ORESTES: Scalability and Low Latency for Polyglot Database Services

Felix Gessert, Norbert Ritter felix.gessert@gmail.com

About me

Since 2 years: PhD student (advisor: Norbert Ritter)
Since 1 year: CEO of Bagend



Outline



Motivation



Cache Sketch Approach



Polyglot Persistence Mediator



Wrap-up and future work

- Cloud Databases
- The Latency Problem of the web

Introduction: Which classes of cloud databases are there?















Motivation



Motivation













Redundance

Applications reimplement backend functionality





Redundance

Applications reimplement backend functionality



Database Dependence

Strong coupling between application and database system





Redundance

Applications reimplement backend functionality



Performance & Scalability Error-prone application-specific scaling policies



Database Dependence

Strong coupling between application and database system





Redundance

Applications reimplement backend functionality



Performance & Scalability Error-prone application-specific scaling policies



Database Dependence

Strong coupling between application and database system



Difficult Development Both client and server have to be implemented



With every **100ms** of additional page load time, revenue decreases by 1%.



Study by Amazon



When increasing load time of search results by **500ms**, traffic decreases by 20%.



Study by Google

Motivation

3-Tier Architecture Single Page Application



Study by Google













Automated Choice of Database System (Polyglot Persistence) Extensible High-Level Database *and* Backend **Abstractions**:

- User-Management
- Access Control
- Schema
- Transactions



Automated Choice of Database System (Polyglot Persistence) Extensible High-Level Database *and* Backend **Abstractions**:

- User-Management
- Access Control
- Schema
- Transactions

Outline



Motivation





Polyglot Persistence Mediator

1	

Wrap-up and future work

- REST/HTTP API
- HTTP Caching
- The Cache Sketch
 - Principle
 - Construction
 - Use



Ţ








Unified REST API

orestes : Orestes Methods

Show/Hide List Operations Expand Operations Raw





Unified REST API

Platform-specific interfaces map to unified REST API



 REST API leverages existing HTTP infrastructure Load-Balancer (*stateless communication REST constraint*) Web-Caches (*caching REST constraint*)

Expiration-based Caching



HTTP Caching Model:

- Expiration-based with defined TTL
- Revalidations check freshness at the server

Expiration-based Caching



Research Question:

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



Server/DB





































The Client Cache Sketch

Let c_t be the client Cache Sketch generated at time t, containing the key key_x of every record x that was written before it expired in all caches, i.e. every x for which holds:

 $\exists r(x, t_r, TTL), w(x, t_w) : t_r + TTL > t > t_w > t_r$



The Client Cache Sketch

Let c_t be the client Cache Sketch generated at time t, containing the key key_x of every record x that was written before it expired in all caches, i.e. every x for which holds:

 $\exists r(x, t_r, TTL), w(x, t_w) : t_r + TTL > t > t_w > t_r$



- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency

Solution: Cached Initialization

- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency





Cache





- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency







- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency









- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency









- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency









- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency













- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency











- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency





- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency











- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency





- Clients load the Cache Sketch at connection
- Every non-stale cached record can be reused without degraded consistency











Problem II: Slow CRUD performance

Solution: Δ-Bounded Staleness

- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity


- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity



- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity



- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity



- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity



- \circ Clients refresh the Cache Sketch so its age never exceeds Δ
- → *Consistency guarantee*: ∆-atomicity



- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts





- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts





- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts





- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts

	Begin Transaction			
	Bloom Filter Cache		REST-Server	
Client	Reads Cache Writes Cache	-2+	REST-Server	Writes (Hidden)
	Commit: read- & write-set versions Committed <i>OR</i> aborted + stale object	3) cts	REST-Server	



- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts



- Solution: Conflict-Avoidant Optimistic Transactions
 - Cache Sketch fetched with transaction begin
 - Cached reads \rightarrow Shorter transaction duration \rightarrow less aborts



The Server Cache Sketch

- Goal: Efficient Generation of Cache Sketch and Invalidation Minimization
- ► Counting Bloom Filter and key → expiration mapping



https://github.com/Baqend/Orestes-Bloomfilter

The Server Cache Sketch

- Goal: Efficient Generation of Cache Sketch and Invalidation Minimization
- ► Counting Bloom Filter and key → expiration mapping



https://github.com/Baqend/Orestes-Bloomfilter

TTL Estimation

Problem: if TTL << time to next write, then it is contained in Cache Sketch unnecessarily long
TTL Estimator: finds "best" TTL



TTL Estimation

Problem: if TTL << time to next write, then it is contained in Cache Sketch unnecessarily long
TTL Estimator: finds "best" TTL



