



Towards Scalable Cloud Data Management

Felix Gessert, Norbert Ritter
gessert/ritter@informatik.uni-hamburg.de

3rd Workshop on Scalable Cloud Data Management

Co-located with the IEEE BigData Conference.
Santa Clara, CA, October 29th 2015.
Starting at **8am in Ballroom C**.

[Workshop Schedule](#)

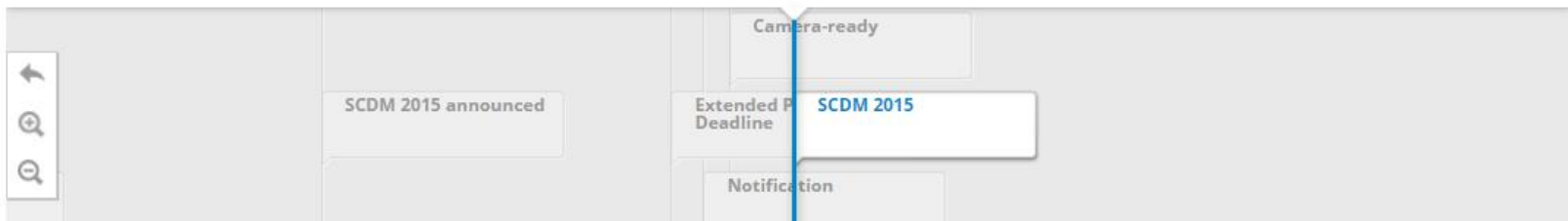
October 29, 2015

SCDM 2015



OCTOBER 5,
2015
Camera-ready

SCDM 2015 Workshop in [Santa Clara, CA](#). The [preliminary schedule](#) is online. SCDM starts on Oct 29, 8am (Ballroom C) with a keynote by Russel Sears on "Purity and the future of scalable storage".





Outline



Motivation



ORESTES: a Cloud-Database Middleware



Solving Latency and Polyglot Storage



Wrap-up

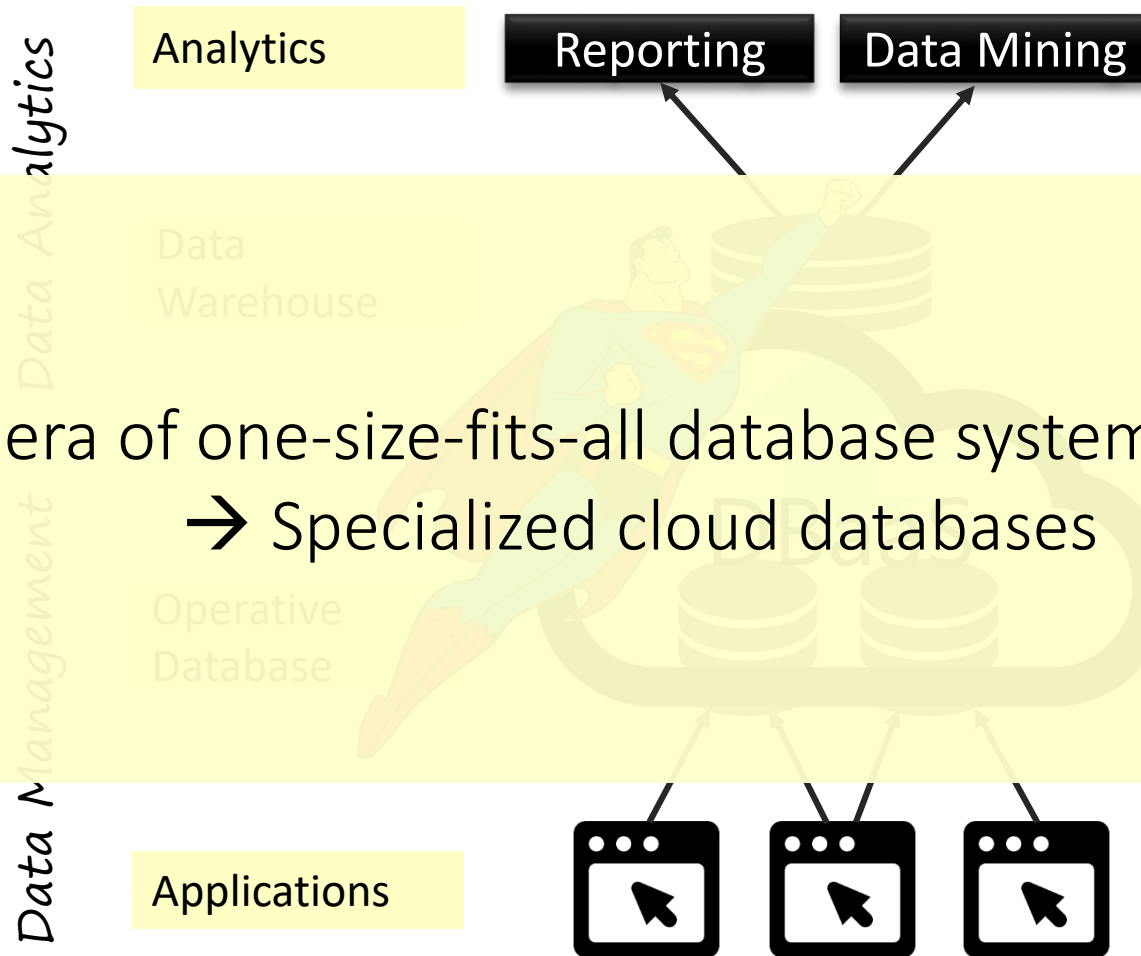
- *Cloud Data Management*
- Cloud Database Models
- Research Challenges



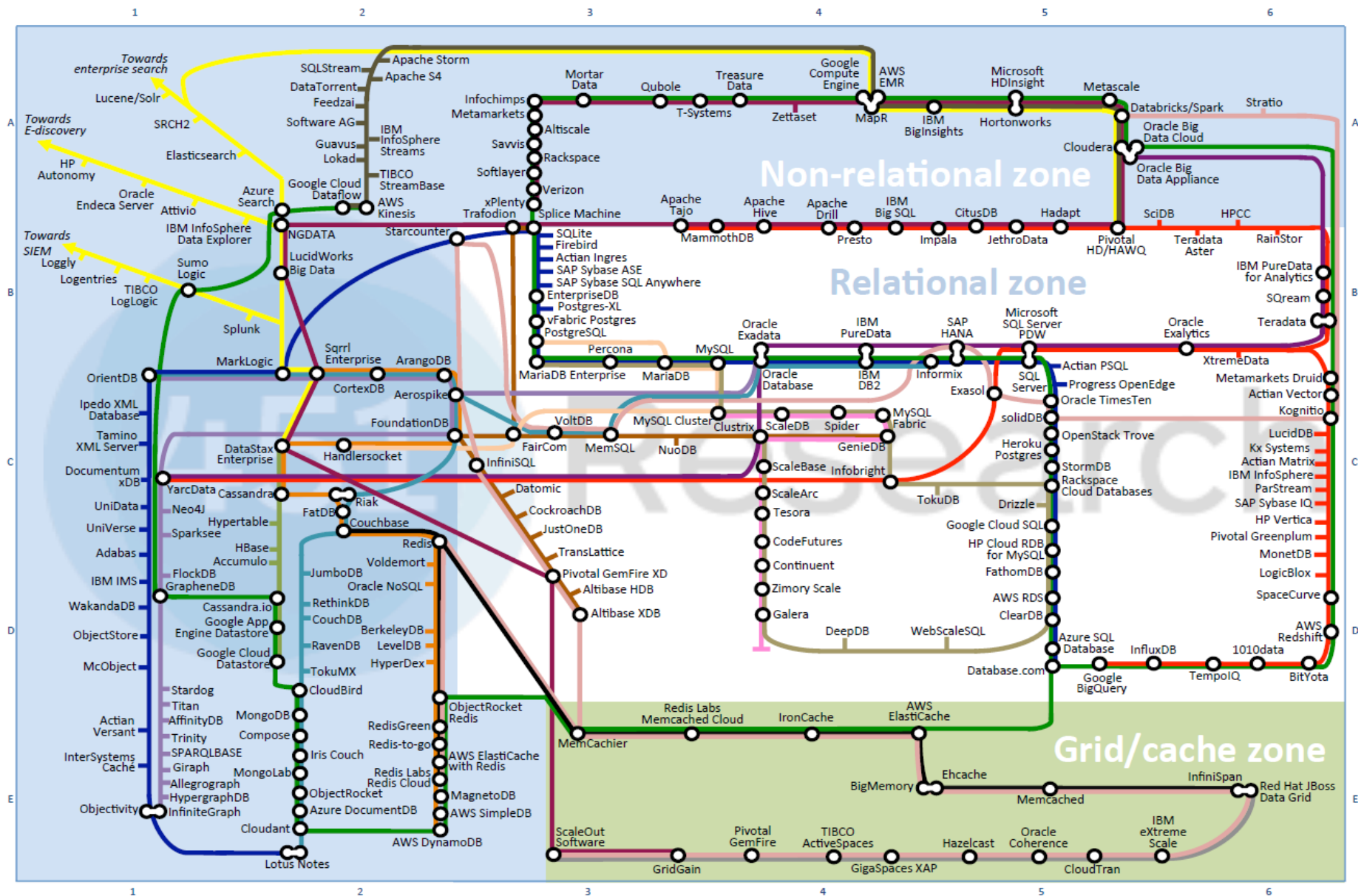
Introduction: What are the challenges in Cloud Data Management?

Architecture

Typical Data Architecture:



The era of one-size-fits-all database systems is over
→ Specialized cloud databases



Database 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

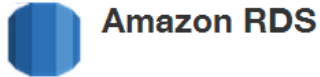
Massive user-
generated content

Cloud-Database Sweetspots



Realtime BaaS

Communication and
collaboration



Managed RDBMS

General-purpose
ACID transactions



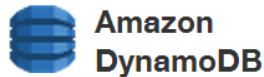
Managed Cache

Caching and
transient storage



Wide-Column Store

Very large tables



Wide-Column Store

Massive user-
generated content



Backend-as-a-Service

Small Websites
and Apps



Managed NoSQL

Full-Text Search



Object Store

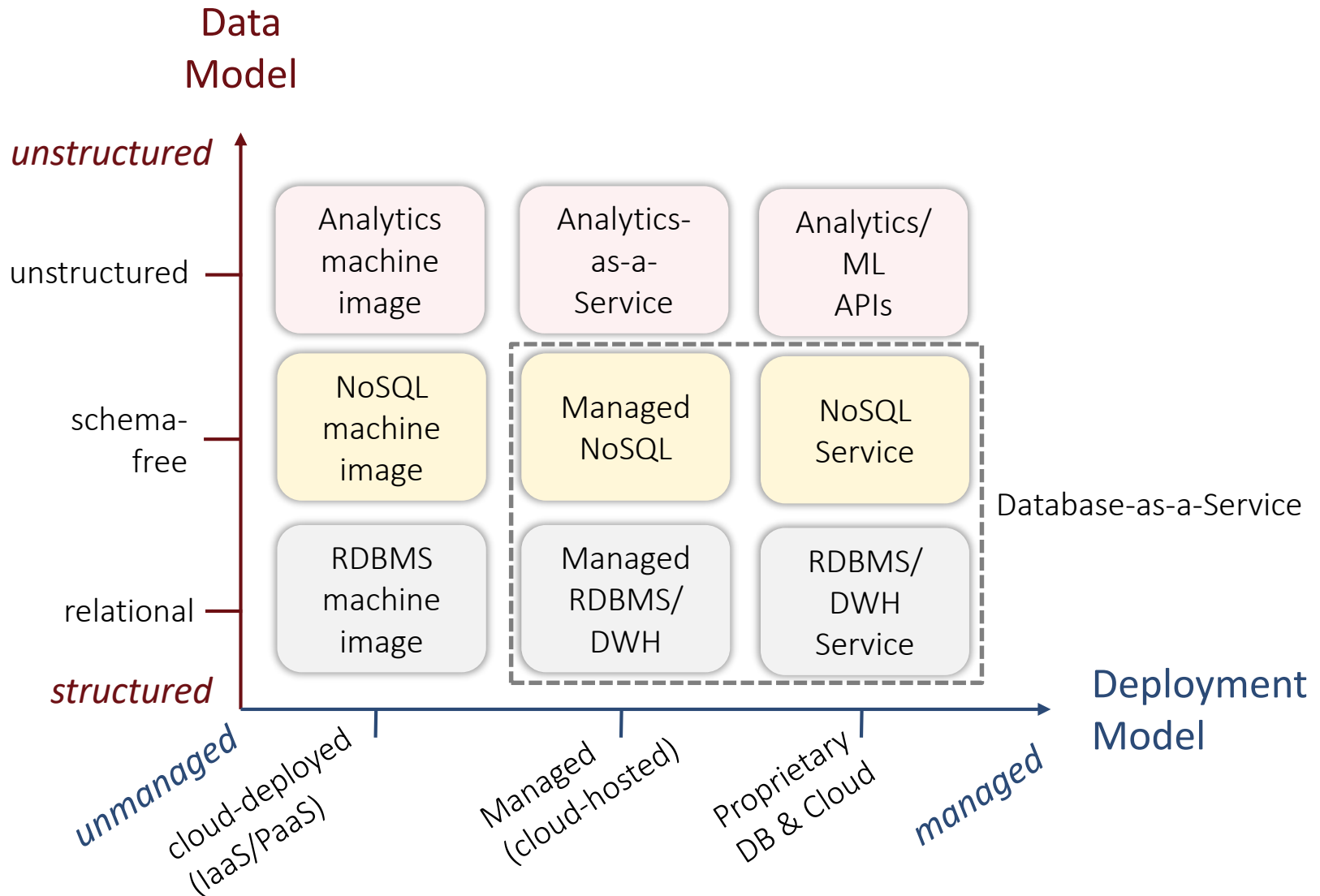
Massive File
Storage



Hadoop-as-a-Service

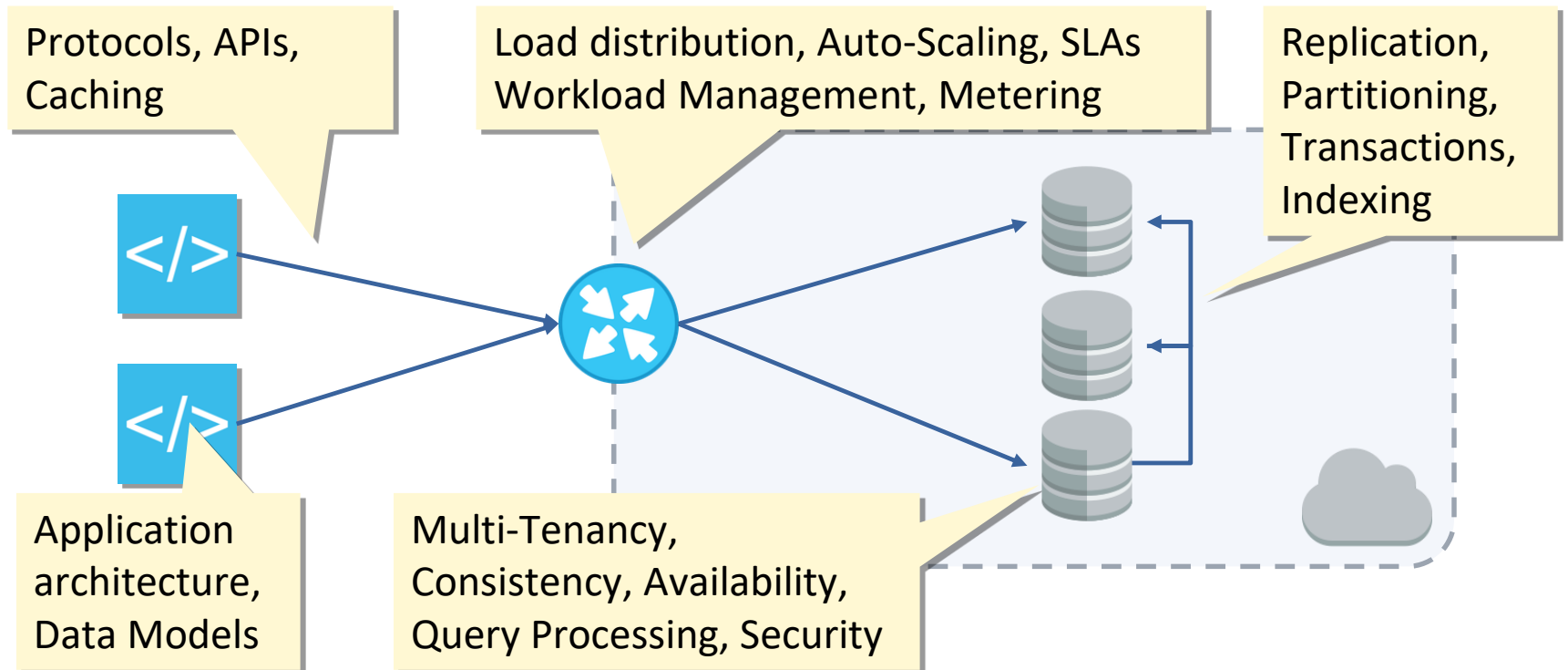
Big Data Analytics

Cloud-Database Models



Cloud Data Management

- ▶ Research field tackling the *design, implementation, evaluation* and *application implications* of **database systems** in cloud environments:

