Functional | Techniques | Non-Functional |
|---|---|---|
| Scan Queries | **Sharding** | Data Scalability |
| ACID Transactions | Range-Sharding<br>Hash-Sharding<br>Entity-Group Sharding<br>Consistent Hashing<br>Shared-Disk | Write Scalability |
| Conditional or Atomic Writes | | Read Scalability |
| Joins | **Replication** | Elasticity |
| Sorting | Commit/Consensus Protocol<br>Synchronous<br>Asynchronous<br>Primary Copy<br>Update Anywhere | Consistency |
| Filter Queries | | Write Latency |
| Full-text Search | **Storage Management** | Read Latency |
| Aggregation and Analytics | Logging<br>Update-in-Place<br>Caching<br>In-Memory Storage<br>Append-Only Storage | Write Throughput |
| | **Query Processing** | Read Availability |
| | Global Secondary Indexing<br>Local Secondary Indexing<br>Query Planning<br>Analytics Framework<br>Materialized Views | 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 (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

### 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 changable

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) d
- **Pro**: Even distribution
- **Contra**: No data locality

### Range-based Sharding
- Assigns ranges defined over fie
- **Pro**: Enables *Range Scans* and *S*
- **Contra**: Repartitioning/balancing

### Entity-Group Sharding
- Explicit data co-location for sin
- **Pro**: Enables *ACID Transactions*
- **Contra**: Partitioning not easily

**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

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 + Failure 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 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), unavaialable under certain network partitions

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

# Replication: When

**Asynchronous** (lazy)

- Writes are acknowledged imm...
- Performed through *log shippi...
- **Pro**: Fast writes, no coordinati...
- **Contra**: Replica data potential...

**Implemented in**

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

**Synchronous** (eager)

- The node accepting writes syn...tes updates/transactions before a...
- **Pro**: Consistent
- **Contra**: needs a commit prot... unavaialable under certain network partitions

**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

# NoSQL Storage Management
## In a Nutshell

**Typical Uses in DBMSs:**

**Speed, Cost** (vertical axis, double arrow)

**Size** (vertical)

**Volatile** / **RAM**

| RR | SR |
|----|----|
| RW | SW |

- Caching
- Primary Storage
- Data Structures

**Durable**

**SSD**

| RR | SR |
|----|----|
| RW | SW |

- Caching
- Logging
- Primary Storage

**HDD**

| RR | SR |
|----|----|
| RW | SW |

- Logging
- Primary Storage

Data

RAM → In-Memory/ Caching

Data

Update-In-Place

Append-Only I/O

Log → Logging

Persistent Storage

Low Performance
High Performance

**RR**: Random Reads  **SR**: Sequential Reads
**RW**: Random Writes  **SW**: Sequential Writes

# NoSQL Storage Management
## In a Nutshell

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

# 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

# Summary

- High-Level Database Categories:
  - Relational, Key-Value, Wide-Column, Document, Graph
  - Two out of {Consistent, Available, Partition Tolerant}
- The (**No)SQL Toolbox**: systems use similar techniques that promote certain capabilities

**Techniques**
*Sharding, Replication,*
*Storage Management,*
*Query Processing*

promote

**Functional**
Requirements

**Non-functional**
Requirements

- **Decision Tree**: maps requirements to concrete systems

# Active Database Features

# Databases are <u>Passive</u>
## Challenge: How to Build <u>Reactive</u> Applications?

circular shapes

**Change discovery** through periodic polling
→ Inefficient
→ Slow

# Active Database Features
## Modeling Behavioral Domain Aspects

**Triggers**: simple action-mechanisms
- Use cases:
  - (Referential) integrity
  - Change data capture

**ECA rules**: <u>E</u>vent-<u>C</u>ondition-<u>A</u>ction
- Captures **composite events**
- More expressive than triggers (**rule languages**)
- Advanced use cases:
  - Materialized view maintenance
  - Pattern recognition
  - (complex) event processing

# View Maintenance
Keeping Track of Query Results

**Materialized Views:** precomputed query results
- ◦ Used to speed up pull-based queries, e.g in data warehouses
- ◦ Implementation aspects:
  - Eager vs. lazy
  - Incremental vs. recomputation-based
  - Partial maintenance vs. full maintenance
  - Self-maintainability vs. expressiveness

**Change Notification Mechanisms**: inform subscribers of possibly invalidated query results
- ◦ Used to invalidate caches in the middle tier (cf. 3-tier stack)

# View Maintenance By Example
## Matching Every Query Against Every Update

Similar processing for:
- Triggers
- ECA rules

→ Potential *bottlenecks*:
- *Number of queries/triggers/rules*
- *Write throughput*
- *Complexity*

Is match?

yes     no

Was match?          Was match?

yes     no          yes     no

change    add    remove    none

EVOLVING DOMAINS

# Data Stream Management

# Push-Based Access For Evolving Domains

## Continuous Queries Over Data Streams

Find people in Room B:

```
SELECT name, x, y
  FROM People
  WHERE x BETWEEN 0 AND 25
    AND y BETWEEN 0 AND 15
  ORDER BY name ASC
```

1. 🔴 Erik (5/10)
2. 🟢 Wolle (21/4)
3.

# Data Stream Management Systems

## High-Level Architecture

stream query
processor

archive
(offline)

working memory

database

# Typical Stream Operators
## Examples

**Filter & Transform**



Filter          Map

**Group**



GroupByKey

**Aggregates**

SUM()
COUNT()

**Windows**



Tumbling



Sliding

# Complex Event Processing
## Detecting Patterns

▸ **Abstraction** from raw event streams

▸ Detection of **relationships** between events

▸ Often modeled in abstraction **hierarchies**

▸ Techniques:

  ◦ Transformation, filtering

  ◦ Correlation, aggregation, …

  ◦ Pattern detection
    → **composite events**



complex events

low-level events

abstraction

event patterns

# Notions of Time
## Arrival Time vs. Event Time

▸ **Arrival time**: When was the event <u>received</u>?
▸ **Event time**: When did the event <u>occur</u>?
▸ **Clock Skew**: difference between arrival and event time

# Approximation & Load Shedding
Provide the „Best" Answer While Avoiding to Fall Behind



raw stream

Prohibitive!

# Approximation & Load Shedding
Provide the „Best" Answer While Avoiding to Fall Behind

▸ **Sampling**: can be optimized for different things, e.g.
  ◦ Position stream (e.g. „select every 10th item")
  ◦ Value (e.g. hash partitioning)
  ◦ Semantic criteria

raw stream

Sampled stream

# Summary

| | Database | Stream |
|---|---|---|
| Update rate | Low | High, bursty |
| Primitive | Persistent collections | Transient streams |
| Temporal scope | Historical | Windowed |
| Access | random | sequential |
| Queries | One-time | Continuous |
| Query Plans | Static | Dynamic |
| Precision | Accurate | Approximate |

# Outline

$\sum$ **Introduction**
Where From? Where To?

⚙ **Stream Processing**
Big Data + Low Latency

🗄 **Real-Time Databases**
Push-Based Collections

💡 **Future Directions**
Current Research & Outlook

- **Big Picture:**
  - Processing Pipelines
  - Stream vs. Batch
  - Lambda vs. Kappa Architecture
- **System Survey:**
  - Storm/Trident
  - Samza
  - Spark Streaming
  - Flink
- **Discussion:**
  - Comparison Matrix
  - Other Systems

OVERVIEW

Scalable Data
Processing

# A Data Processing Pipeline

We are here!