TTL Estimation

Problem: if TTL << time to next write, then it is contained in Cache Sketch unnecessarily long
TTL Estimator: finds "best" TTL



YCSB Monte Carlo Caching Simulator (YMCA)

- Goal: Analysis of arbitrary caching architectures using the standard YCSB benchmark
 - **Metrics**: Latency, TP, Cache Hits, Stale Reads, Invalidations
 - Training of TTL Estimator: Hill Climber finds optimal params



Results: Simulation & real-world



Outline



Motivation



Cache Sketch Approach



Polyglot Persistence Mediator



Wrap-up and future work

- Idea
- Process
- Evaluation

- Fully transparent choice of DB, based on requirements (SLAs)
- 3-step-process:
 - 1. Requirements
 - 2. Resolution
 - 3. Mediation





- Fully transparent choice of DB, based on requirements (SLAs)
 - 3-step-process:
 - 1. Requirements
 - 2. Resolution
 - 3. Mediation





- Fully transparent choice of DB, based on requirements (SLAs)
 - 3-step-process:
 - 1. Requirements
 - 2. Resolution
 - 3. Mediation





- Fully transparent choice of DB, based on requirements (SLAs)
- 3-step-process:
 - 1. Requirements
 - 2. Resolution
 - 3. Mediation

	Annotation	Туре	Annotated at
	Read Availability	Continuous	*
	Write Availability	Continuous	*
	Read Latency	Continuous	*
	Write Latency	Continuous	*
/ 1.	Write Throughput	Continuous	*
/ sc	Data Vol. Scalability	Non-Functional	Field/Class/DB
1	Write Scalability	Non-Functional	Field/Class/DB
	Read Scalabilty	Non-Functional	Field/Class/DB
↓	Elasticity	Non-Functional	Field/Class/DB
Databas	Durability	Non-Functional	Field/Class/DB
/	Replicated	Non-Functional	Field/Class/DB
	Linearizability	Non-Functional	Field/Class
	Read-your-Writes	Non-Functional	Field/Class
Class	Causal Consistency	Non-Functional	Field/Class
	Writes follow reads	Non-Functional	Field/Class
/	Monotonic Read	Non-Functional	Field/Class
Field Fiel	Monotonic Write	Non-Functional	Field/Class
	Scans	Functional	Field
	Sorting	Functional	Field
annc annc	Range Queries	Functional	Field
Inhe	Point Lookups	Functional	Field
anno	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

- The Provider resolves the requirements
- RANK algorithm recursively analyzes schema annotations
- For each schema element:
 - Find each DB that satisfies all binary requirements
 - Pick the one that has the best score according to utility function of historic or predicted metrics





- The Provider resolves the requirements
- RANK algorithm recursively analyzes schema annotations
- For each schema element:
 - Find each DB that satisfies all binary requirements
 - Pick the one that has the best score according to utility function of historic or predicted metrics





- The Provider resolves the requirements
- RANK algorithm recursively analyzes schema annotations
- For each schema element:
 - Find each DB that satisfies all binary requirements
 - Pick the one that has the best score according to utility function of historic or predicted metrics





- The Provider resolves the requirements
- RANK algorithm recursively analyzes schema annotations
- For each schema element:
 - Find each DB that satisfies all binary requirements
 - Pick the one that has the best score according to utility function of historic or predicted metrics





PPM Step 3: Mediation

- The PPM routes data and operations to the chosen DBs
- In the *primary database* model, it triggers periodic materializations
- Metrics (latency, availability, etc.) are reported to the resolver





Evaluation

Scenario: news articles with impression counts Objective: low-latency top-k querys, high-throughput counts



Found Resolution

Evaluation

Scenario: news articles with impression counts Objective: low-latency top-k querys, high-throughput counts



Outline



Motivation



Cache Sketch Approach



Polyglot Persistence Mediator



Wrap-up and future work

- Summary
- Research areas
- Baqend

Summary



- Cache Sketch: dual approach to web caching for database services
 - Consistent (Δ-atomic) *expiration-based* caching
 - Invalidation-based caching with minimal purges



Polyglot Persistence Mediator: SLA-driven, fine-grained selection of database systems



Future & Current Work

Cache Sketch:

- Online learning of best TTL estimation
- Quantify COT properties
- *Query* result caching
- Extend YMCA to replication and sharding architectures

Polyglot Persistence:

- Common requirements/SLA "library"
- Implementation of Resolution
- Live Migration
- Scalable *metrics* aggregation











Future Work: Query Caching



Future Work: Query Caching



Future Work: Query Caching


Future Work: Query Caching



Future Work: Query Caching



Bacend Build faster Apps faster.

Baqend

Orestes as a startup, funded since March 2014



Thank you

Contact:

gessert@informatik.uni-hamburg.de
http://orestes.info
http://baqend.com