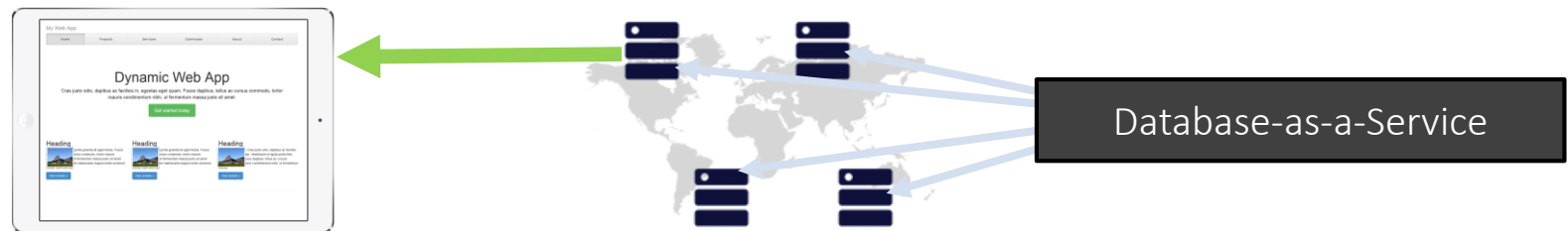


Open Research Questions

Performance & Latency



- ▶ How can database systems support **novel application architectures** (e.g., single-page or real-time apps)?
- ▶ Can the **functionality-performance trade-off** popularized by the NoSQL movement be turned into a tunable runtime configuration?
- ▶ How can a DBaaS deliver **low latency** in face of distributed storage and application tiers?



Open Research Questions

Consistency & Transactionality



- ▶ Which consistency and transaction guarantees can be provided across **(geo-)replicated, partitioned**, possibly heterogeneous/polyglot database systems?
- ▶ How can the **consistency-latency-availability trade-off** be best exposed to applications and developers?
- ▶ Can the existing methods (*quorum-based, consensus-based, master-slave*, etc.) be **reconciliated** into a single approach and tied to application requirements?
- ▶ How can we **replace CAP** by a more fine-grained and nuanced consistency classification scheme?

Open Research Questions

Service-Level Agreements



- ▶ How can database SLAs be guaranteed in a virtualized, multi-tenant **cloud environment**?
- ▶ Can we derive Service-Level-Objectives that are easy enough to understand and maintain to be **practical**?

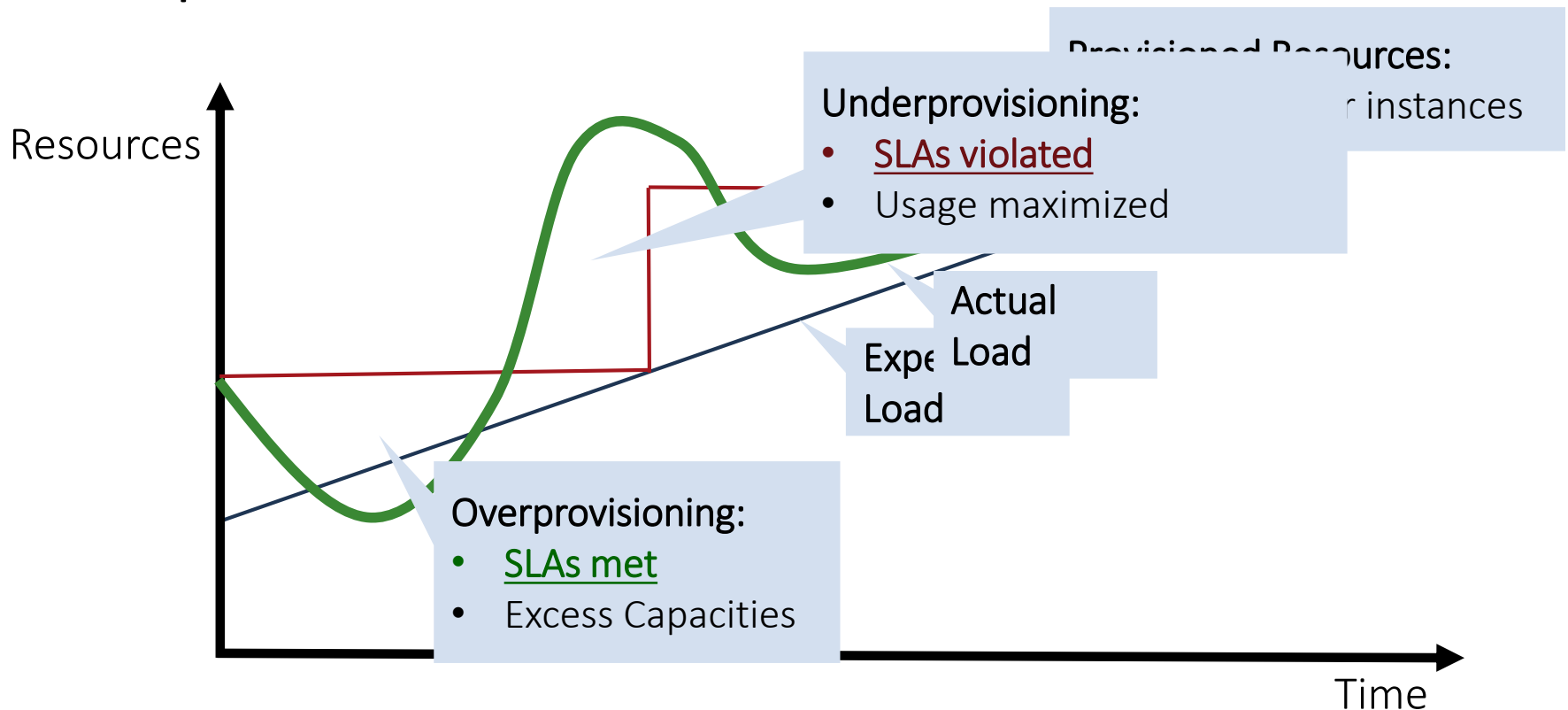
	Model	CAP	SLAs
SimpleDB	Table-Store	CP	✗
DynamoDB	Table-Store	CP	✗
Azure Tables	Table-Store	CP	99.9% uptime
AE/Cloud DataStore	Entity-Group Store	CP	✗
S3, Az. Blob, GCS	Object-Store	AP	99.9% uptime (S3)

Open Research Questions

Service-Level Agreements



- ▶ How can SLAs be incorporated in **autoscaling** to optimize costs and minimize SLA violations?

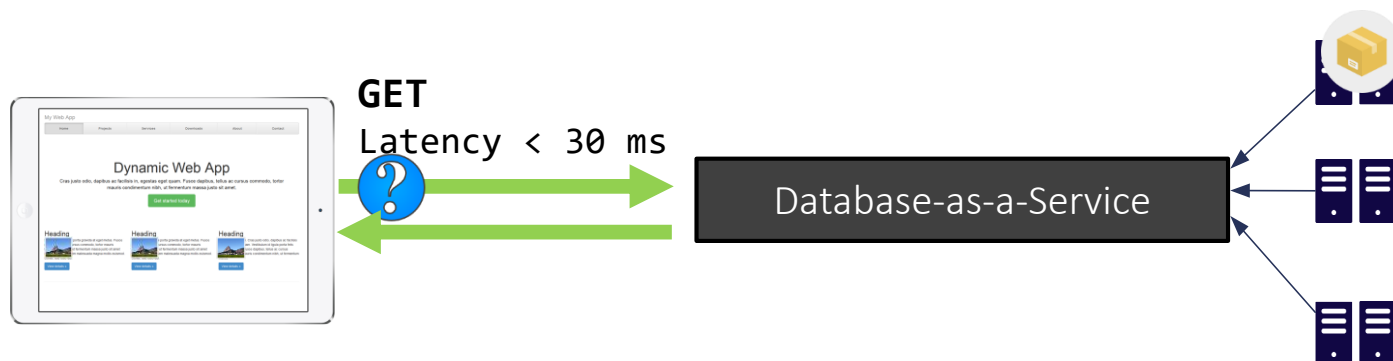


Open Research Questions

Poylgot Persistence



- ▶ Can the data system functions of **storage, search, streaming and analytics** be unified and integrated?
- ▶ Is it possible to **automate**, optimize and learn the **best choice** of given database systems?
- ▶ How can queries and data be **routed to databases**, so that SLAs & performance characteristics are met?



Outline



Motivation



ORESTES: a Cloud-Database Middleware



Solving Latency and Polyglot Storage



Wrap-up

- Two problems:
 - Latency
 - Polyglot Storage
- Vision: Orestes Middleware



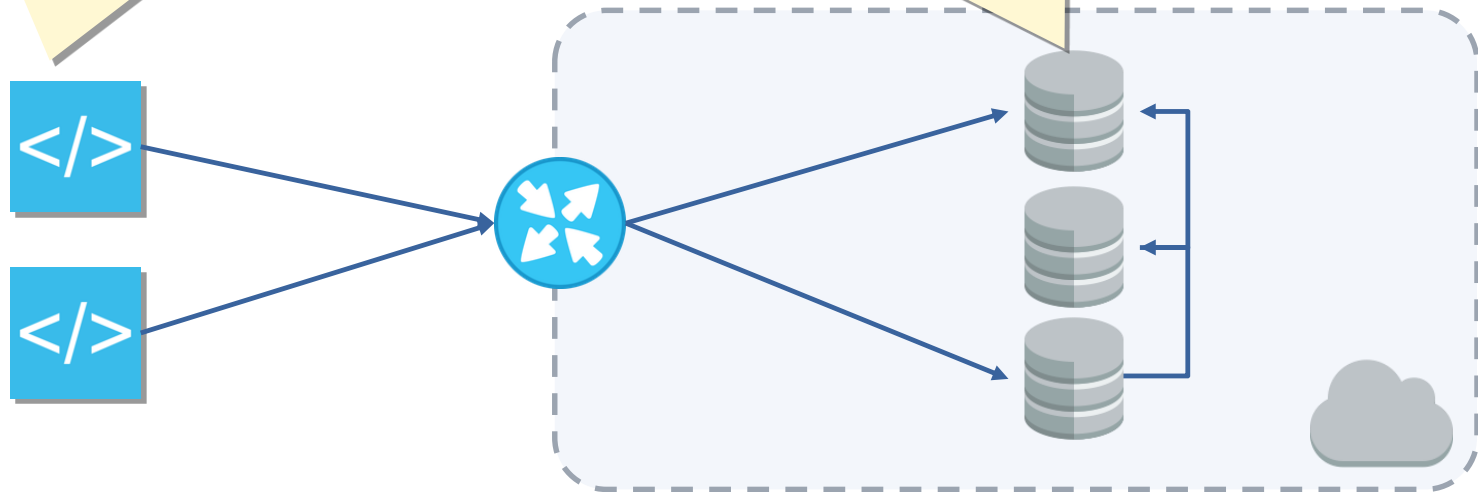
Latency & Polyglot Storage

Two central problems

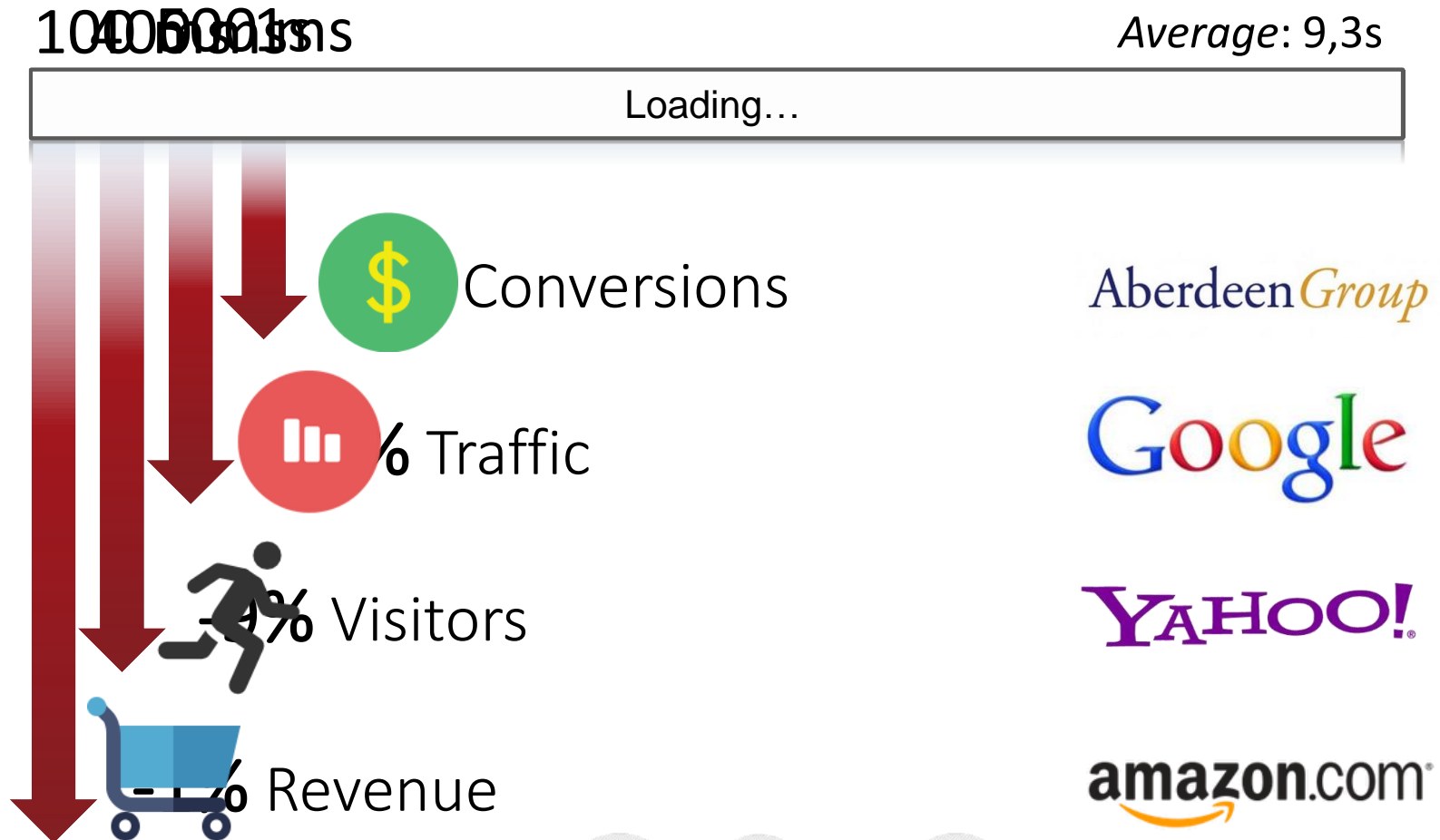
- ▶ Goal of ORESTES: Solve both problems through a scalable **cloud-database middleware**