**Persistence/ Streaming**

**Processing**

**Serving**

**Application**

# Data Processing Frameworks
## 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 Frameworks

## What are your options?

Spark Streaming

Google Dataflow

HERON

Spark

STORM

IBM InfoSphere Streams

Amazon Elastic MapReduce

STORM Trident

Flink

APEX

kafka streams

samza

hadoop

concord

# Big Data Processing Frameworks
## What are your options?

# Big Data Processing Frameworks
## What are your options?

CONCEPTS

# Batch vs. Stream Processing

# Batch Processing

„Volume"

- **Cost-effective** & Efficient
- **Easy to reason about**: operating on complete data

But:

- <span style="color:red">**High latency**</span>: periodic jobs (e.g. during night times)

Persistence
(e.g. HDFS)

Batch
(e.g. MapReduce)

Serving
(e.g. HBase)

Application

# Stream Processing
„Velocity"

- Low end-to-end latency
- Challenges:
  - **Long-running jobs** - no downtime allowed
  - **Asynchronism** - data may arrive delayed or out-of-order
  - **Incomplete input** - algorithms operate on partial data
  - More: fault-tolerance, state management, guarantees, …

**Streaming**
(e.g. Kafka, Redis)

**Real-Time**
(e.g. Storm)

**Serving**

**Application**

# Lambda Architecture

$$\text{Batch}(D_{old}) + \text{Stream}(D_{\Delta now}) \approx \text{Batch}(D_{all})$$

- **Fast** output (real-time)
- Data retention + reprocessing (batch)
  → **„eventually accurate"** merged views of real-time & batch
  Typical setups: Hadoop + Storm ($\rightarrow$ Summingbird), Spark, Flink
- **High complexity** 2 code bases & 2 deployments



**Streaming**
(e.g. Kafka, Redis)  **Persistence**  **Batch**  **Serving**  **Application**

**Real-Time**

Nathan Marz, *How to beat the CAP theorem* (2011)
http://nathanmarz.com/blog/how-to-beat-the-cap-theorem.html

# Kappa Architecture

Stream($D_{all}$) = Batch($D_{all}$)

- **Simpler** than Lambda Architecture
- **Data retention** for history
- Reasons against Kappa:
  - Existing **legacy batch system**
  - **Special tools** only for a particular batch processor
  - Only **incremental** algorithms

*replay*

**Streaming + retention**
(e.g. Kafka, Kinesis)

**Real-Time**

**Serving**

**Application**

Jay Kreps, *Questioning the Lambda Architecture* (2014)
https://www.oreilly.com/ideas/questioning-the-lambda-architecture

# Wrap-up
## Data Processing

- Processing frameworks abstract from **scaling issues**

**Batch processing**
- easy to reason about
- extremely efficient
- huge input-output latency

**Stream processing**
- quick results
- purely incremental
- potentially complex to handle

- **Lambda Architecture**: batch + stream processing
- **Kappa Architecture**: stream-only processing

SURVEY

# Popular Stream Processing Systems

# Processing Models
Batch vs. Micro-Batch vs. Stream

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



low latency                                      high throughput

# Storm
## „Hadoop of real-time"

**Overview**

- **First** production-ready, well-adopted stream processor
- **Compatible**: native Java API, Thrift, distributed RPC
- **Low-level**: no primitives for joins or aggregations
- **Native stream processor**: latency < 50 ms feasible
- **Big users**: Twitter, Yahoo!, Spotify, Baidu, Alibaba, …

**History**

- **2010**: developed at BackType (acquired by Twitter)
- **2011**: open-sourced
- **2014**: Apache top-level project

# Dataflow

**STORM**

Directed Acyclic Graphs (DAG):

- **Spouts**: pull data into topology
- **Bolts**: do processing, emit data
- Asynchronous
- Lineage can be tracked for each tuple
  → At-least-once has **2x messaging overhead**

# Cluster Architecture
## How Storm Scales



Submit Topology

**Nimbus**

Zookeeper

Supervisor
| Worker | Worker |
| Worker | Worker |

**Storm Slave**

Supervisor
| Worker | Worker |
| Worker | Worker |

**Storm Slave**

# Cluster Architecture
## How Storm Scales



Submit Topology

Scheduling & Monitoring

**Nimbus**

Handles coordination

Zookeeper

Supervisor

| Worker | Worker |
|--------|--------|
| Worker | Worker |

**Storm Slave**

Supervisor

| Worker | Worker |
|--------|--------|
| Worker | Worker |

**Storm Slave**

JVM for each worker (runs spouts and bolts as tasks)

# State Management
## Recover State on Failure

- **In-memory** or **Redis**-backed reliable state
- *Synchronous state communication* on the critical path
  → infeasible for large state

# Back Pressure
Throttling Ingestion on Overload

**1.** too many tuples → **2.** tuples time out and fail

**3.** tuples get replayed

**Approach**: monitoring bolts' inbound buffer
1. Exceeding **high watermark** $\rightarrow$ throttle!
2. Falling below **low watermark** $\rightarrow$ full power!

# Trident
## Stateful Stream Joining on Storm

## Overview:

- Abstraction layer on top of Storm
- Released in 2012 (Storm 0.8.0)
- **Micro-batching**
- **New features**:
  - High-level API: aggregations & joins
  - Strong ordering
  - Stateful exactly-once processing
    → Performance penalty

# Trident
## Partitioned Micro-Batching

# Samza
## Real-Time on Top of Kafka



## **Overview**

- ◦ Co-developed with **Kafka**
  → **Kappa Architecture**
- ◦ **Simple**: only single-step jobs
- ◦ Local state
- ◦ Native stream processor: low latency
- ◦ **Users**: LinkedIn, Uber, Netflix, TripAdvisor, Optimizely, …

## **History**

- ◦ Developed at **LinkedIn**
- ◦ **2013**: open-source (Apache Incubator)
- ◦ **2015**: Apache top-level project

# Dataflow
## Simple By Design



- **Job**: processing step (≈ Storm bolt)
  → Robust
  → But: often several jobs
- **Task**: job instance (parallelism)
- **Message**: single data item
- **Output persisted** in Kafka
  → Easy data sharing
  → Buffering (no back pressure!)
  → But: Increased latency
- **Ordering** within partitions
- Task = Kafka partitions: not-elastic on purpose

Martin Kleppmann, *Turning the database inside-out with Apache Samza* (2015)
https://www.confluent.io/blog/turning-the-database-inside-out-with-apache-samza/ (2017-02-23)

# Samza
Local State

Advantages of local state:

- **Buffering**
  - → No back pressure
  - → At-least-once delivery
  - → Simple recovery
- **Fast lookups**

# Dataflow
## Example: Enriching a Clickstream



**Example**: the *enriched clickstream* is available to every team within the organization

# State Management
## Straightforward Recovery

# Spark
## „MapReduce successor"

## Overview

◦ **High-level API**: immutable collections (RDDs)

| Core | SQL | MLlib | GraphX | Spark Streaming |
|------|-----|-------|--------|-----------------|

◦ **Community**: 1000+ contributors in 2015

◦ **Big users**: Amazon, eBay, Yahoo!, IBM, Baidu, …

## History

◦ **2009**: developed at UC Berkeley

◦ **2010**: open-sourced

◦ **2014**: Apache top-level project

# Spark Streaming

**Overview**

- High-level API: DStreams (~Java 8 Streams)
- Micro-Batching: seconds of latency
- Rich features: stateful, exactly-once, elastic

**History**

- 2011: start of development
- 2013: Spark Streaming becomes part of Spark Core

# Spark Streaming
## Core Abstraction: DStream

**Resilient Distributed Data set** (RDD)

- Immutable collection & **deterministic** operations

- Lineage tracking:
  → state can be reproduced
  → periodic checkpoints reduce recovery time

**DStream:** Discretized RDD

- **RDDs are processed in order**: no ordering within RDD

- RDD scheduling ~50 ms → latency >100ms

# Example
## Counting Page Views

```
pageViews = readStream("http://...", "1s")
ones = pageViews.map(event => (event.url, 1))
counts = ones.runningReduce((a, b) => a + b)
```



Zaharia, Matei, et al. "Discretized streams: Fault-tolerant streaming computation at scale." *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*. ACM, 2013.