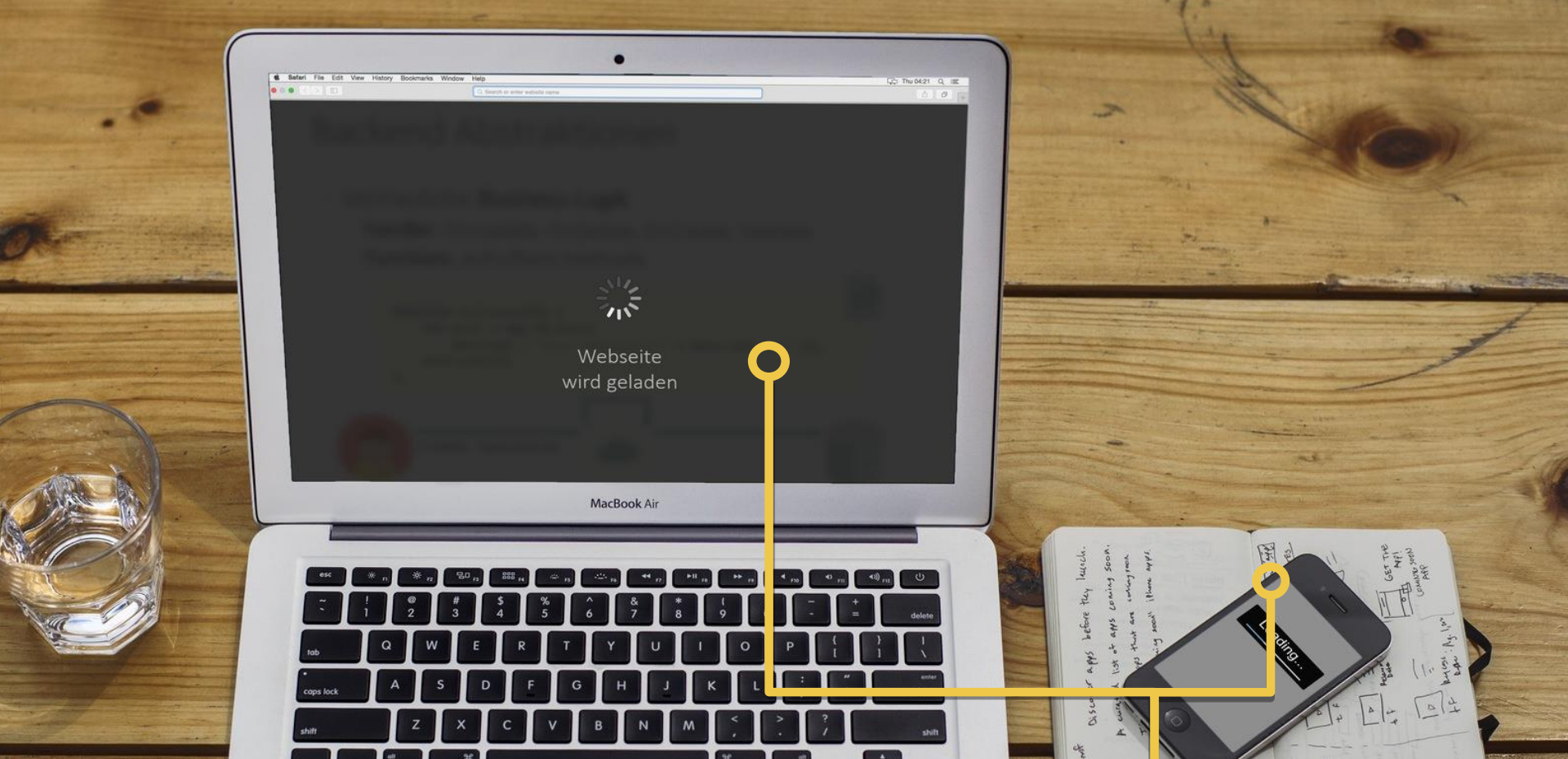
If the application is *geographically distributed*, how can we guarantee **fast database access**?

If one size *doesn't* fit all – how can **polyglot persistence** be leveraged on a declarative, automated basis?



Problem I: Latency





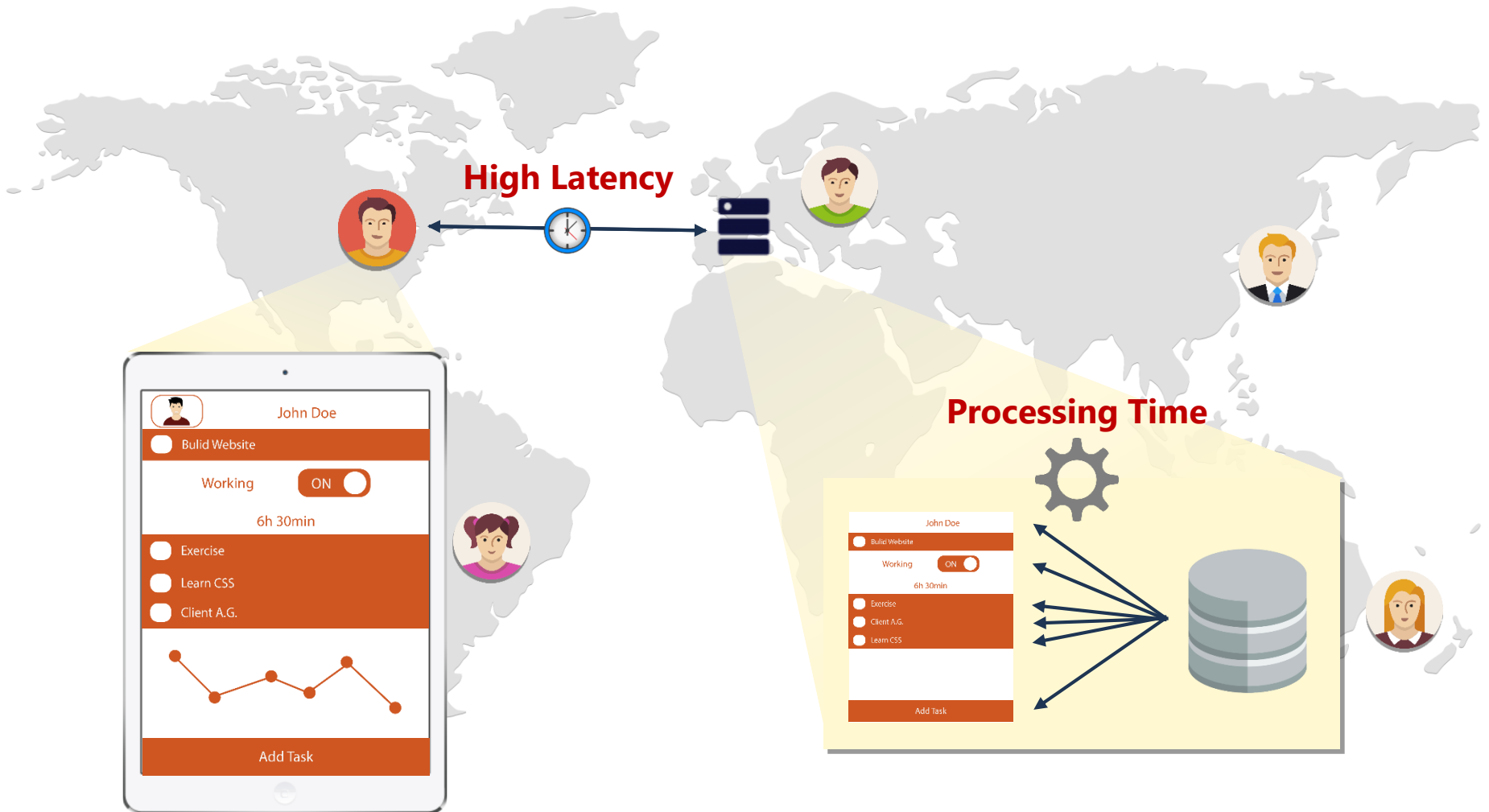
If perceived speed is such an important factor



...what causes slow page load times?

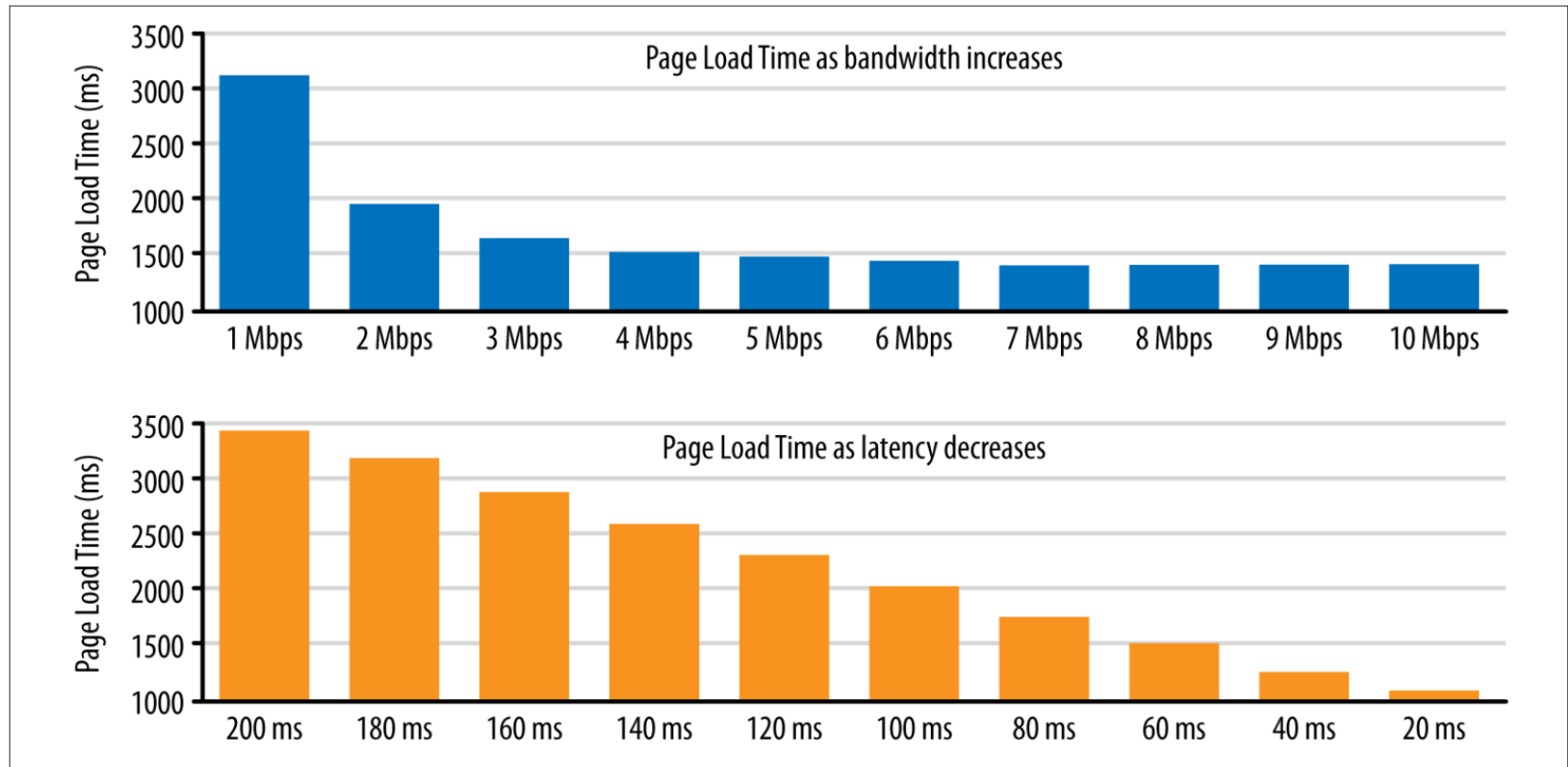
State of the art

Two bottlenecks: latency und processing



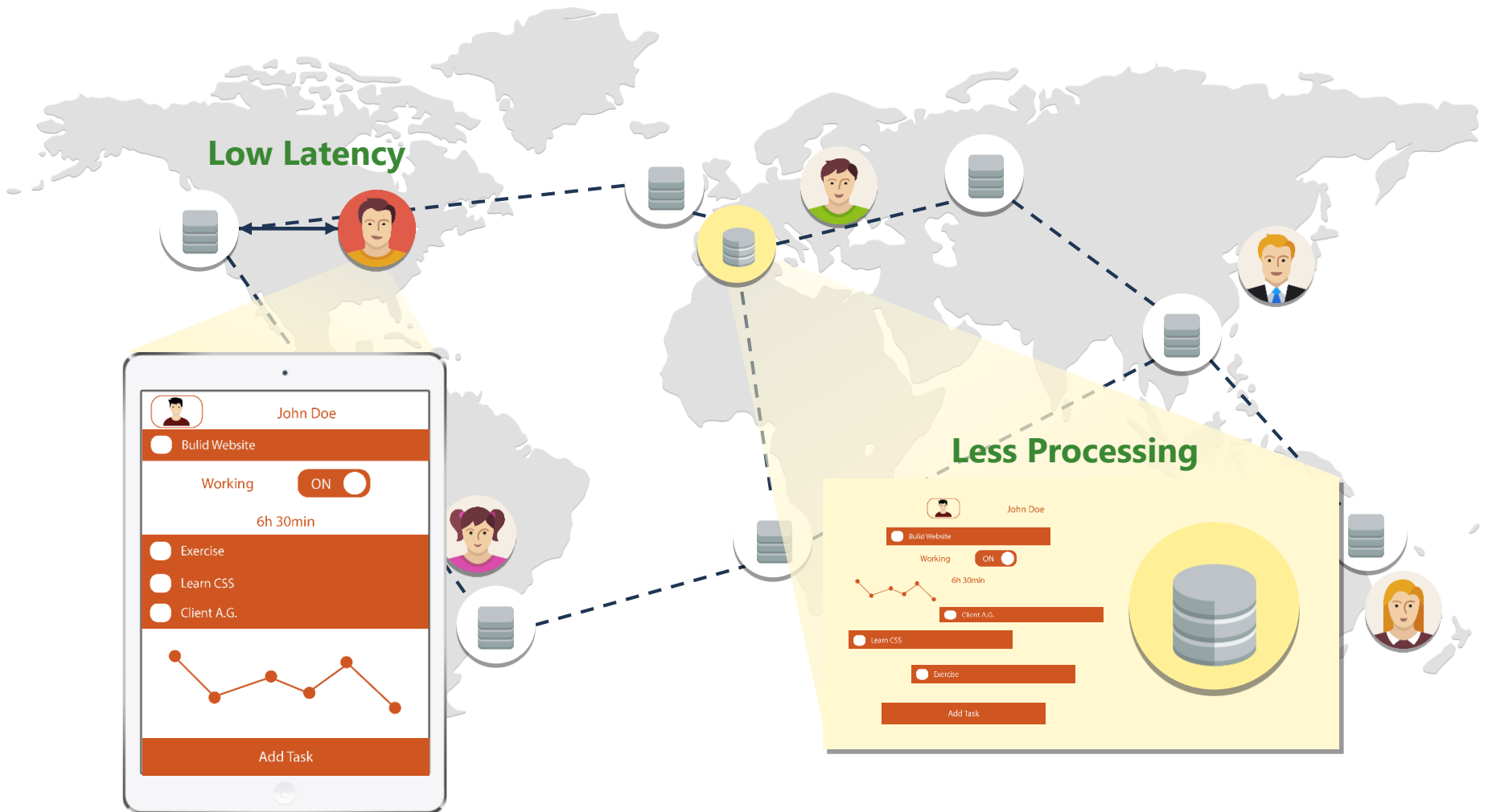
Network Latency

The underlying problem of high page load times



The low-latency vision

Data is served by ubiquitous web-caches



The web's caching model

Staleness as a consequence of scalability



Expiration-based

Every object has a defined Time-To-Live (TTL)



Revalidations

Allow clients and caches to check freshness at the server



Research Question:

Can database services leverage the web caching infrastructure for low latency with rich consistency guarantees?



Stale
Data

Problem II: Polyglot Persistence

Current best practice

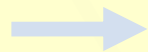


Research Question:

Can we automate the

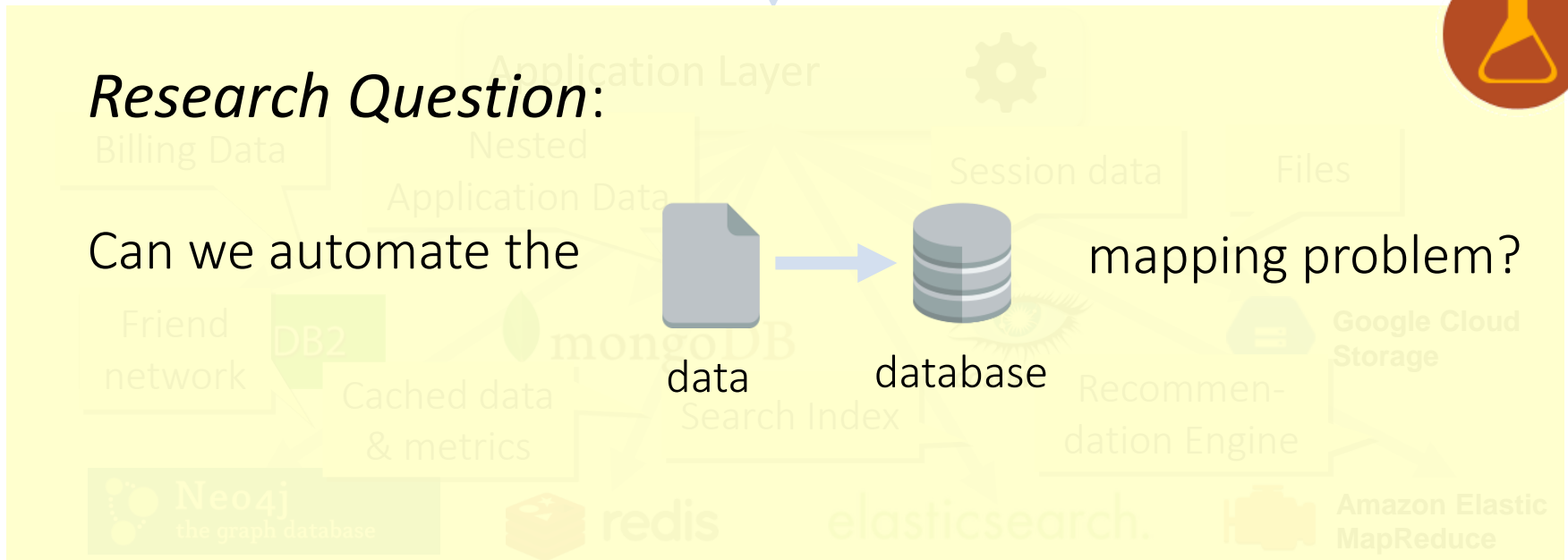


data



database

mapping problem?



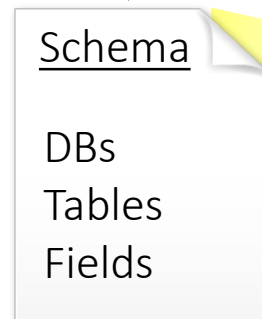
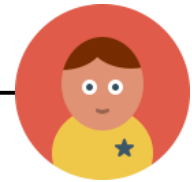
Felix Gessert, Norbert Ritter:
Polyglot Persistence. Datenbankspektrum.

Vision

Schemas can be annotated with requirements

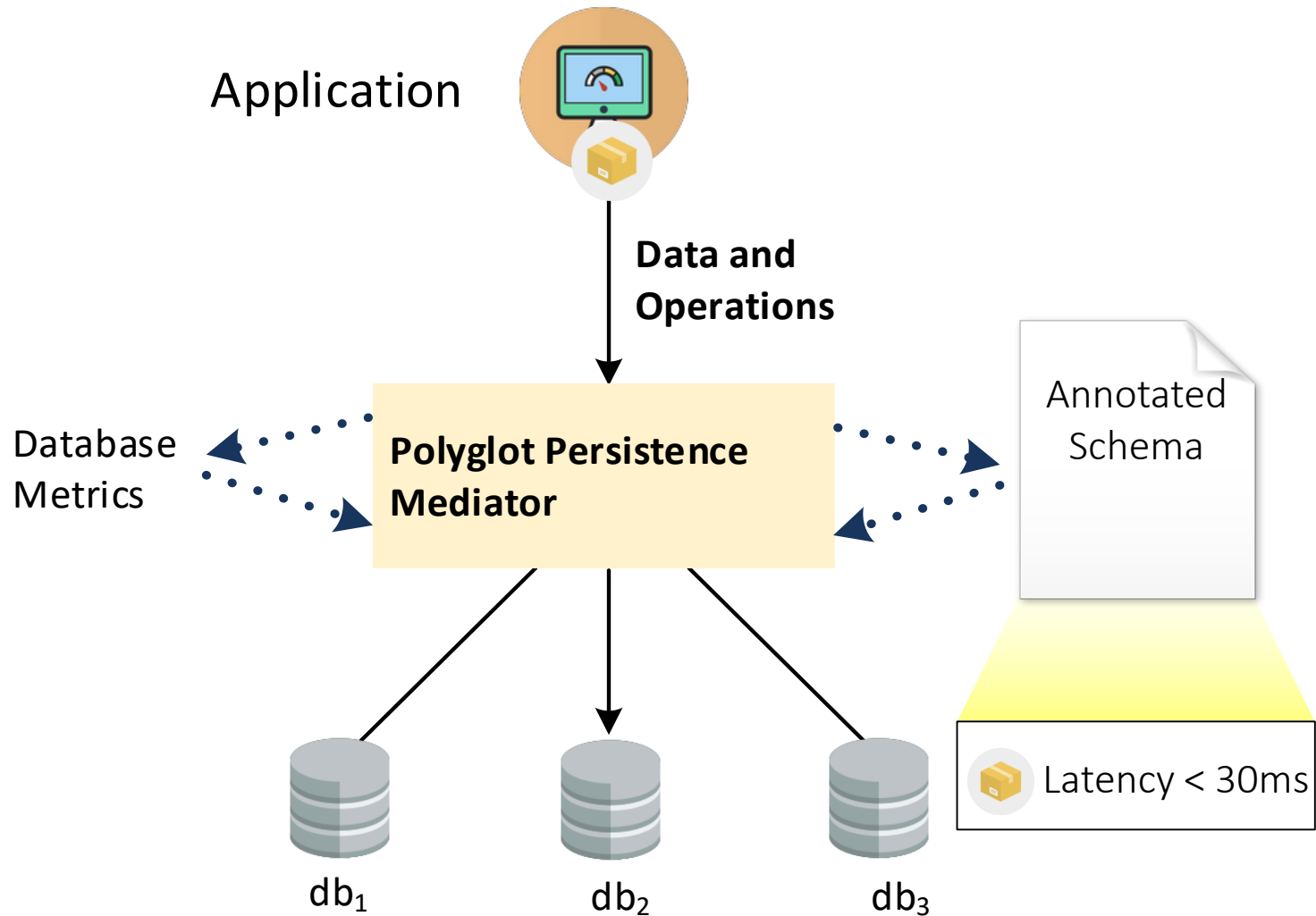


- Write Throughput > 10,000 RPS
- Read Availability > 99.9999%
- Scans = **true**
- Full-Text-Search = **true**
- Monotonic Read = **true**



Vision

The Polyglot Persistence Mediator chooses the database

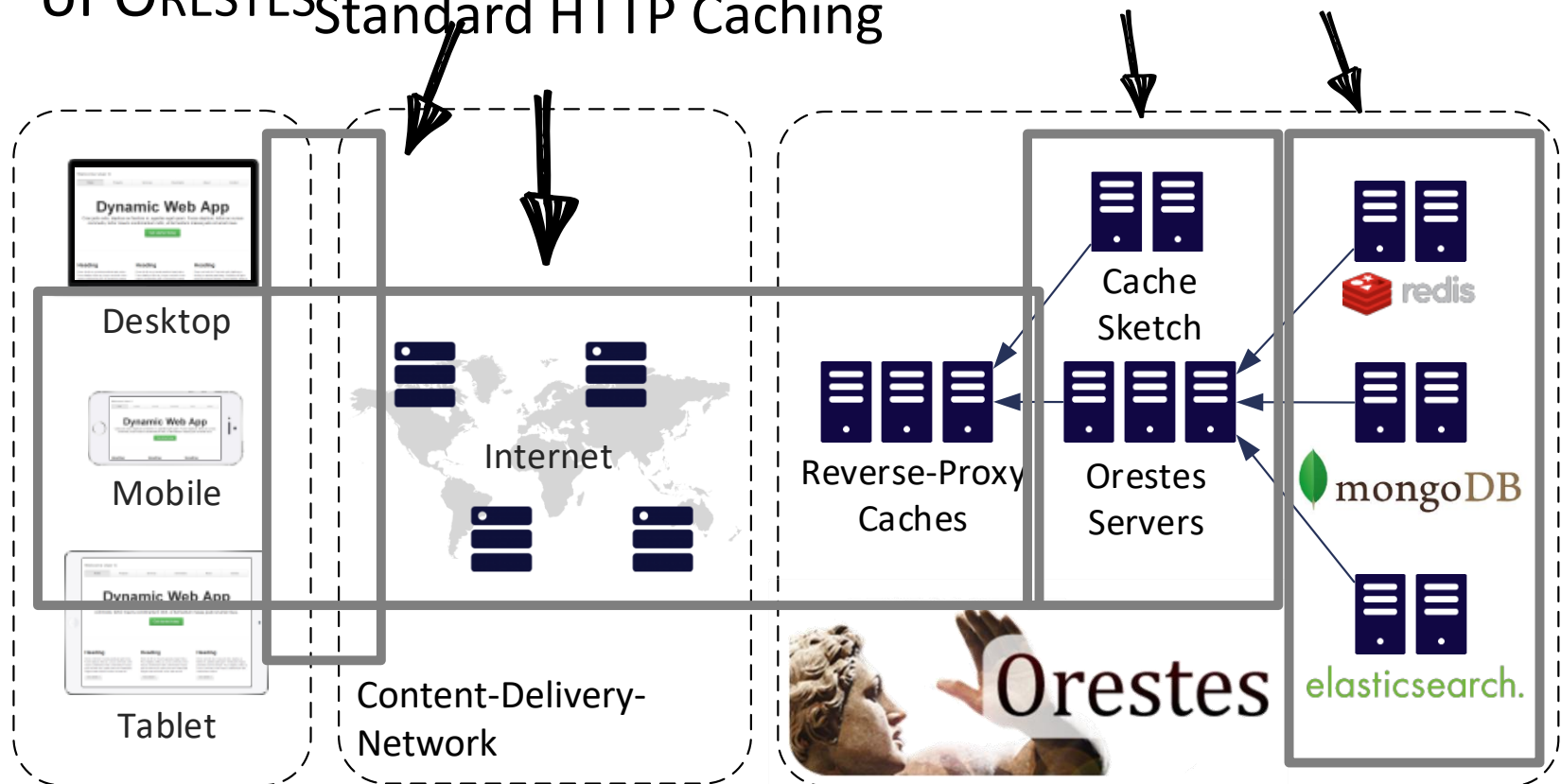


The Big Picture

Implementation in ORESTES

Database-as-a-Service Middleware:

- Polyglot Storage and Low-Latency are the central goals of ORESTES
- Unified REST API
- Standard HTTP Caching
- Caching, Transactions, Schemas, Authorization, Multi-Tenancy



Outline



Motivation



ORESTES: a Cloud-Database Middleware



Solving Latency and Polyglot Storage

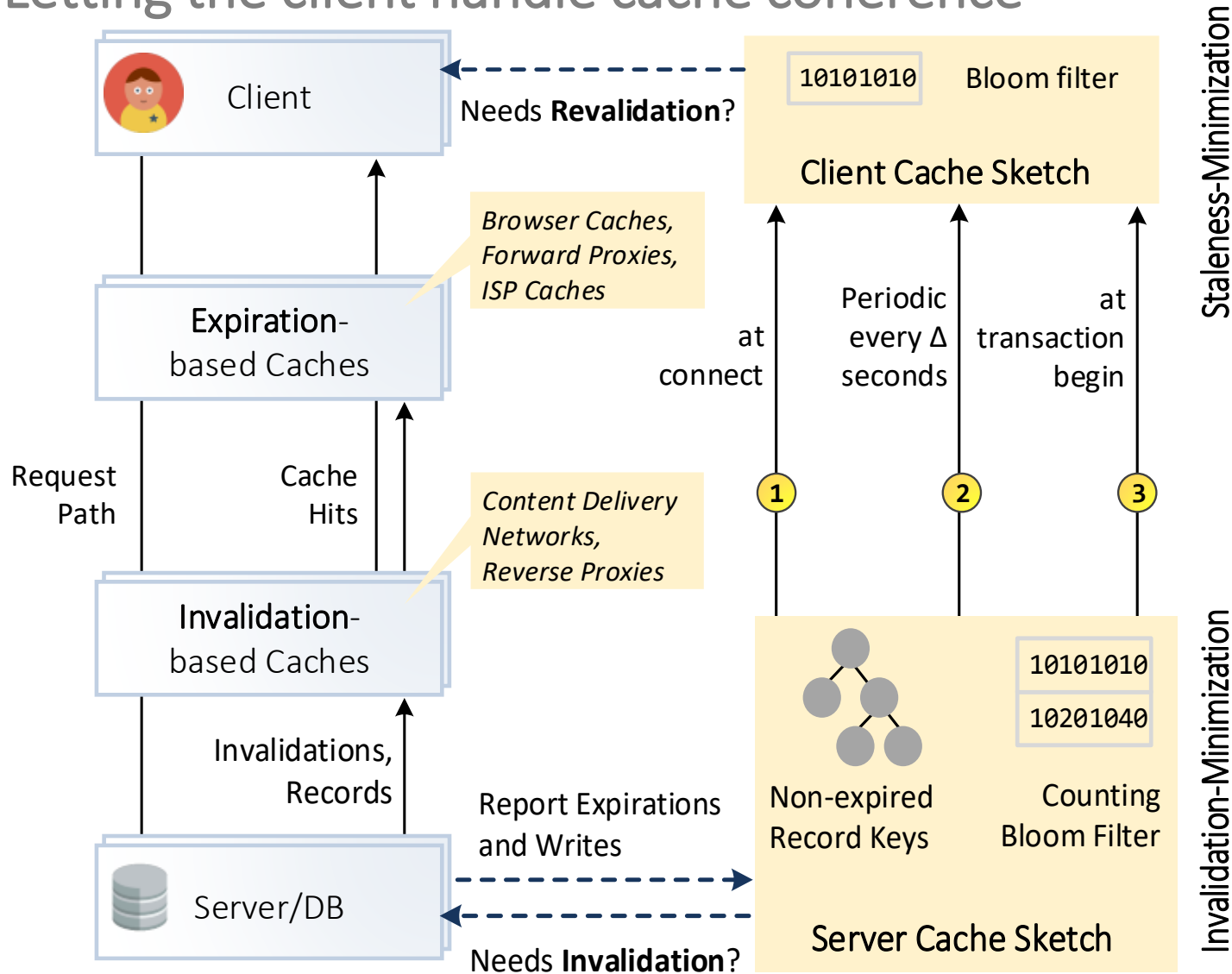


Wrap-up

- Cache Sketch Approach
 - Caching Objects
 - Caching Query Results
 - Continuous Queries
- Polyglot Persistence Mediator
 - Resolution
 - Mediation
 - Polyglot Materialized Views