# Flink

## Overview

- ◦ **Native stream processor:** Latency <100ms feasible
- ◦ **Abstract API** for stream and batch processing, stateful, exactly-once delivery
- ◦ **Many libraries**: Table and SQL, CEP, Machine Learning , Gelly…
- ◦ **Users**: Alibaba, Ericsson, Otto Group, ResearchGate, Zalando…

## History

- ◦ **2010**: start as **Stratosphere** at TU Berlin, HU Berlin, and HPI Potsdam
- ◦ **2014**: Apache Incubator, project renamed to Flink
- ◦ **2015**: Apache top-level project

# Architecture
## Streaming + Batch



| CEP | Table | Storm API | Apache Beam | SAMOA | | Hadoop M/R | Table | Gelly | ML | Apache Beam | Cascading | Zeppelin |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

| DataStream (Java / Scala) | | DataSet (Java/Scala) |
|---|---|---|

Streaming dataflow runtime

| YARN | Cluster | Local |
|---|---|---|

# Managed State
Streaming + Batch

- Automatic **Backups** of local state
- Stored in **RocksDB,** Savepoints written to **HDFS**

# Highlight: Fault Tolerance
## Distributed Snapshots



data stream

← *newer records*

*older records* →

checkpoint barrier *n*

checkpoint barrier *n-1*

*stream record (event)*

part of checkpoint *n+1*

part of checkpoint *n*

part of checkpoint *n-1*

- **Ordering** within stream partitions
- Periodic **checkpoints**
- Recovery:
    1. *reset state* to checkpoint
    2. *replay data* from there

→ **Exactly-once**

# Side-by-side comparison

# 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** | ✓ | ✓ | ✗ | ✓ | ✗ |

# Performance
## Yahoo! Benchmark

▸ Based on **real use case**:
  ◦ Filter and count ad impressions
  ◦ 10 minute windows

"**Storm** […] and **Flink** […] show **sub-second latencies** at relatively high throughputs with **Storm** having the **lowest 99th percentile** latency. **Spark** streaming […] supports high **throughputs**, but at a relatively **higher latency**."

From https://yahooeng.tumblr.com/post/135321837876/benchmarking-streaming-computation-engines-at

# Other Systems

**Heron**

**Apex**

**Dataflow**

**Beam**

**Kafka Streams**

**IBM InfoSphere Streams**

**And even more**: Kinesis, Gearpump, MillWheel, Muppet, S4, Photon, ...

# Summary

▸ Stream Processors:



STORM     Flink     samza     Spark Streaming

latency                                          throughput

▸ **Many Dimensions of Interest**: consistency guarantees, state management, backpressure, ordering, elasticity, …

# Outline

**Introduction**
Where From? Where To?

**Stream Processing**
Big Data + Low Latency

**Real-Time Databases**
Push-Based Collections

**Future Directions**
Current Research & Outlook

- **Big Picture:**
  - **Why** Push-Based Database Queries?
  - **Where** Do Real-Time Databases Fit in?
- **System Survey:**
  - Meteor
  - RethinkDB
  - Parse
  - Firebase
- **Discussion:**
  - Comparison Matrix
  - Other Systems

REAL-TIME DBS

# Making Databases Push-Based

# Traditional Database Access
## No Request? No Data!



**Query maintenance:** periodic polling
→ Inefficient
→ Slow

# Quick Comparison

DBMS vs. RT DB vs. DSMS vs. Stream Processing



**Database Management**

**Real-Time Databases**

**Data Stream Management**

**Stream Processing**

static collections

evolving collections

persistent/ ephemeral streams

ephemeral streams

pull-based

push-based

REAL-TIME DBS

# System Survey

# Meteor

**Overview:**

- JavaScript Framework for interactive apps and websites
  - **MongoDB** under the hood
  - **Real-time** result updates, full MongoDB expressiveness
- **Open-source**: MIT license
- **Managed service**: Galaxy (Platform-as-a-Service)

**History:**

- 2011: *Skybreak* is announced
- 2012: Skybreak is renamed to Meteor
- 2015: Managed hosting service Galaxy is announced

# Live Queries
## Poll-and-Diff

- **Change monitoring**: app servers detect relevant changes
  → *incomplete* in multi-server deployment
- **Poll-and-diff**: queries are re-executed periodically
  → **staleness window**
  → **does not scale** with queries



repeat query every 10 seconds

forward
CRUD

monitor
incoming
writes

CRUD

# Oplog Tailing
## Basics: MongoDB Replication

METE☄R

- **Oplog**: rolling record of data modifications
- **Master-slave replication**:
  Secondaries subscribe to oplog

write operation

mongoDB cluster
(3 shards)

Primary A   Primary B   Primary C

apply

propagate change

Secondary C1   Secondary C2   Secondary C3

# Oplog Tailing
## Tapping into the Oplog

- *Every* Meteor server receives *all* DB writes through oplogs

mongoDB cluster (3 shards)

Primary A   Primary B   Primary C

query
(when in doubt)

Oplog broadcast

monitor
oplog

METEOR
App server

METEOR
App server

push relevant events

CRUD

# Oplog Tailing
## Oplog Info is Incomplete

**METE★R**

## What game does Bobby play?
→ if baccarat, he takes first place!
→ if something else, nothing changes!

*Partial* update from oplog:
`{ name: „Bobby", score: 500 } // game: ???`

Baccarat players sorted by high-score

**METE★R**
```
1. { name: „Joy", game: „baccarat", score: 100 }
2. { name: „Tim", game: „baccarat", score: 90 }
3. { name: „Lee", game: „baccarat", score: 80 }
```

# Oplog Tailing
## Tapping into the Oplog



- *Every* Meteor server receives *all* DB writes through oplogs
→ does not scale

# RethinkDB

## Overview:

- „**MongoDB done right**": comparable queries and data model, but also:
  - **Push-based queries** (filters only)
  - **Joins** (non-streaming)
  - **Strong consistency**: linearizability
- **JavaScript SDK** (*Horizon*): open-source, as managed service
- **Open-source**: Apache 2.0 license

## History:

- 2009: RethinkDB is founded
- 2012: RethinkDB is open-sourced under AGPL
- 2016, May: first official release of Horizon (JavaScript SDK)
- 2016, October: RethinkDB announces shutdown
- 2017: RethinkDB is relicensed under Apache 2.0

# RethinkDB
## Changefeed Architecture

- Range-sharded data
- **RethinkDB proxy**: support node without data
  - Client communication
  - Request routing
  - Real-time query matching

- *Every* proxy receives *all* database writes
  → does not scale

RethinkDB storage cluster

RethinkDB proxy

App server

RethinkDB proxy

App server

*Bottleneck!*

William Stein, *RethinkDB versus PostgreSQL: my personal experience* (2017)
http://blog.sagemath.com/2017/02/09/rethinkdb-vs-postgres.html (2017-02-27)

Daniel Mewes, *Comment on GitHub issue #962: Consider adding more docs on RethinkDB Proxy* (2016)
https://github.com/rethinkdb/docs/issues/962 (2017-02-27)

# Parse

## Overview:

- **Backend-as-a-Service** for mobile apps
  - **MongoDB:** largest deployment world-wide
  - **Easy development**: great docs, push notifications, authentication, ...
  - **Real-time** updates for most MongoDB queries
- **Open-source**: BSD license
- **Managed service**: discontinued

## History:

- 2011: Parse is founded
- 2013: Parse is acquired by Facebook
- 2015: more than 500,000 mobile apps reported on Parse
- 2016, January: Parse shutdown is announced
- 2016, March: **Live Queries** are announced
- 2017: Parse shutdown is finalized

# Parse
## LiveQuery Architecture

- **LiveQuery Server**: no data, real-time query matching
- *Every* LiveQuery Server receives
  *all* database writes
  → does not scale



*Bottleneck!*

# Firebase



**Overview:**

- **Real-time state synchronization** across devices
- **Simplistic data model:** nested hierarchy of lists and objects
- **Simplistic queries**: mostly navigation/filtering
- **Fully managed**, proprietary
- **App SDK** for App development, mobile-first
- **Google services integration**: analytics, hosting, authorization, ...