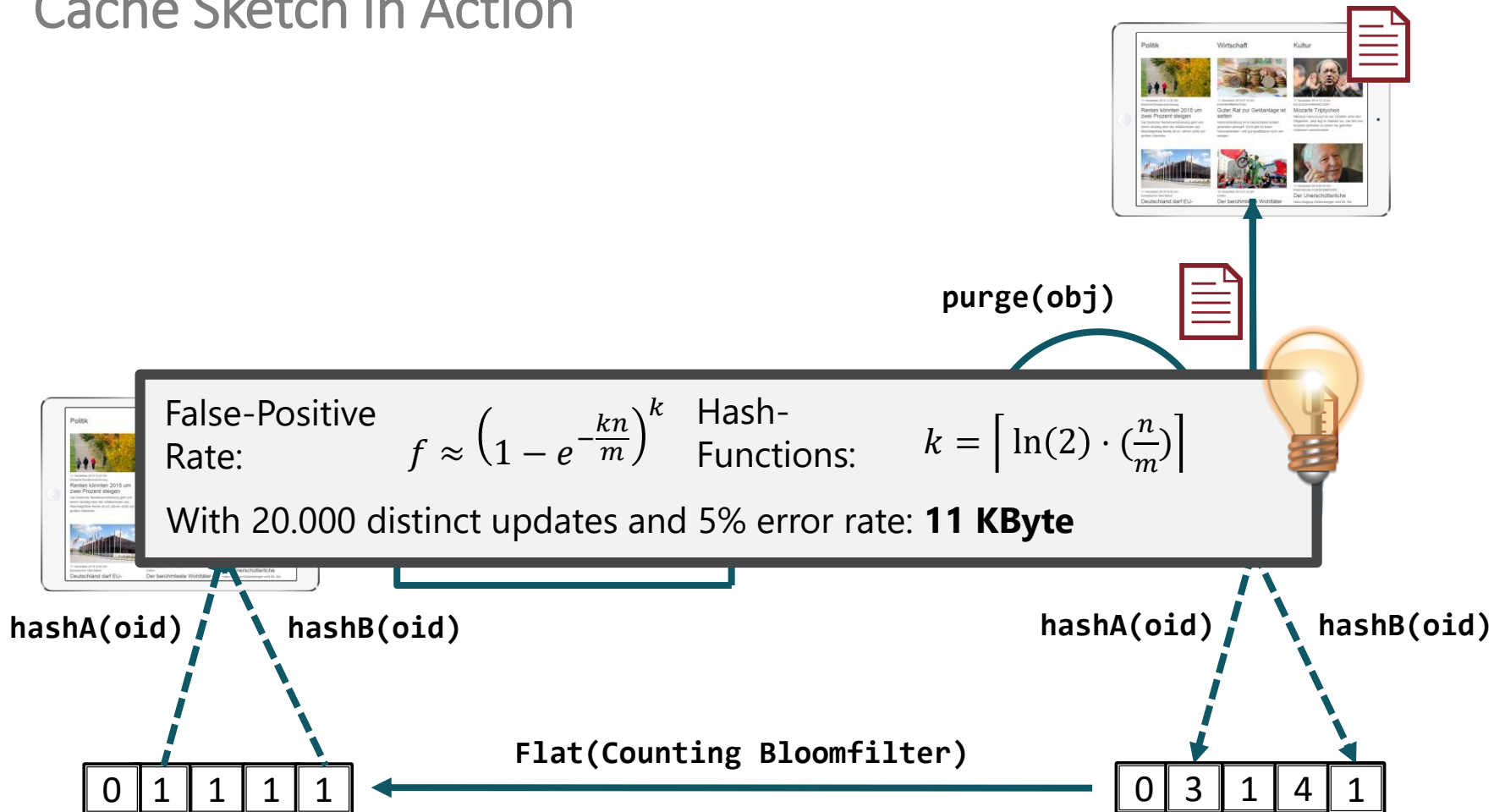
The Cache Sketch approach

Letting the client handle cache coherence



Visually Explained

Cache Sketch in Action



Object Caching

Summary of Properties

- ▶ *Consistency guarantee: Δ -atomicity*
- ▶ **Modes:**
 - ▶ *Cached initialization:* piggybacked Cache Sketch enables fast page loads
 - ▶ *Bounded Staleness:* application refreshes Cache Sketch in fixed intervals
 - ▶ *Conflict-Avoidant Optimistic Transactions:* guarantee ACID despite cached reads
- ▶ TTL Estimator: learns and (statistically) estimates appropriate expirations



Felix Gessert, Florian Bucklers, Norbert Ritter:
Orestes: A scalable Database-as-a-Service architecture for low latency. CloudDB2014@ICDE.



Felix Gessert, Michael Schaarschmidt, Wolfram Wingerath, Steffen Friedrich, Norbert Ritter: *The Cache Sketch: Revisiting Expiration-based Caching in the Age of Cloud Data Management.* BTW 2015

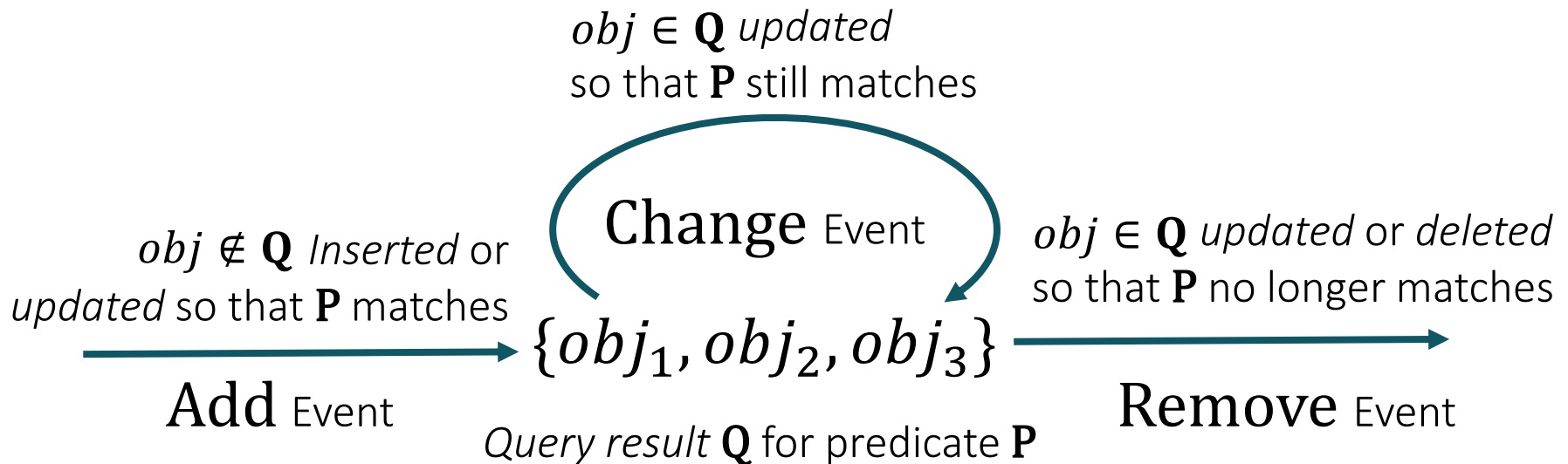
From Object Caching to Query Caching

Generalizing the Cache Sketch to query results

► **Main challenge:** when to invalidate?

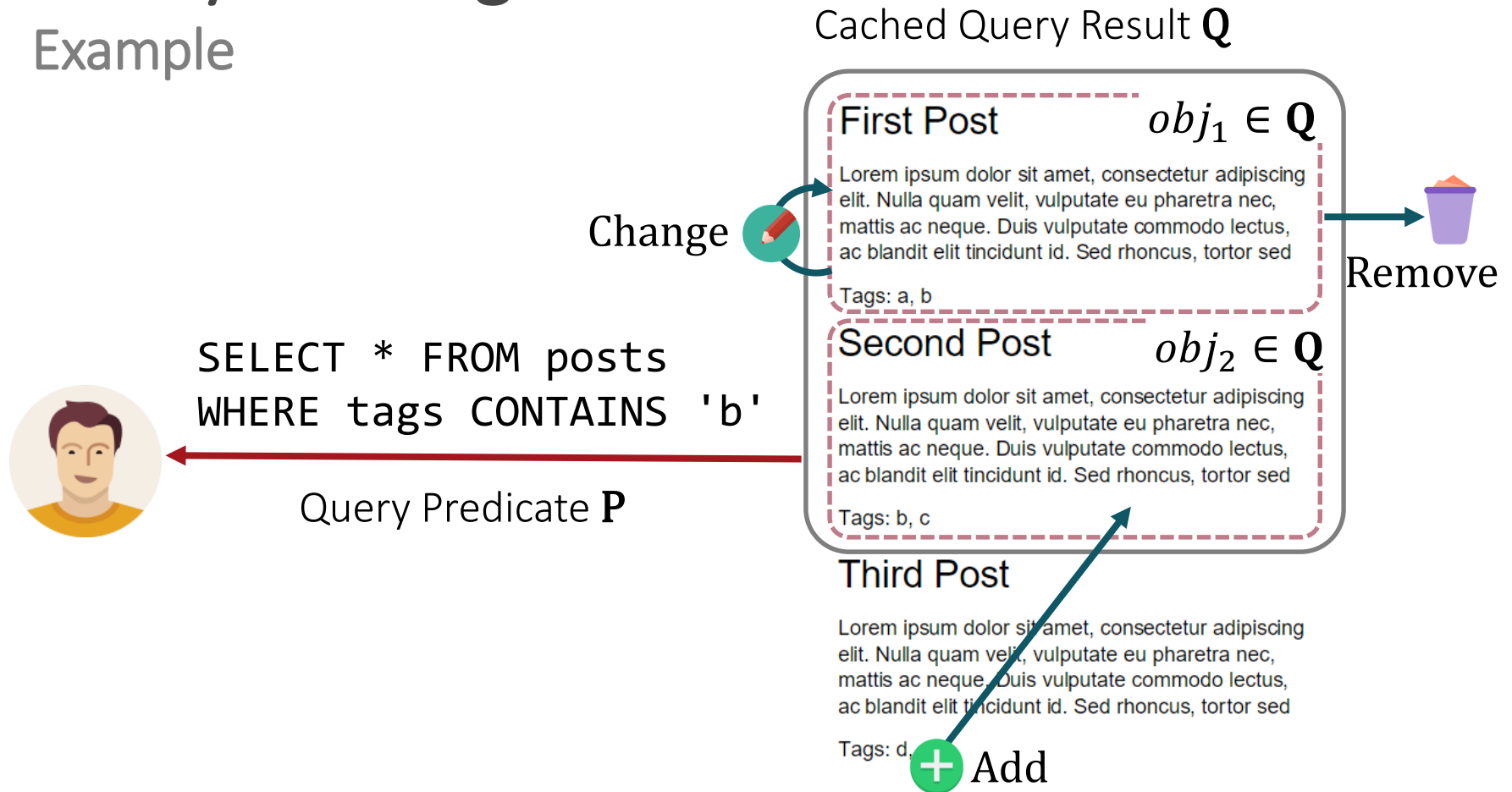
- **Objects:** for every update and delete
- **Queries:** when the query result changes

→ How to detect query **result changes** in real-time?



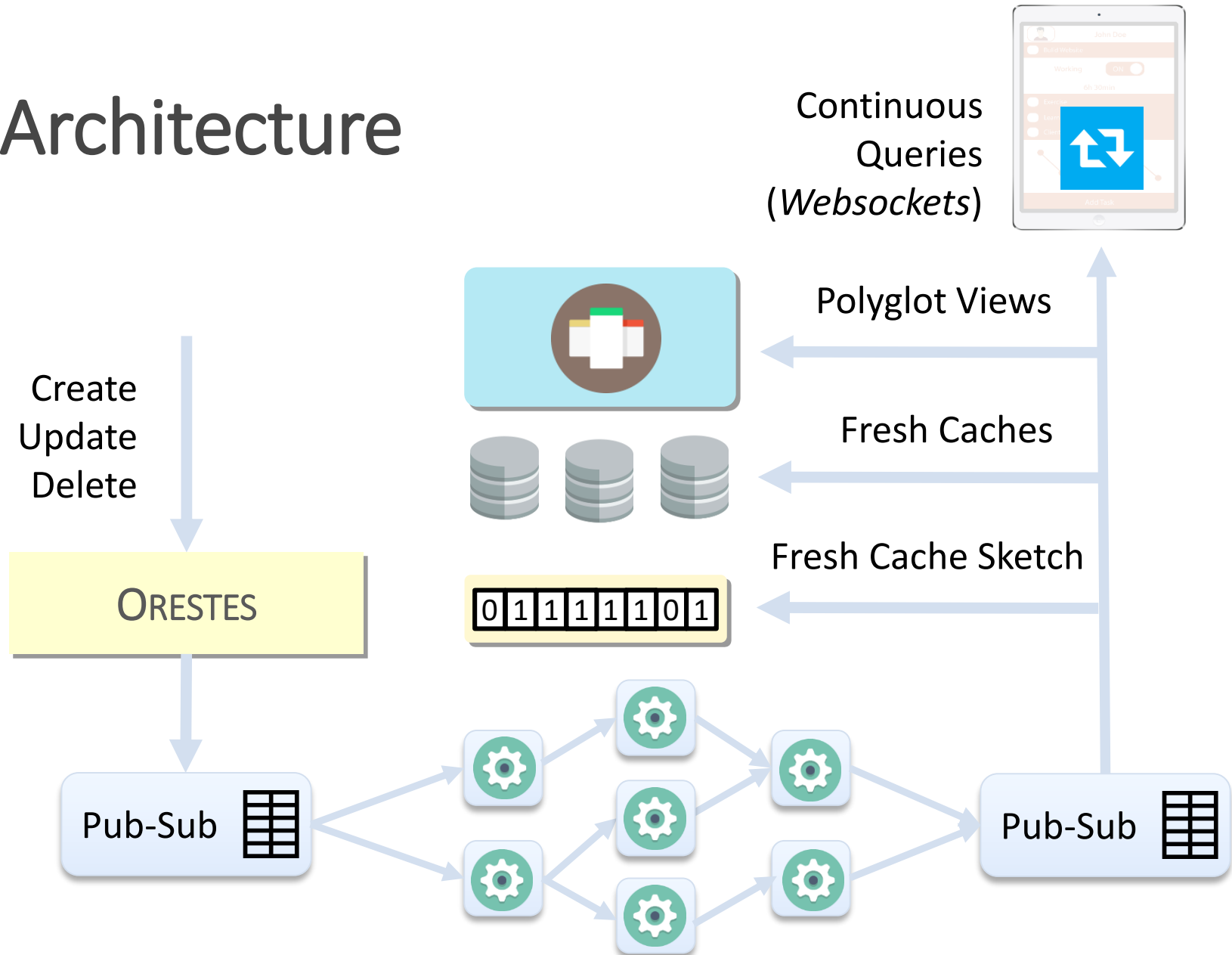
Query Caching

Example



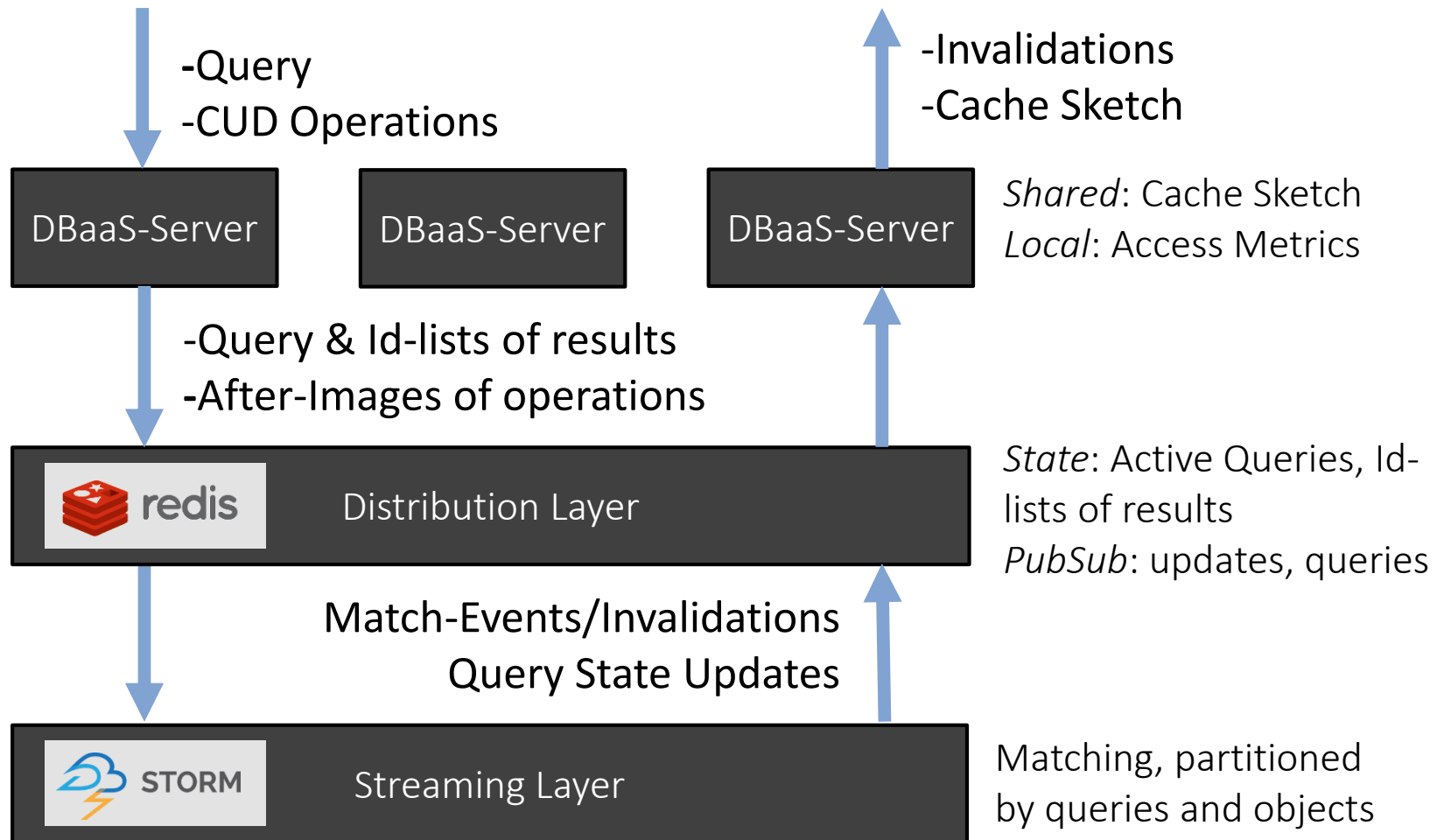
- ▶ **Add, Change, Remove** all entail an invalidation and addition to the cache sketch

Architecture



Architecture

Generalizing the Cache Sketch to Query Results

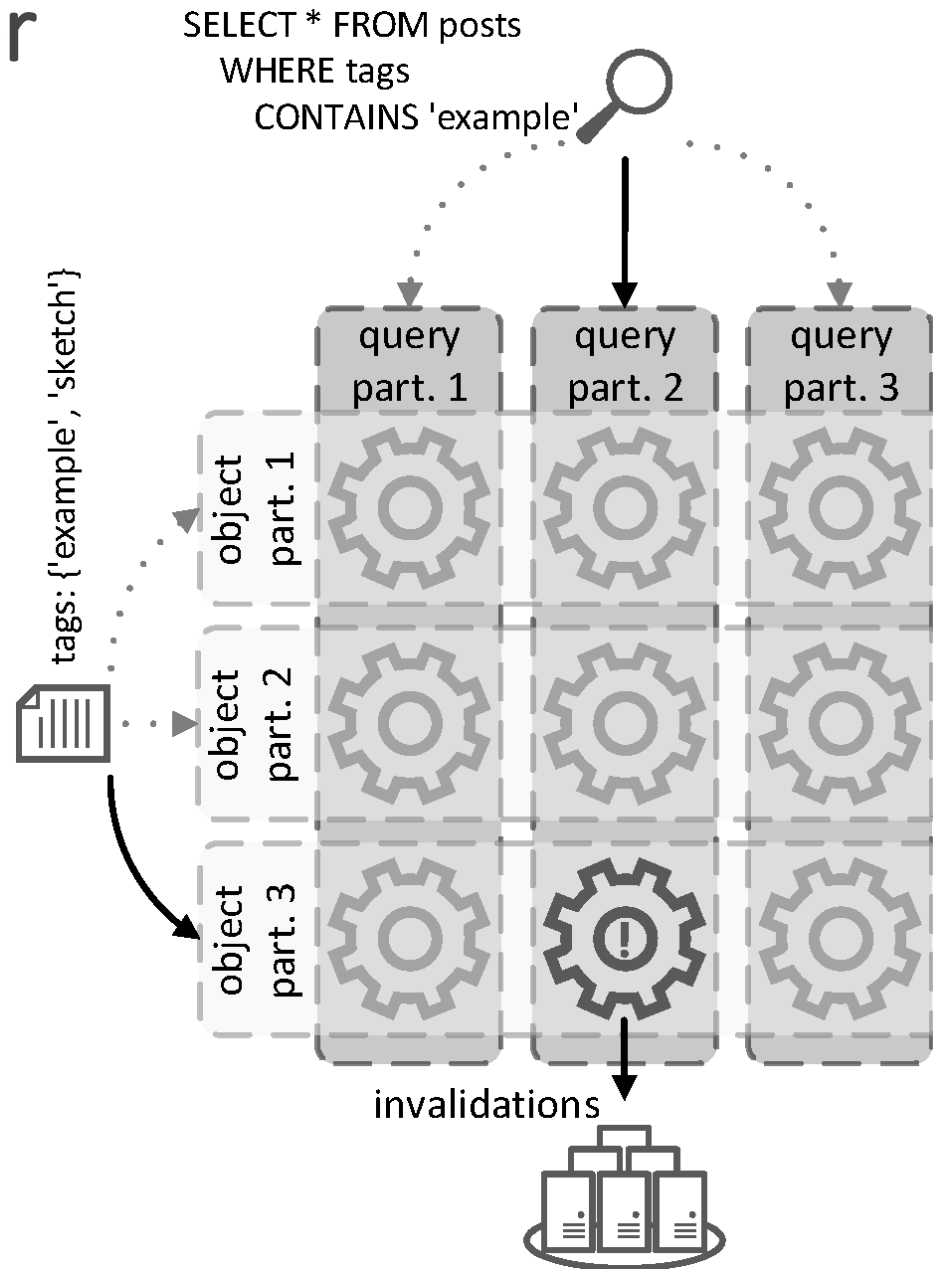


Streaming Layer

Query Matching

Design goals:

- Scalability
- Elasticity
- Low Latency



Optimal Query Representation

Id-Lists

$\{id_1, id_2, id_3\}$

Invalidated by: Add, Remove
→ less invalidations

Performance: at least two
network round-trips

Object-Lists

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

Invalidated by: Add, Remove, Change

Performance: one round-trip
→ lower latency

Cost-based decision model:

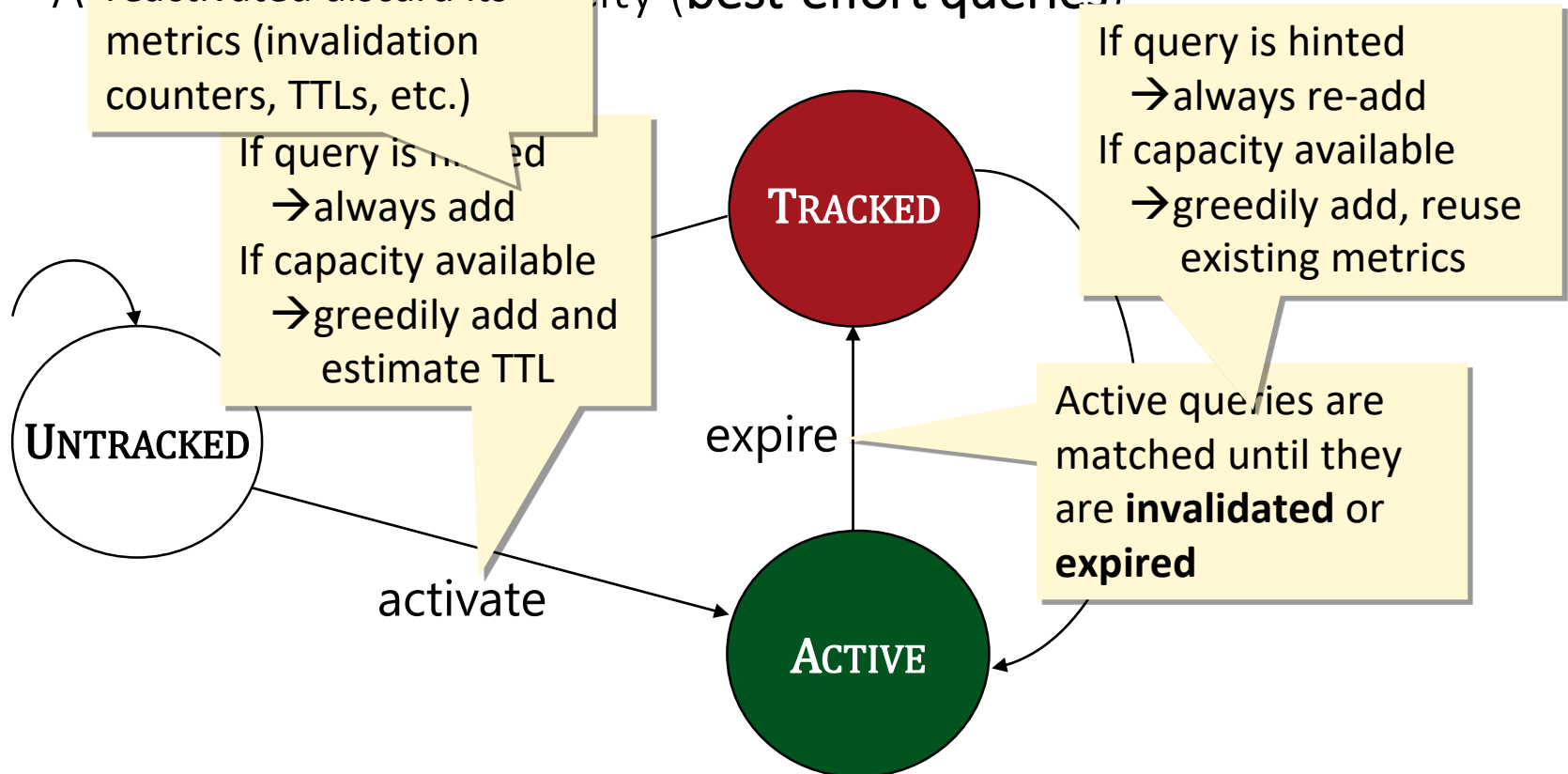
$$w \left(\frac{\text{changes}}{\text{removes} + \text{adds} + \text{changes}} \right) > 1 - \frac{1}{1 + \lceil \text{resultsize} / \text{connections} \rceil}$$

Fraction of avoided invalidations

avoided round-trips

Disitributed Capacity Management

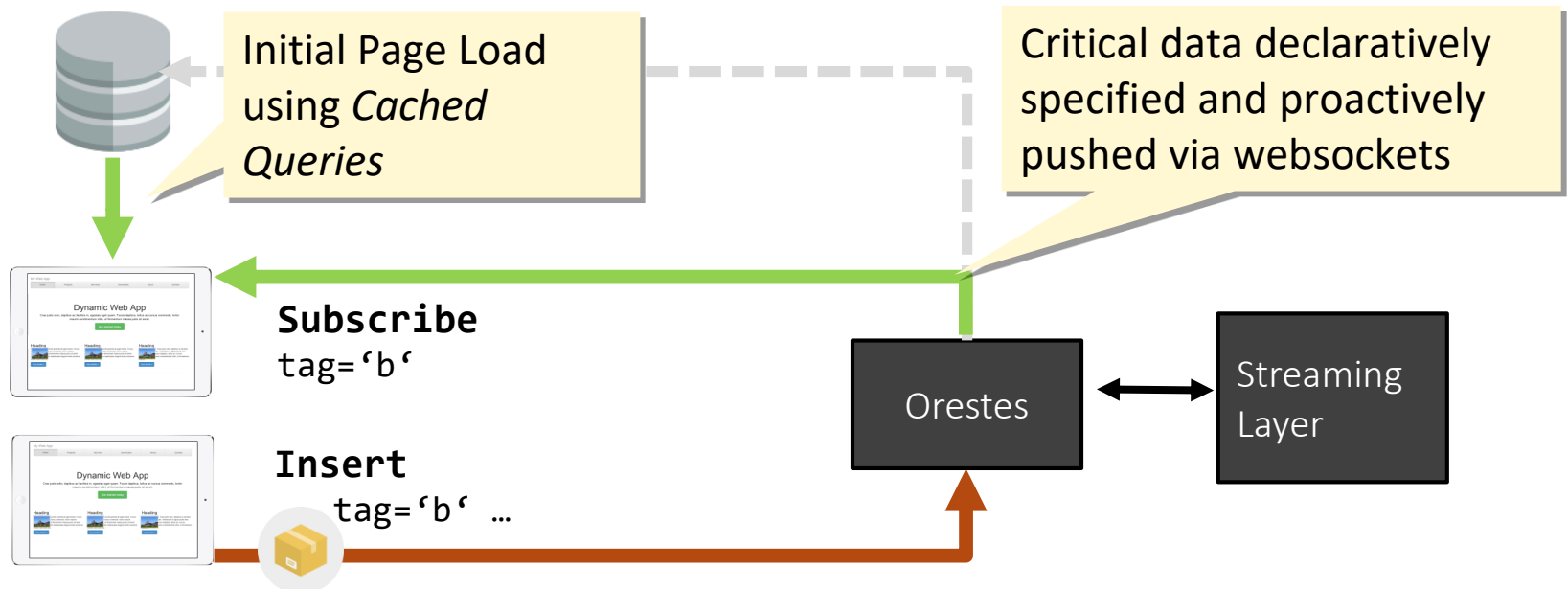
- ▶ Matching capacity is limited
 - A reactivated query discards its capacity (best-effort queries)
 - A reactivated query discards its capacity (best-effort queries)