**History:**

- 2011: chat service startup Envolve is founded
  → was often used for cross-device state synchronization
  → state synchronization is separated (Firebase)
- 2012: Firebase is founded
- 2013: Firebase is acquired by Google
- 2017, October: Firestore is released

# Firebase
## Real-Time State Synchronization

- **Tree data model**: application state ~JSON object
- **Subtree synching**: push notifications for specific keys only
  → Flat structure for fine granularity

→ *Limited expressiveness!*

# Firebase
## Query Processing in the Client

- Push notifications for **specific keys** only
  - Order by a **single attribute**
  - Apply a **single filter** on that attribute

- Non-trivial query processing in client
  → does not scale!

Jacob Wenger, on the Firebase Google Group *(2015)*
https://groups.google.com/forum/#!topic/firebase-talk/d-XjaBVL2Ko (2017-02-27)

Illustration taken from: Frank van Puffelen, *Have you met the Realtime Database? (2016)*
https://firebase.googleblog.com/2016/07/have-you-met-realtime-database.html (2017-02-27)

# Firebase
## Hard Scaling Limits

"Scale to around **100,000 concurrent connections** and **1,000 writes/second** in a single database. Scaling beyond that requires sharding your data across multiple databases."

# Firebase
## Firestore: New Model



documents

collections

references

# Firebase
## Firestore: New Model

finer access granulates

tree-like structure

# Firebase
## Firestore: Summary

- More specific data selection
- Logical AND for some filter combinations

... But:

- Still **Limited Expressiveness**
  - No logical OR
  - No logical AND for many filter combinations
  - No content-based search (regex, full-text search)
- Still **Limited Write Throughput**:
  - *500 writes/s* per collection
  - *1 writes/s* per document

Firebase, *Firestore: Quotas and Limits (2018)*
https://firebase.google.com/docs/firestore/quotas (2018-03-10)

# Honorable Mentions

Other Systems With Real-Time Features

# Summary & Discussion

# Wrap-Up
## Direct Comparison

| | Meteor | | RethinkDB | Parse | Firebase | Baqend |
|---|---|---|---|---|---|---|
| | Poll-and-Diff | Oplog Tailing | | | | |
| **Scales with write TP** | ✔ | ✘ | ✘ | ✘ | ✘ | ✔ |
| **Scales with no. of queries** | ✘ | ✔ | ✔ | ✔ | ? (100k connections) | ✔ |
| Composite queries (AND/OR) | ✔ | ✔ | ✔ | ✔ | ◯ (AND In Firestore) | ✔ |
| Sorted queries | ✔ | ✔ | ✔ | ✘ | ◯ (single attribute) | ✔ |
| Limit | ✔ | ✔ | ✔ | ✘ | ✔ | ✔ |
| Offset | ✔ | ✔ | ✘ | ✘ | ◯ (value-based) | ✔ |

# Summary
## Real-Time Databases: Major challenges

**Scalability**:
- ▶ Handle increasing throughput
- ▶ Handle additional queries

**Expressiveness**:
- ▶ Content-based search? Composite filters?
- ▶ Ordering? Limit? Offset?

**Legacy Support**:
- ▶ Real-time queries for *existing databases*?
- ▶ *Decouple* OLTP from real-time workloads?

# Outline

**Introduction**
Where From? Where To?

**Stream Processing**
Big Data + Low Latency

**Real-Time Databases**
Push-Based Collections

**Future Directions**
Current Research & Outlook

- **Caching Dynamic Data:**
  - Why is the Web Slow?
  - Caching to the Rescue!
  - Query Caching
- **Real-Time Queries:**
  - Scalability
  - Expressiveness
  - Legacy Compatibility
  - Use Cases
- **Open Challenges:**
  - TTLs & Transactions
  - Polyglot Persistence
- **Summary**

# Our 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



`0 1 1 0 1 0 0 1`

# Consistent Web Caching

The Cache Sketch



| 0 | 2 | 1 | 4 | 0 |

# Consistent Web Caching

The Cache Sketch



Browser Cache

CDN

| 0 | 2 | 1 | 4 | 0 |

# Consistent Web Caching

The Cache Sketch



Browser Cache

CDN

| 0 | 2 | 1 | 4 | 0 |
|---|---|---|---|---|

# Consistent Web Caching
The Cache Sketch



**purge(obj)**

Browser Cache

CDN

**hashA(oid)**          **hashB(oid)**

| 0 | 3 | 1 | 4 | 1 |

# Consistent Web Caching
## The Cache Sketch



Browser Cache

CDN

**Flat(Counting Bloomfilter)**

| 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|

← 

| 0 | 3 | 1 | 4 | 1 |
|---|---|---|---|---|

# Consistent Web Caching

The Cache Sketch



hashA(oid)    hashB(oid)

| 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|

Browser Cache

CDN

| 0 | 3 | 1 | 4 | 1 |
|---|---|---|---|---|

# Consistent Web Caching

The Cache Sketch



**hashA(oid)**    **hashB(oid)**

| 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|

| 0 | 3 | 1 | 4 | 1 |
|---|---|---|---|---|

# Consistent Web Caching

The Cache Sketch



Browser Cache

CDN

| 0 | 1 | 1 | 1 | 1 |

| 0 | 3 | 1 | 4 | 1 |

# Consistent Web Caching

The Cache Sketch



Browser Cache

CDN

**hashA(oid)**        **hashB(oid)**

| 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|

| 0 | 2 | 1 | 4 | 0 |
|---|---|---|---|---|

# Consistent Web Caching
The Cache Sketch

With 20.000 distinct updates and 5% error rate: **11 KByte**

hashA(oid)          hashB(oid)

| 0 | 1 | 1 | 1 | 1 |
|---|---|---|---|---|

| 0 | 2 | 1 | 4 | 0 |
|---|---|---|---|---|

# How to <u>Invali</u>date <u>DB</u> Query Results?

# InvaliDB
## Invalidating DB Queries



How to **detect changes to query results**:
*"Give me the most popular products that are in stock."*

Add

Change

Remove

# InvaliDB
## Invalidating DB Queries

Real-Time
Queries
(*Websockets*)