Continuous Queries

Complementing Cached Queries

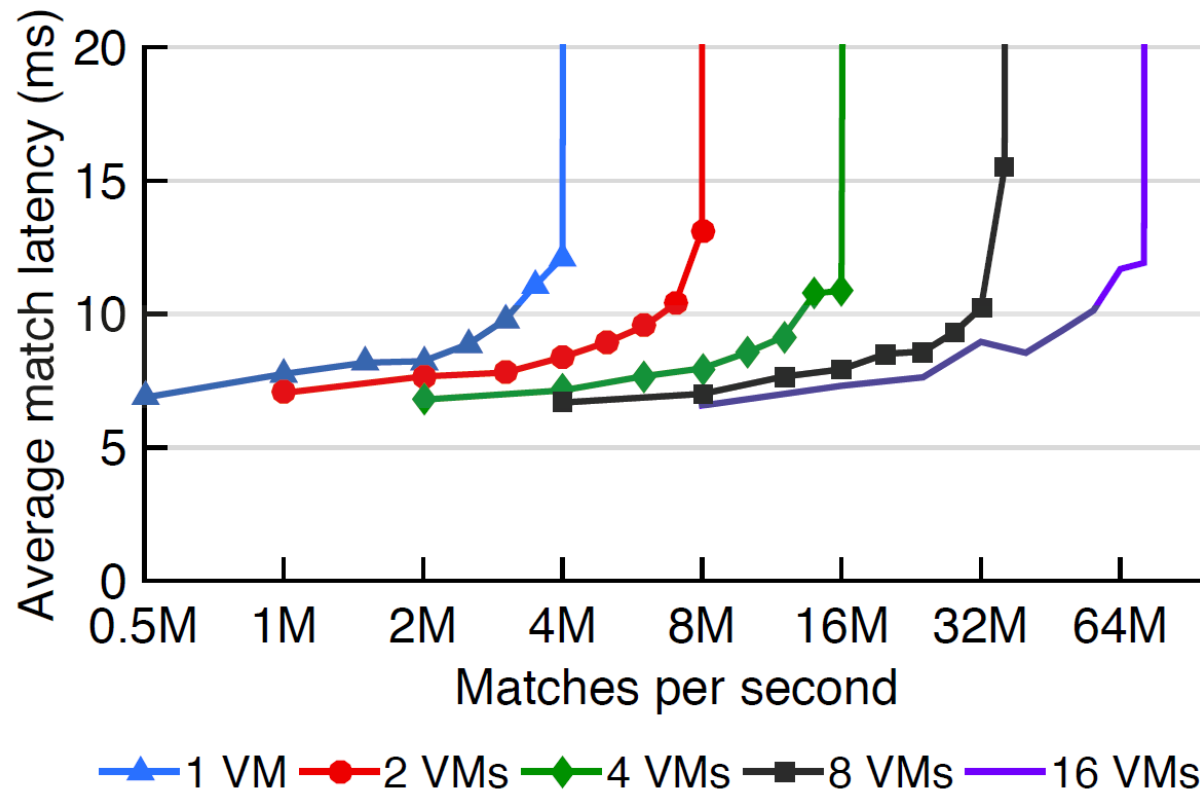
- ▶ Same streaming architecture can similarly notify applications (browsers) about query result changes
- ▶ **Application Pattern:**



Matching Performance

Latency of detecting invalidations

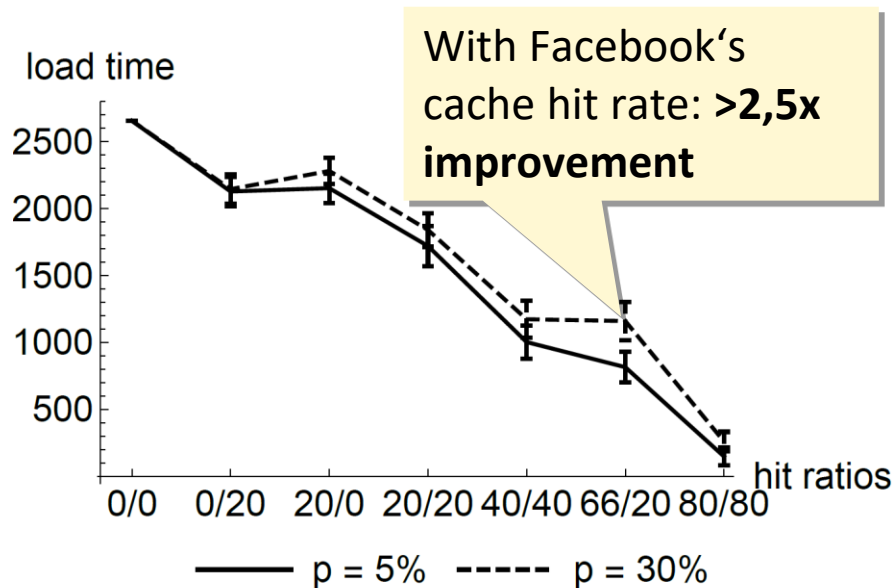
- ▶ Latency mostly < 15ms, scales linearly w.r.t. number of servers and number of tables



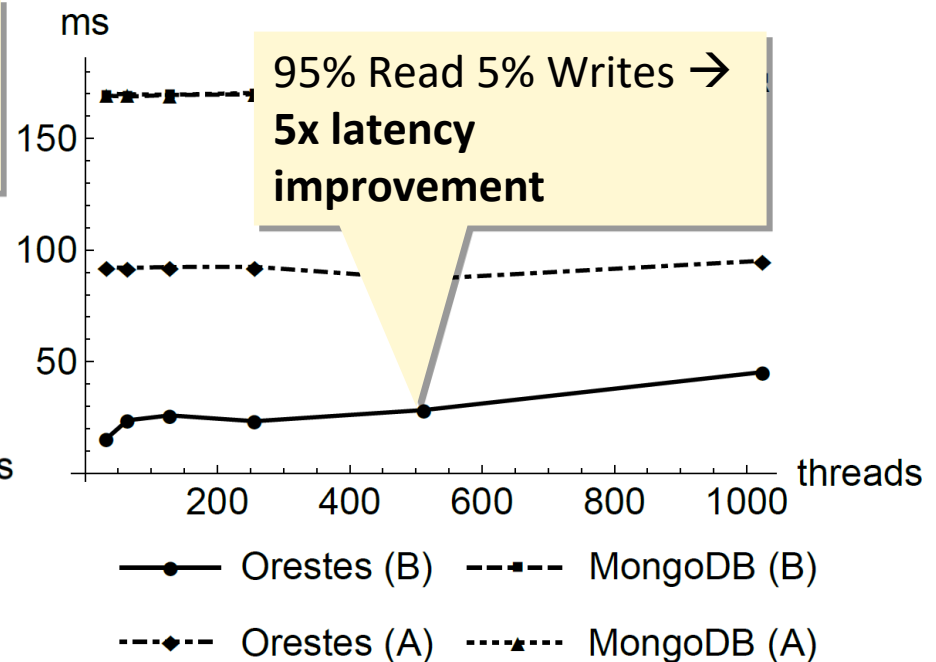
Performance



Page load times with **cached initialization** (simulation):



Average Latency for YCSB Workloads A and B (real):

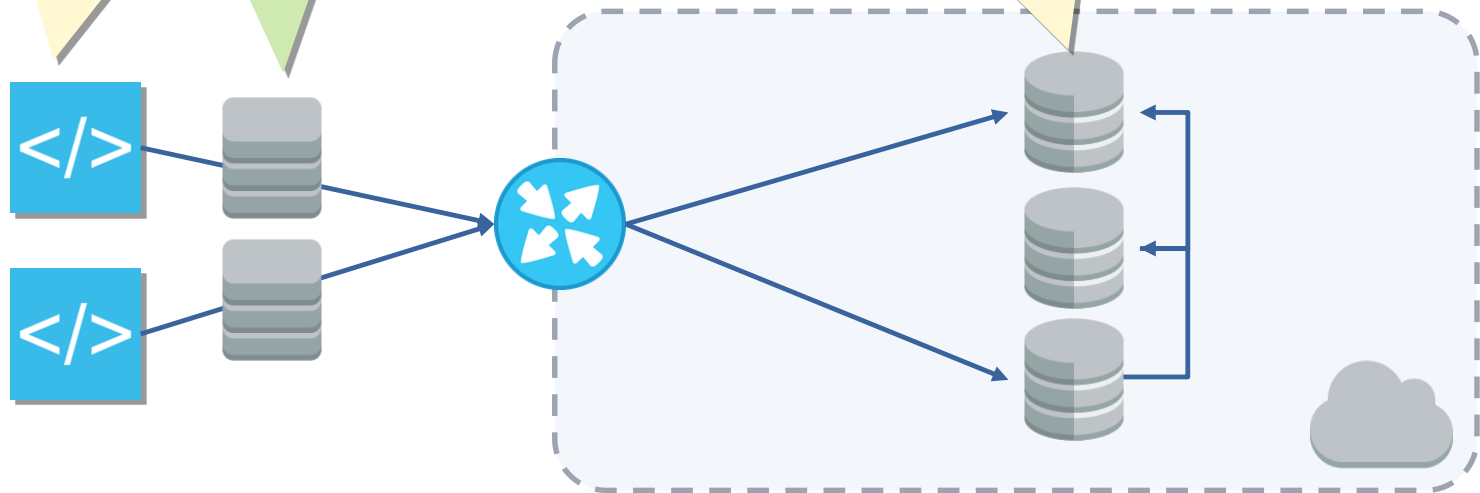


Low Latency

If the appl
distributed
fast datab

Transparent **end-to-end**
caching using the Cache
Sketch.

If one size *doesn't* fit all – how can
polyglot persistence be leveraged
on a declarative, automated basis?



Towards Automated Polyglot Persistence

Necessary steps

► Goal:

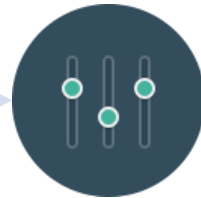
- Extend classic workload management to *polyglot persistence*
- Leverage heterogeneous (NoSQL) databases

1. Requirements



Tenant specifies requirements as Service-Level-Agreements

2. Resolution



Find or provision a suitable combination of databases

3. Mediation



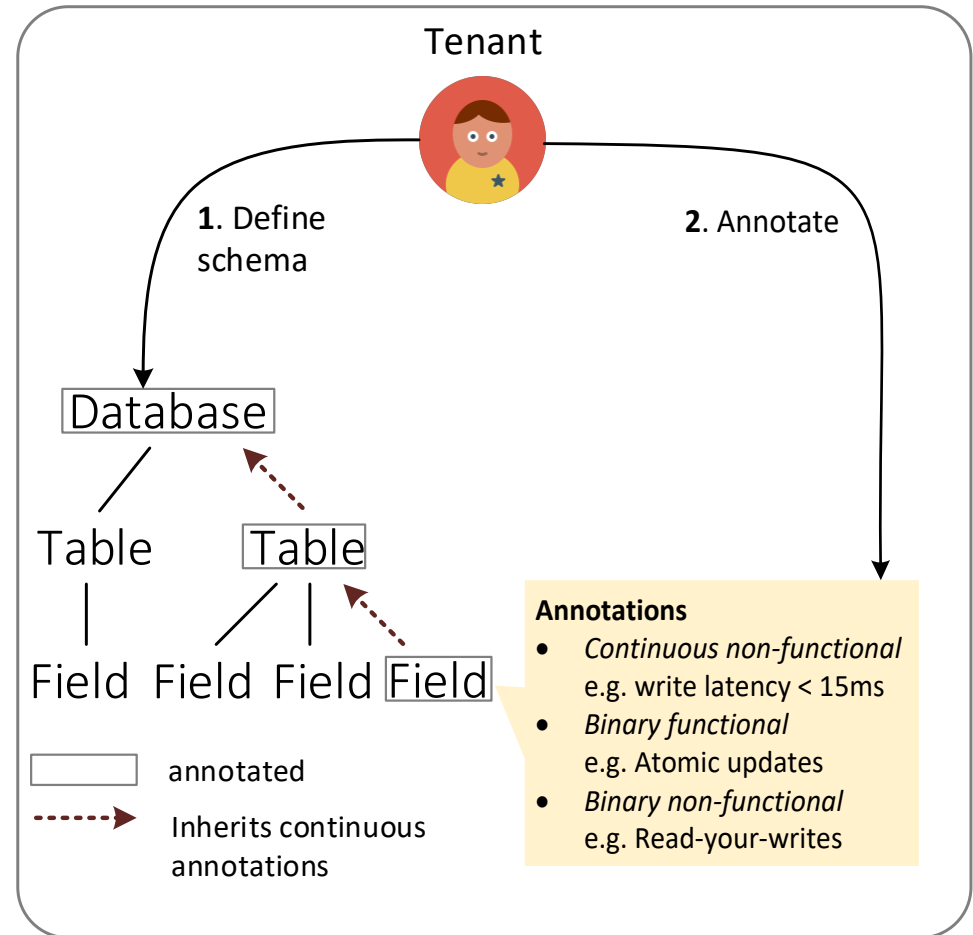
Mediate data and database operations



Step I - Requirements

Expressing the application's needs

Annotation	Type	Annotated at
Read Availability	Continuous	*
Write Availability	Continuous	*
Read Latency	Continuous	*
Write Latency	Continuous	*
Write Throughput	Continuous	*
Data Vol. Scalability	Non-Functional	Field/Class/DB
Write Scalability	Non-Functional	Field/Class/DB
Read Scalability	Non-Functional	Field/Class/DB
Elasticity	Non-Functional	Field/Class/DB
Durability	Non-Functional	Field/Class/DB
Replicated	Non-Functional	Field/Class/DB
Linearizability	Non-Functional	Field/Class
Read-your-Writes	Non-Functional	Field/Class
Causal Consistency	Non-Functional	Field/Class
Writes follow reads	Non-Functional	Field/Class
Monotonic Read	Non-Functional	Field/Class
Monotonic Write	Non-Functional	Field/Class
Scans	Functional	Field
Sorting	Functional	Field
Range Queries	Functional	Field
Point Lookups	Functional	Field
ACID Transactions	Functional	Class/DB
Conditional Updates	Functional	Field
Joins	Functional	Class/DB
Analytics Integration	Functional	Field/Class/DB
Fulltext Search	Functional	Field
Atomic Updates	Functional	Field/Class