Create
Update
Delete

Fresh Caches

Server

Pub-Sub ▦
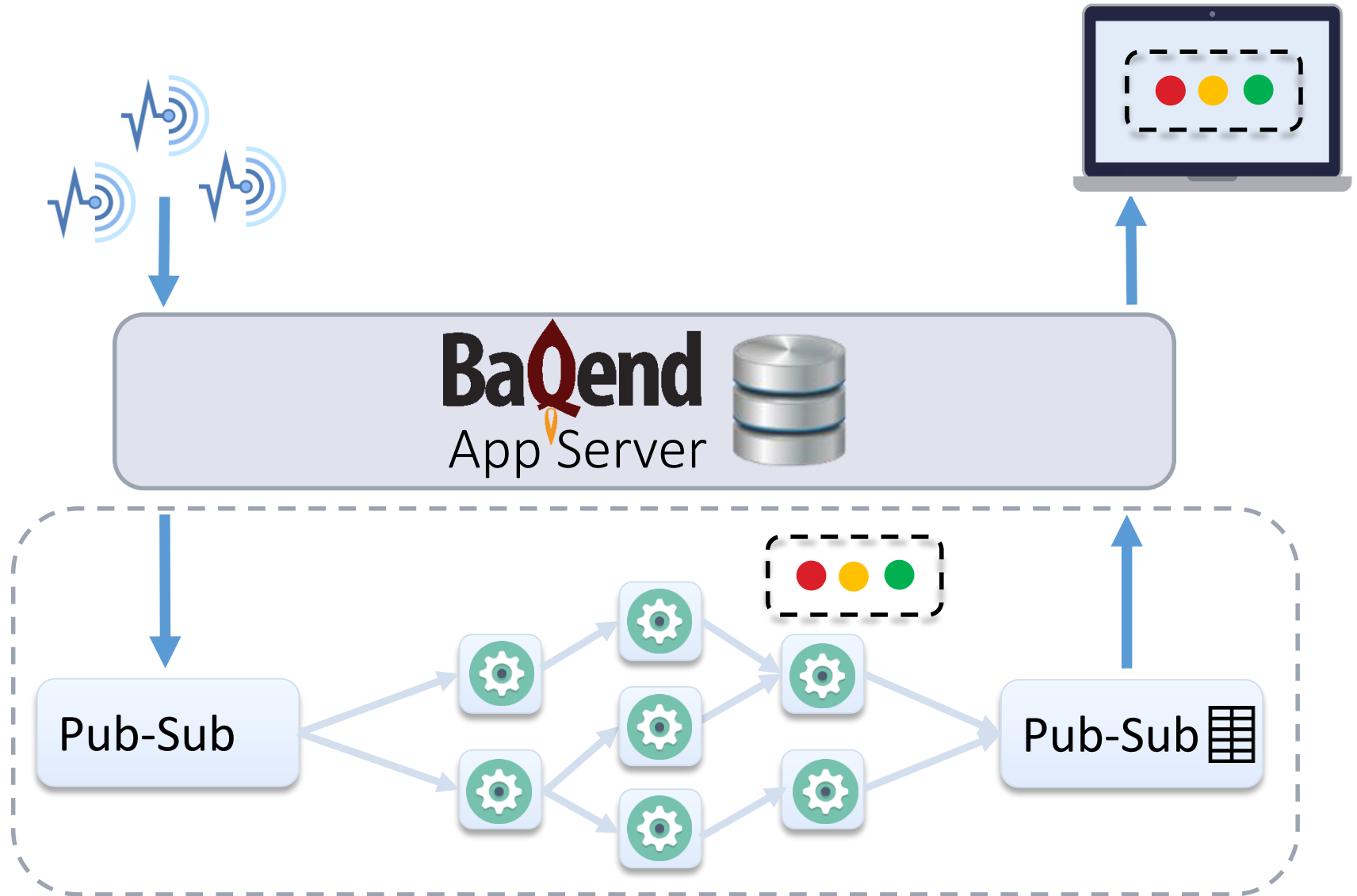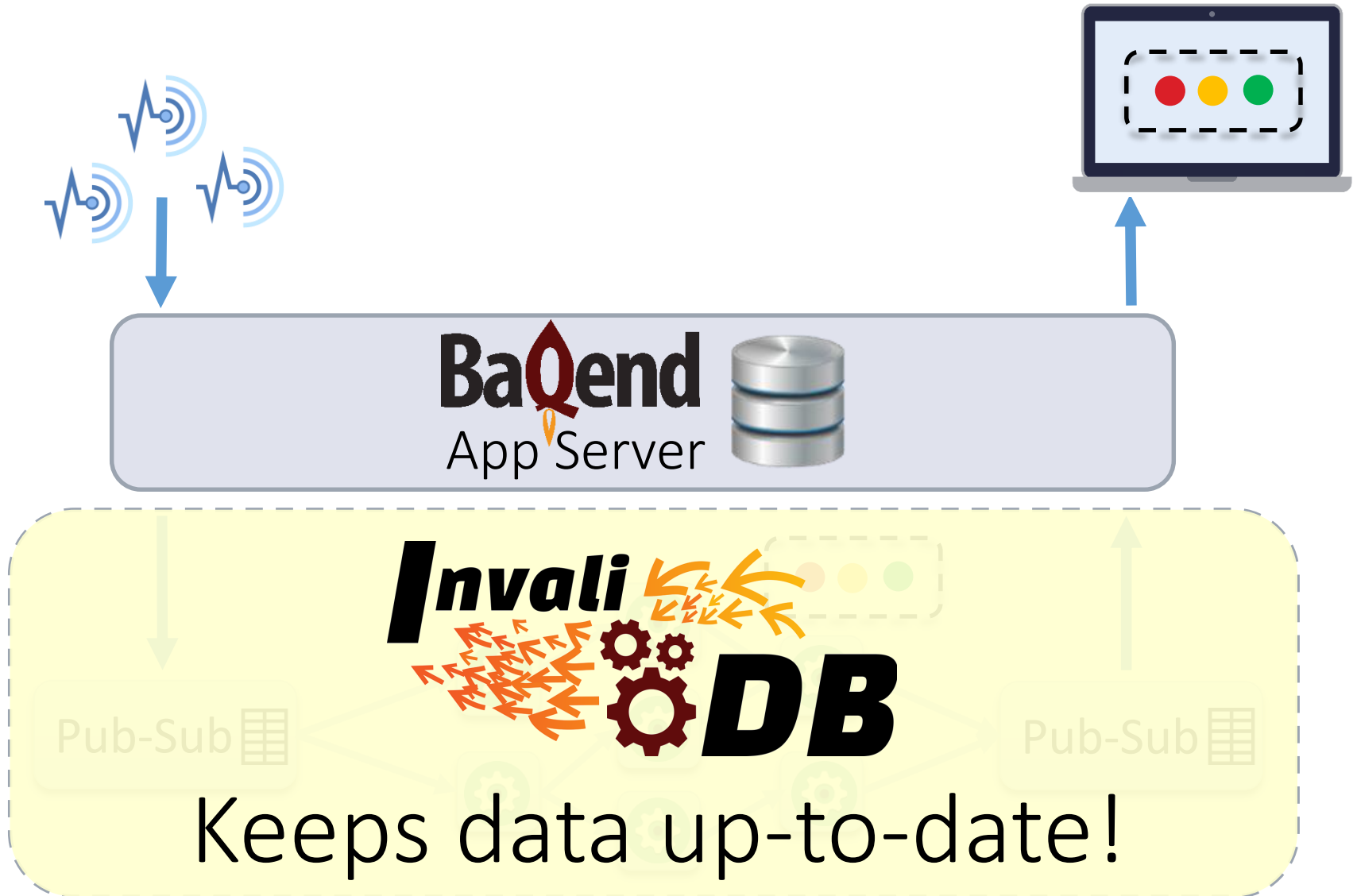
Pub-Sub ▦

*Invali* **DB**

# Baqend Real-Time Queries
## Realtime Decoupled

# Baqend Real-Time Queries
## Realtime Decoupled

# Baqend Real-Time Queries
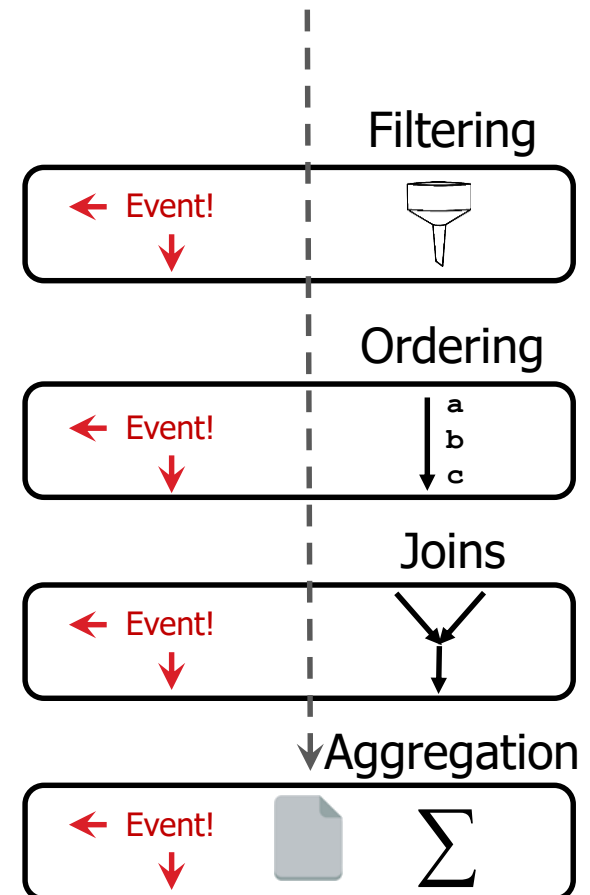## Realtime Decoupled



Keeps data up-to-date!

# Baqend Real-Time Queries
## Staged Real-Time Query Processing

Change notifications go through different query processing stages:

1. **Filter queries**: track matching status
   $\rightarrow$ *before-* and after-images
2. **Sorted queries**: maintain result order
3. **Joins**: combine maintained results
4. **Aggregations**: maintain aggregations

# Baqend Real-Time Queries
## Filter Queries: Distributed Query Matching

**Two-dimensional partitioning**:

- *by Query*
- *by Object*
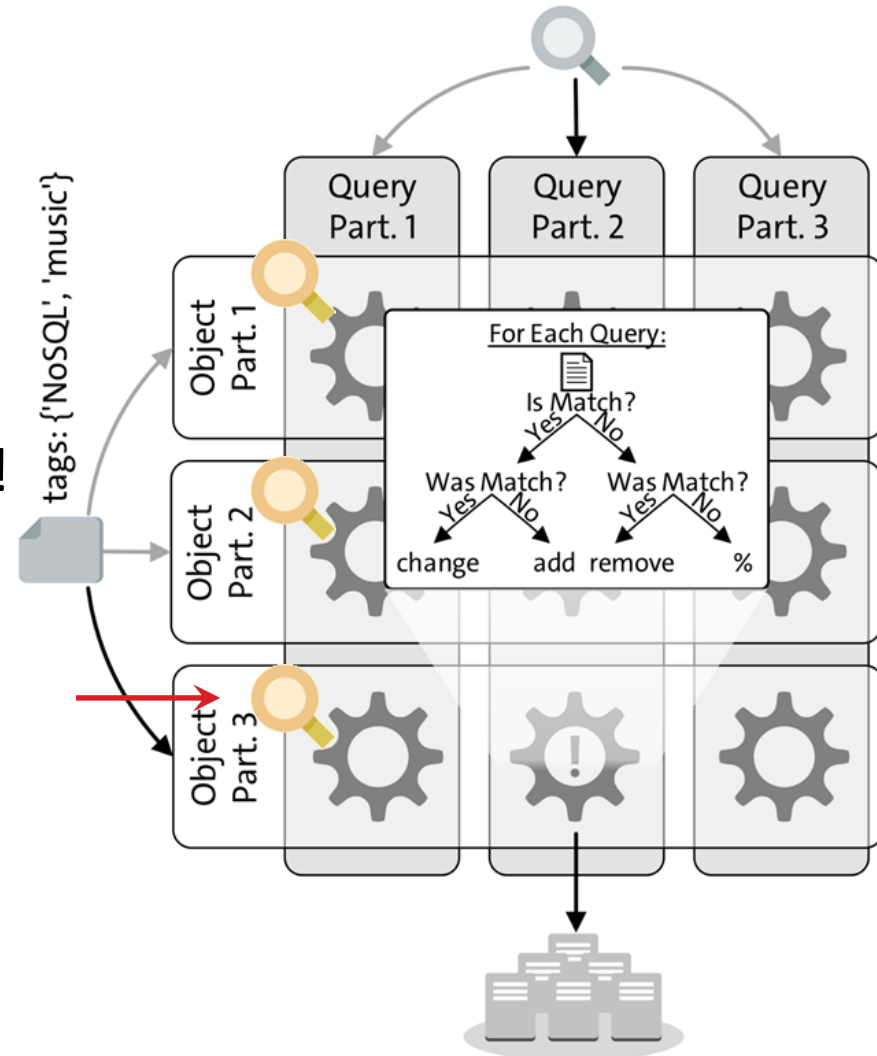
→ scales with queries and writes

Implementation:

- Apache Storm
- Topology in Java
- MongoDB query language
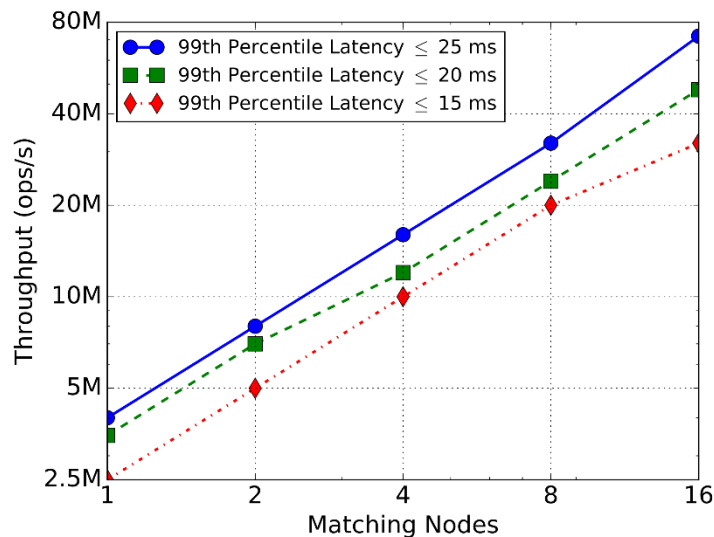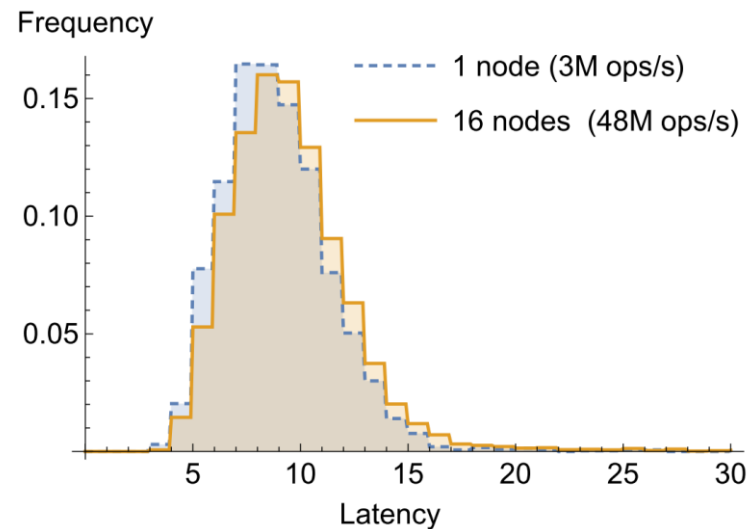- **Pluggable query engine**

→ legacy-compatible

SELECT * FROM posts WHERE tags CONTAINS 'NoSQL'

Query Part. 1

Query Part. 2

Query Part. 3

Object Part. 1

Object Part. 2

Object Part. 3

tags: {'NoSQL', 'music'}

Write op!