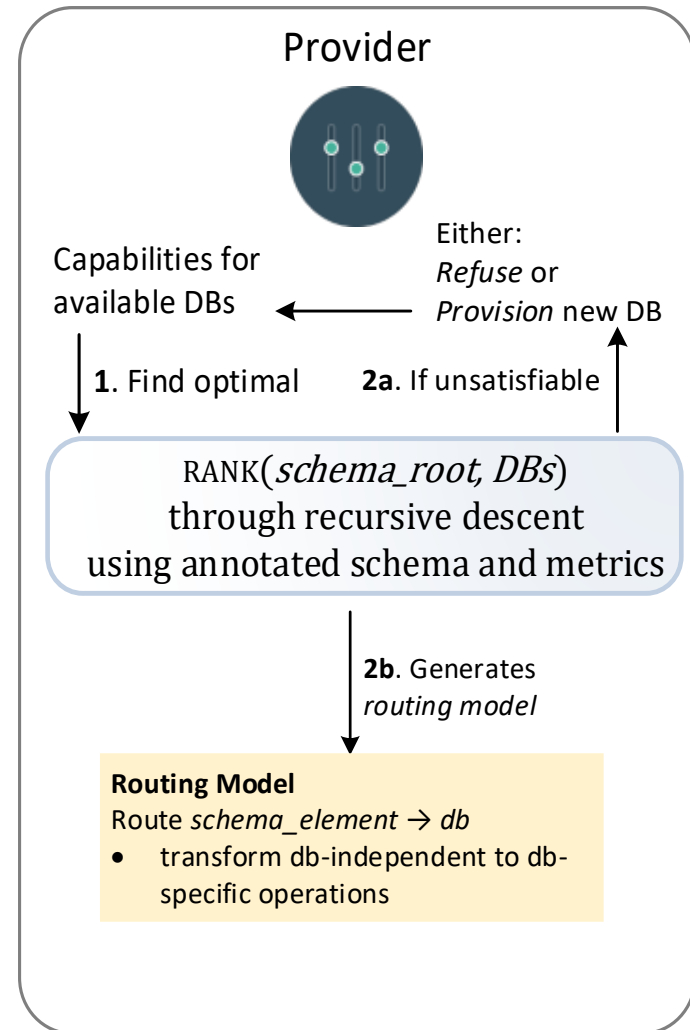


1 Requirements

Step II - Resolution

Finding the best database

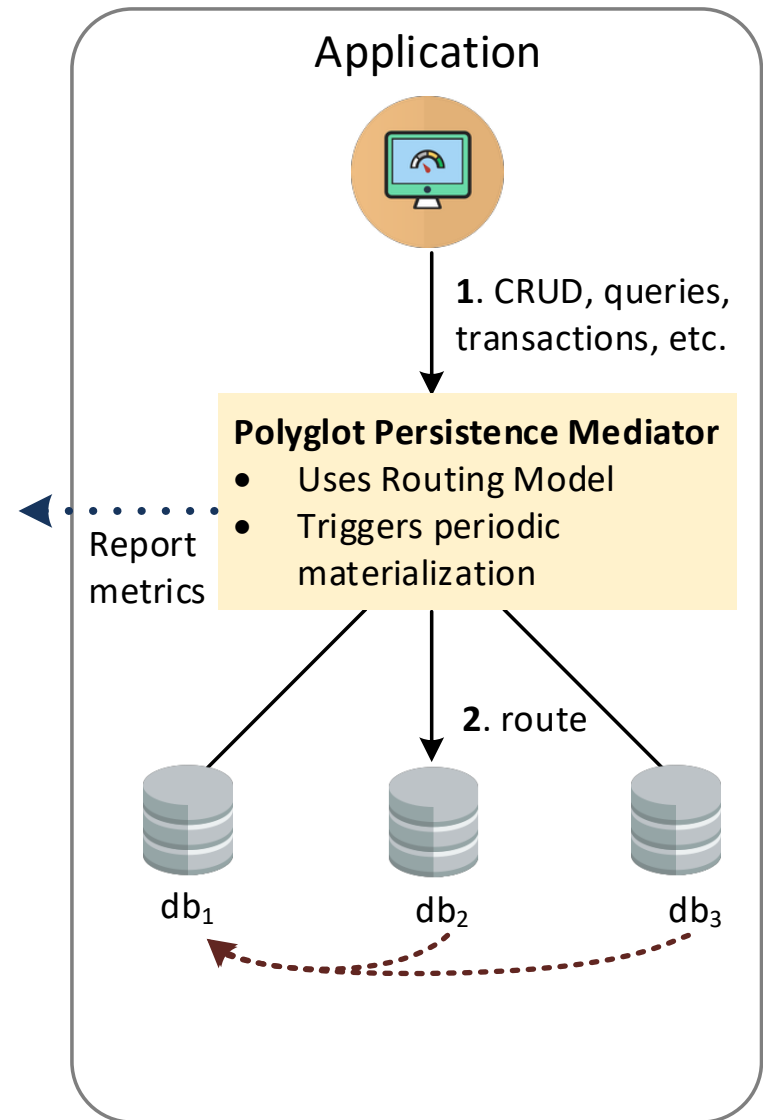
- ▶ The Provider resolves the requirements
- ▶ **RANK**: scores available database systems
- ▶ **Routing Model**: defines the optimal mapping from schema elements to databases



Step III - Mediation

Routing data and operations

- ▶ The PPM routes data
- ▶ **Operation Rewriting:** translates from abstract to database-specific operations
- ▶ **Runtime Metrics:** Latency, availability, etc. are reported to the resolver
- ▶ **Primary Database Option:** All data periodically gets materialized to designated database



Evaluation: News Article

Prototype of Polyglot Persistence Mediator in ORESTES

Scenario: news articles with impression counts

Objectives: low-latency top-k queries, high-throughput on counter, article-queries



Evaluation: News Article

Prototype built on ORESTES

Scenario: news articles with impression counts

Objectives: low-latency top-k queries, high-throughput on counter, article-queries



Counter updates kill performance

Evaluation: News Article

Prototype built on ORESTES

Scenario: news articles with impression counts

Objectives: low-latency top-k queries, high-throughput on counter, article-queries



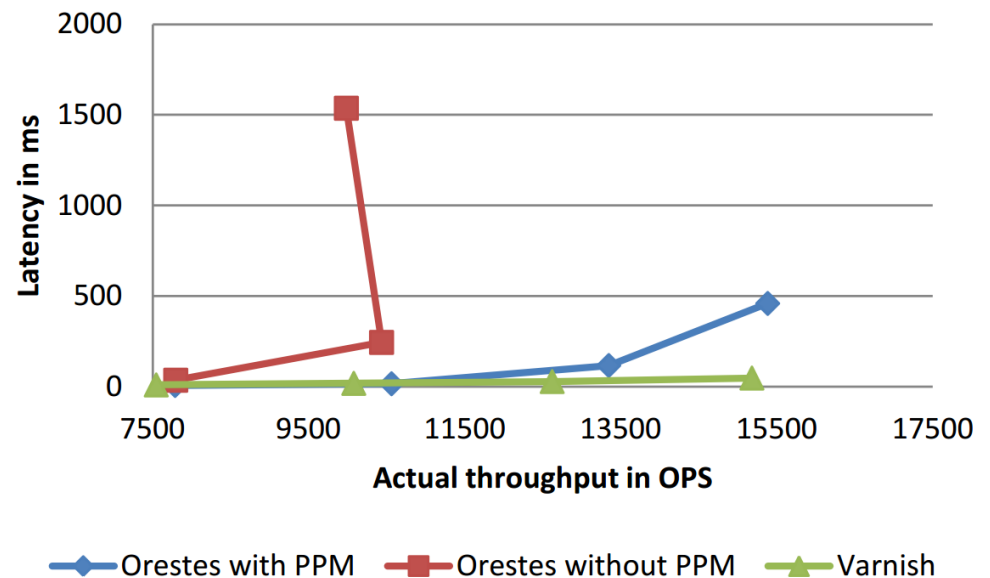
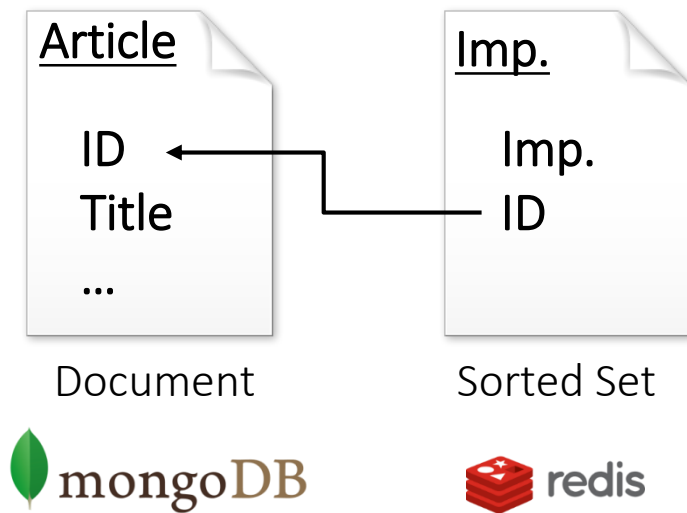
No powerful queries

Evaluation: News Article

Prototype built on ORESTES

Scenario: news articles with impression counts

Objectives: low-latency top-k queries, high-throughput on counter, article-queries



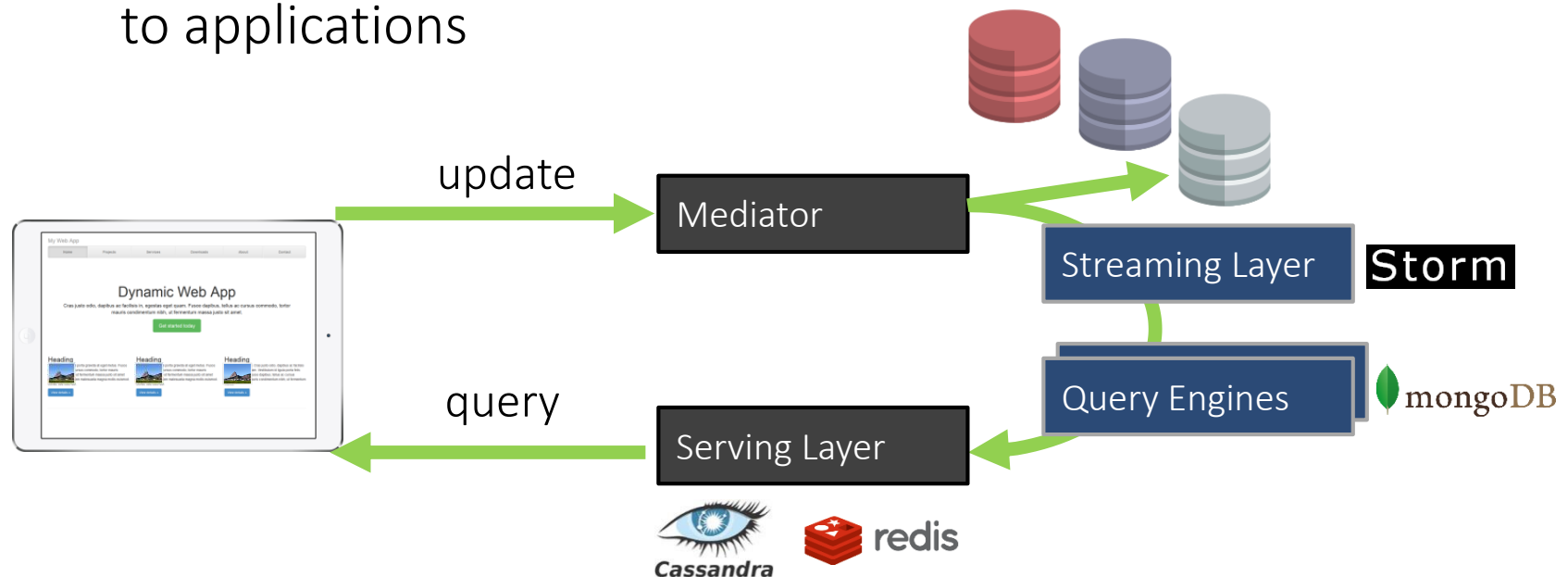
Found Resolution

Polyglot Materialized Views

Arbitrary Queries over arbitrary databases

► Approach:

- **Mediator** emits change data stream (after-images)
- **Streaming layer** maintains registered materialized views using pluggable query engines
- **Serving layer** stores materialized views and serves them to applications



Outline



Motivation



ORESTES: a Cloud-
Database Middleware



Solving Latency and
Polyglot Storage



Wrap-up

- Current/Future Work
- Summary
- Putting ORESTES into practice

Summary



- ▶ **Cache Sketch** (web caching for database services):
 - Consistent (Δ -atomic) *expiration-based* caching
 - *Invalidation-based* caching with minimal purges
- ▶ **Query Caching:**
 - Invalidations and Cache Sketch updates in real-time
 - Cache-optimal representation of results
- ▶ **Continuous & Materialized Queries**
 - Real-time updates to query results
- ▶ **Polyglot Persistence Mediator:**
 - SLA-based routing of queries and data to appropriate database systems

A dramatic photograph of a space shuttle launching, with a massive plume of white smoke and fire at the base. The shuttle is white with a gold-colored nose cone. The sky is blue with scattered white clouds.

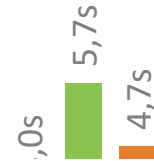
BaQend

Build faster Apps **faster.**

www.baqend.com

Page-Load Times

What impact does the Cache Sketch have?



Politik



11. November 2014 12:42 Uhr
Deutsche Rentenversicherung

Renten könnten 2015 um zwei Prozent steigen

Die Deutsche Rentenversicherung geht von einem Anstieg über der Inflationsrate aus. Abschlagsfreie Rente ab 63 Jahren stößt auf großes Interesse.

Wirtschaft



11. November 2014 07:15 Uhr
HONORARBERATUNG

Guter Rat zur Geldanlage ist selten

Honorarberatung ist in Deutschland endlich gesetzlich geregelt. Doch gibt es kaum Honorarberater. Und gut qualifizierte noch viel weniger.

Kultur



11. November 2014 10:14 Uhr
NICOLAUS HARNONCOURT

Mozarts Triptychon

Nikolaus Harnoncourt ist der Detektiv unter den Dirigenten. Jetzt legt er Indizien vor, wie drei von Mozarts Sinfonien zu einem nie gehörten Oratorium verschmelzen.



11. November 2014 10:05 Uhr
Europäischer Gerichtshof

Deutschland darf EU-Ausländern Hartz IV verweigern

Der Europäische Gerichtshof hat entschieden: Deutschland kann arbeitslose Zuwanderer aus der EU von Sozialleistungen ausschließen. Das Urteil könnte ein Signal sein.



10. November 2014 21:32 Uhr
CHINA

Der berühmteste Wohltäter Chinas – nach eigenen Angaben

Der chinesische Unternehmer Chen Guangbiao hat sich als Wohltäter sehr bekannt gemacht. Er schenkt Geldbündeln und zertrümmert öffentlich Luxusautos.



11. November 2014 06:39 Uhr
HANS MAGNUS ENZENSBERGER

Der Unerschütterliche

Hans Magnus Enzensberger wird 85. Ein Besuch bei dem herzlich eigenwilligen Intellektuellen. Mit Tumult hat er gerade ein erstaunlich persönliches Buch veröffentlicht.



11. November 2014 08:48 Uhr
APEC-GIPFELTREFFEN

Obama besänftigt China

Die USA wollen China in der APEC-Präsident Obama vor den Staatschef Xi. Der plädiert für mehr wirtschaftliche Verflechtung.



10. November 2014 21:32 Uhr
CHINA

China's Wachstum

China hat im Oktober ein Rekordtempo von 7,4 Prozent erreicht. Doch diese Zahlen sind fragwürdig. Die Statistik des Tsinghua-Instituts für Wirtschaftswissenschaften ist unser



19. November 2014 um 18:25 Uhr
DOR-DESIGN

Sandmännchen und Stasi-Mikrofone

Das größte Museum für DOR-Design steht ausgerechnet in Los Angeles. Ein Buch über das Wende Museum zeigt, welche Schätze und Abgründe es dort zu entdecken gibt.



10. November 2014 19:17 Uhr
ISRAEL

Keiner will von Intifada sprechen

Messerattacken auf Israelis. Krawalle auf dem Tempelberg. Schirmutzel im Gassengewirr.



10. November 2014 19:17 Uhr
ISRAEL

Keiner will von Intifada sprechen

Messerattacken auf Israelis. Krawalle auf dem Tempelberg. Schirmutzel im Gassengewirr.



19. November 2014 um 15:25 Uhr
AZEALIA BANKS

Klare Ansage aus Harlem

Erst galt Azealia Banks als großes Raptalent, dann als streitsüchtig und selbstverliebt. Ihr seit Jahren erwartetes Debut zeigt jetzt, wie gut das eine zum anderen passt.

+156%

BaQend

0,5s

1,3s

FRANKFURT

Thanks a lot!

gessert/ritter@informatik.uni-hamburg.de

baqend.com, orestes.info