For Each Query:

Is Match?
Yes / No

Was Match?
Yes / No

Was Match?
Yes / No

change    add    remove    %

# Baqend Real-Time Queries
## Low Latency + Linear Scalability



Linear Scalability

Stable Latency Distribution

*Quaestor: Query Web Caching for Database-as-a-Service Providers*
VLDB '17

# Programming Real-Time Queries
## JavaScript API

```javascript
var query = DB.Tweet.find()
            .matches('text', /my filter/)
            .descending('createdAt')
            .offset(20)
            .limit(10);
```

**Static Query**

```javascript
query.resultList(result => ...);
```

Google

**Real-Time Query**

```javascript
query.resultStream(result => ...);
```

Twoogle

# Twoogle

Filter word, e.g. "http", "Java", "Baqend" 🔍

**Real-Time**  Static

Last result update at 15:51:21 (less than a second ago)

1. Conju.re (conju_re, 3840 followers) tweeted:
https://twitter.com/conju_re/status/859767327570702336

Congress Saved the Science Budget—And That's the Problem https://t.co/UdrjNidakc https://t.co/xlNjpEpKZG

2. ねぼすけゆーだい (Yuuu___key, 229 followers) tweeted:
https://twitter.com/Yuuu___key/status/859767323384623104

けいきさんと PENGUIN　RESEARCHのけいたくん がリプのやり取りしてる...

3. Whitney Shackley (bschneids11, 5 followers) tweeted:
https://twitter.com/bschneids11/status/859767319534469122

holy...... waiting for it so long🍫 ☺ https://t.co/UdXcHJb7X3

4. Lisa Schmid (LisaMSchmid, 67 followers) tweeted on #teamscs, and #scs...
https://twitter.com/LisaMSchmid/status/859767317311500290

Congrats to Matthew Kent, winner of the 26th #TeamSCS Coding Challenge. https://t.co/vx1o0WgJrZ #SCSchallenge

5. Brian Martin Larson (Brian_Larson, 40 followers) tweeted on #teamscs, a...
https://twitter.com/Brian_Larson/status/859767317303001089

Congrats to Matthew Kent, winner of the 26th #TeamSCS Coding Challenge.

---

# 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 Legacy Websites
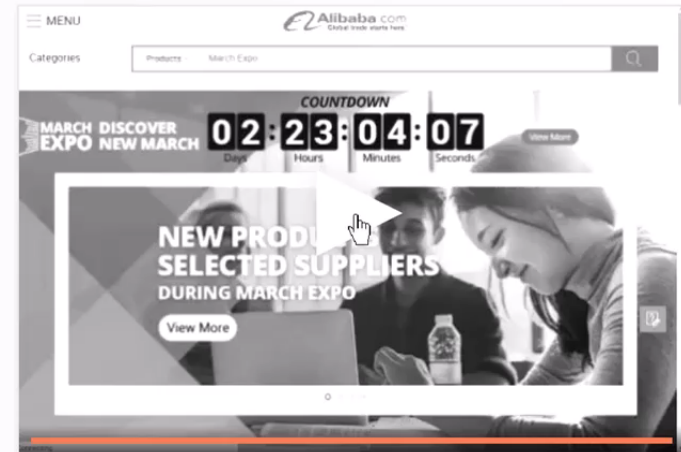
Website

3rd Party
**Services**

Existing
**Backend**

# Speed Kit
## Baqend Caching for Legacy Websites

Website with
**Snippet**

Speed Kit
**Service Worker**

**Baqend**
Service

*Requests*

*Fast Requests*

*other*

Pull

Push

3rd Party
**Services**

Existing
**Backend**

# Open Challenges

# TTL Estimation
Quantifying Cacheability of Dynamic Content

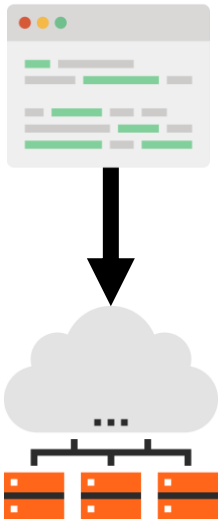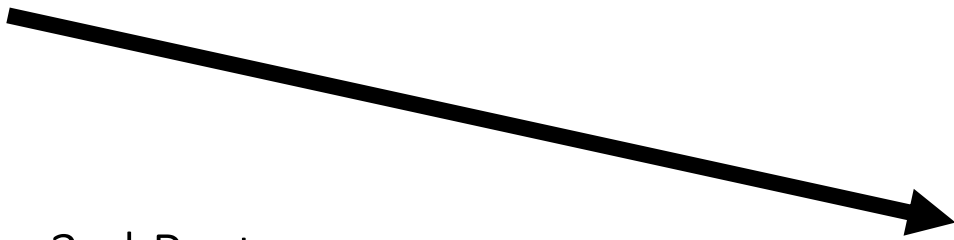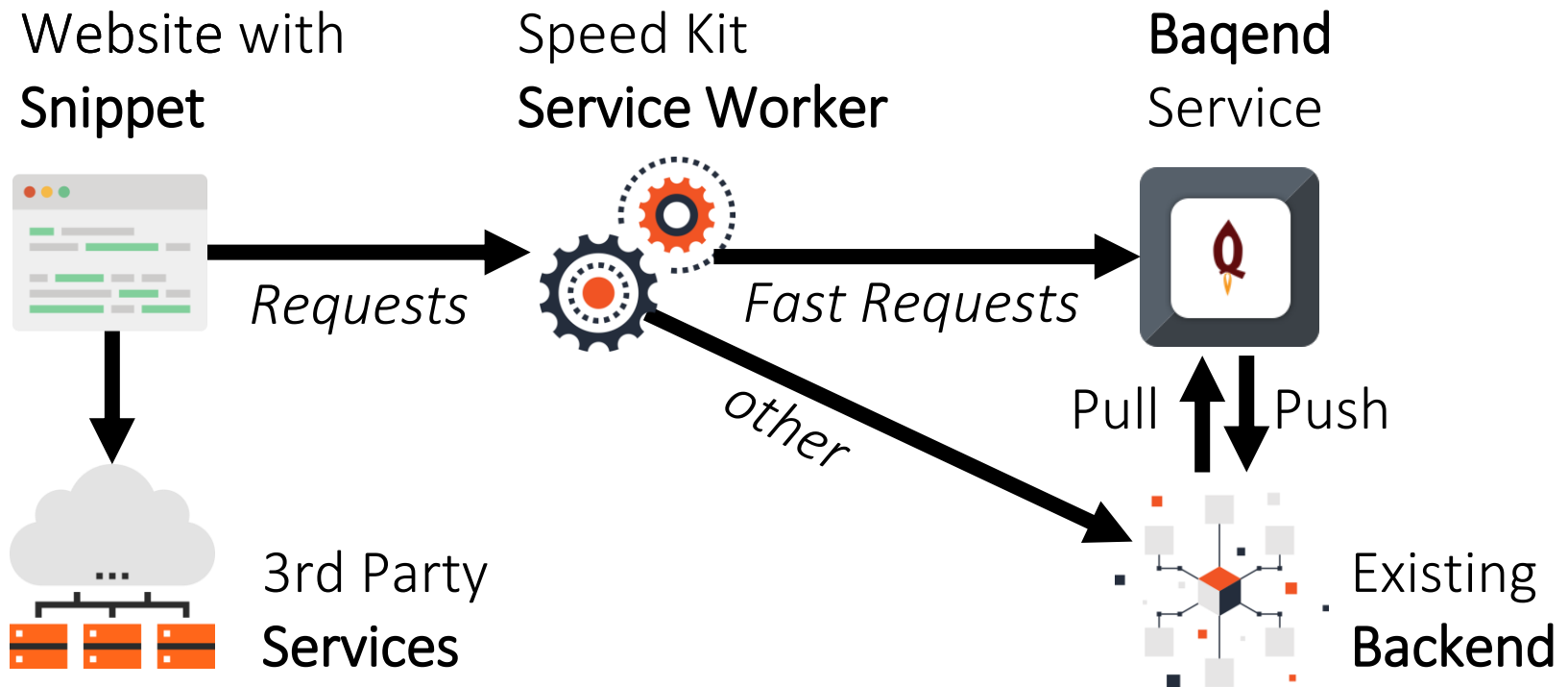▸ **Setting:** server assigns a caching time-to-live (TTL) to each record and query result

▸ **Problem**:

TTLs too short: Bad cache-hit rate

TTLs too large: Bloom filter's false positive rate degrades

▸ **Approach:** Collect access metrics and estimate

**Objects:** calculate the expected value of the time to next write (assuming a poisson process)

**Queries**:

- **Initial estimate**: estimated time until first object in result is updated
- **Refinement**: upon invalidation TTL is adapted towards observed TTL using an EWMA

# TTL Estimation
Learning Representations

**Setting:** query results can either be represented as references (id-list) or full results (object-lists)

Id-Lists

$$\{id_1, id_2, id_3\}$$

Less Invalidations

Object-Lists

$$\{\{id: 1, val: 'a'\}, \{id: 2, val: 'b'\},$$
$$\{id: 3, val: 'c'\}\}$$

Less Round-Trips

**Current Approach**: Cost-based decision model that weighs expected round-trips vs expected invalidations

**Desired:** Adaptive agent that actively explores decisions

# TTL Estimation
## Open Challenge: Learning Workloads

(a) Growing Pattern

(b) On/Off Pattern

(c) Bursty Pattern

(d) Random Pattern

**◀** „**Backwards-oriented**" (*current approach*):

- Mesure & use **moving average** or **newest measurement**
- Cannot react to spikes/fluctuation nor detect patterns

**▶** „**Predictive online-learning**":

- Extrapolate using **regression** (e.g. linear or polynomial) or **time-series models** (Exponential Smoothing, AR, ARIMA, Gaussian Processes, …)
- Resource intensive, very difficult to select & evalute model

**↻** „**Reactive**":

- Use **Reinforcement learning** to automatically explore decisions
- Rewards usually noisy, delayed or hidden (e.g. staleness)

# Polyglot Persistence Mediator
Schemas can be annotated with requirements/SLAs

- Write Throughput > 10,000 RPS
- Read Availability > 99.9999%
- Scans = true
- Full-Text-Search = true
- Monotonic Read = true

Schema

DBs
Tables
Fields

# Polyglot Persistence Mediator
Routing to the „optimal" datbase system

Application

Recursive Ranking Algorithm
for $schemaElemt \rightarrow DB$ mapping

**Data and
Operations**

Database
Metrics

**Polyglot Persistence
Mediator**

Routing
Model

Annotated
Schema

Latency < 30ms

db$_1$            db$_2$            db$_3$

# Polyglot Persistence
## Open Challenges

**Meta-DBaaS**: Mediate over DBaaS-systems unify SLAs

**Live Migration**: adapt to changing requirements

**Database Selection**: Actively minimize SLA violations

**Utility Functions/SLAs**: Capture trade-offs comprehensively

**Workload Management**: Adaptive Runtime Scheduling

# Distributed Transactions

**Transaction Abort Rates:** How to mitigate high abort rates caused by long running transactions?

**Automatic Transaction Protocol Selection**: Can the optimal protocol (2PL, BOCC+, RAMP, ...) be learned and chosen at runtime?

**Transactional Visibility For Real-Time Queries**: How to include transactions without introducing bottlenecks?

# Summary

# Summary
## Real-Time Data Management

pull-based                                                                    push-based

| | **Database Management** | **Real-Time Databases** | **Data Stream Management** | **Stream Processing** |
|---|---|---|---|---|
| **Data** | persistent collections | | persistent/ephemeral streams | |
| **Processing** | one-time | one-time + continuous | continuous | |
| **Access** | random | random + sequential | sequential | |
| **Schema** | structured | | | structured, unstructured |

# Summary
## Real-Time Data Management

| **Database Management** | **Real-Time Databases** | **Data Stream Management** | **Stream Processing** |
|---|---|---|---|
| static collections | evolving collections | structured streams | unstructured streams |

pull-based ← → push-based

## NoSQL Databases: a Survey and Decision Guidance

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.

## NoSQL Databases:
### A Survey and Decision Guidance

### TL;DR

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.

## Scalable Stream Processing: A Survey of Storm, Samza, Spark and Flink

## Scalable Stream Processing:
### A Survey of Storm, Samza, Spark and Flink

With this article, we would like to share our insights on real-time data processing we gained building Baqend. This is an updated version of our most recent stream processor survey which is another cooperation with the University of Hamburg (authors: Wolfram Wingerath, Felix Gessert, Steffen Friedrich and Norbert Ritter). As you may or may not have been aware of, a lot of stream processing is going on behind the curtains at Baqend. In our quest to provide the lowest-possible latency, we have built a system to enable **query caching** and **real-time notifications** (similar to *changefeeds* in RethinkDB/Horizon) and hence learned a lot about the competition in the field of stream processors.

Read them at blog.baqend.com!

# Our Related Publications

## Scientific Papers:

📖 *Quaestor: Query Web Caching for Database-as-a-Service Providers*
VLDB '17

📖 *NoSQL Database Systems: A Survey and Decision Guidance*
SummerSOC '16

📖 *Real-time stream processing for Big Data*
it - Information Technology 58 (2016)

📖 *The Case For Change Notifications in Pull-Based Databases*
BTW '17

## Blog Posts:

📖 *Real-Time Databases Explained: Why Meteor, RethinkDB, Parse and Firebase Don't Scale*
Baqend Tech Blog (2017): https://medium.com/p/822ff87d2f87

**A Real-Time Database Survey: The Architecture of Meteor, RethinkDB, Parse & Firebase**

*Real-time databases make it easy to implement reactive applications, because they keep your critical information up-to-date. But how do they work and how do they scale? In this article, we dissect the real-time query features of Meteor, RethinkDB, Parse and Firebase to uncover scaling limitations inherent to their respective designs. We then go on to discuss and categorize related real-time systems and share our lessons learned in providing real-time queries without any bottlenecks in Baqend.*

**A Real-Time Database Survey:**
The Architecture of Meteor, RethinkDB, Parse & Firebase

Learn more at blog.baqend.com!

# Thank you

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

Blog: blog.baqend.com
Slides: slides.baqend.com

@baqendcom