

Active-Active ecosystem at **Linked**

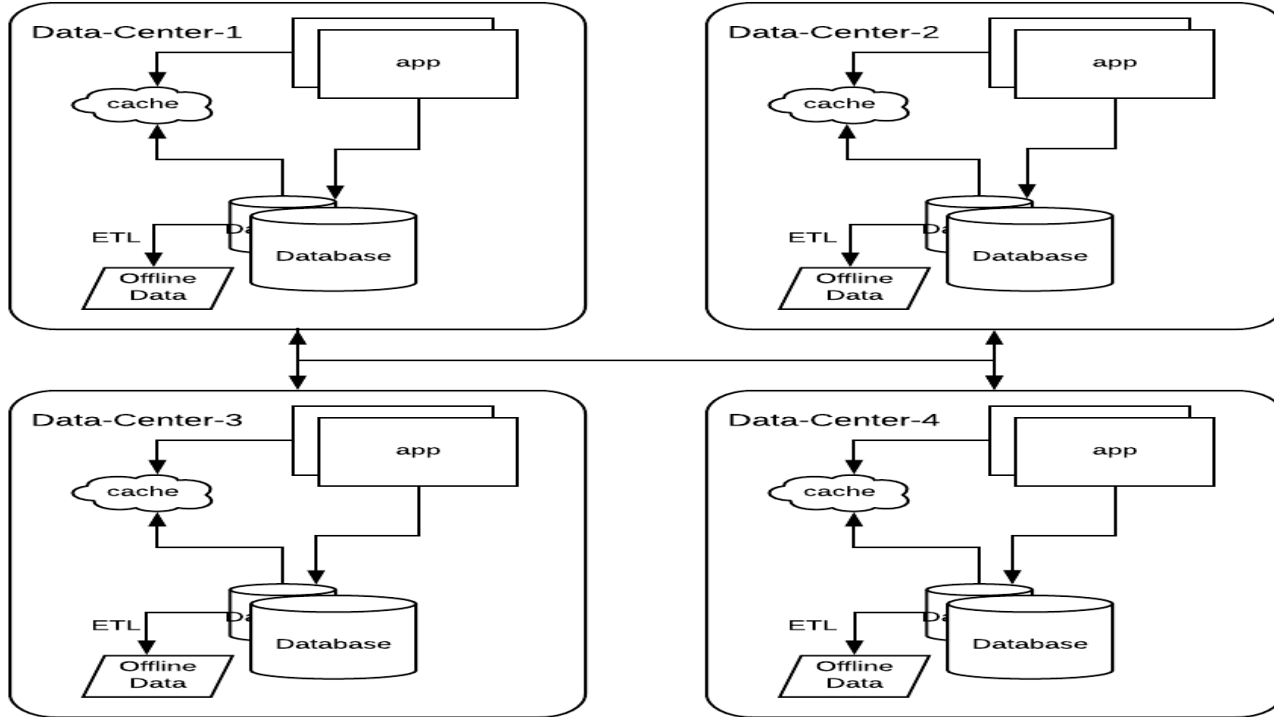
Janardh Bantupalli
Senior Staff Database Engineer

Sai Sundar
Director, Database Engineering

Agenda

- ❑ Oracle Active-Active Replication
- ❑ Incremental Capture and ETL: Data and Schema
- ❑ Reliability Infrastructure: Metrics and Monitoring
- ❑ Data Integrity Validation: Realtime Data Audit
- ❑ Q & A

Active-Active @ LinkedIn



Active-Active: Design

- ❑ Active-active configuration with DDL replication
- ❑ Columns: gg_modi_ts and gg_status, gg_priority (in few cases)
 - ❑ Before triggers to populate timestamp values
- ❑ Default LWW conflict resolution; priority resolution in few cases

Active-Active: Features

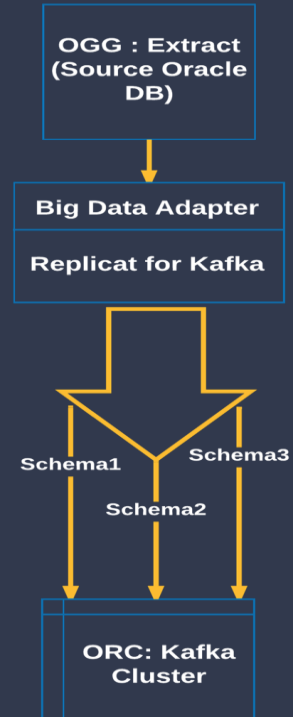
- ❑ Non-overlapping primary keys across DCs
- ❑ Hard deletes to soft delete conversion
 - ❑ Row re-birth prevention, cleanup of soft deletes
- ❑ Update-to-insert conversion; full row capture
- ❑ Scaling using parallelism; deadlock mitigation

Active-Active: Features (continued)

- ❑ OGG process offloading
- ❑ Data encryption and compression
- ❑ Parallel apply and eventual consistency
- ❑ Foreign key, unique key, novalidate constraint handling

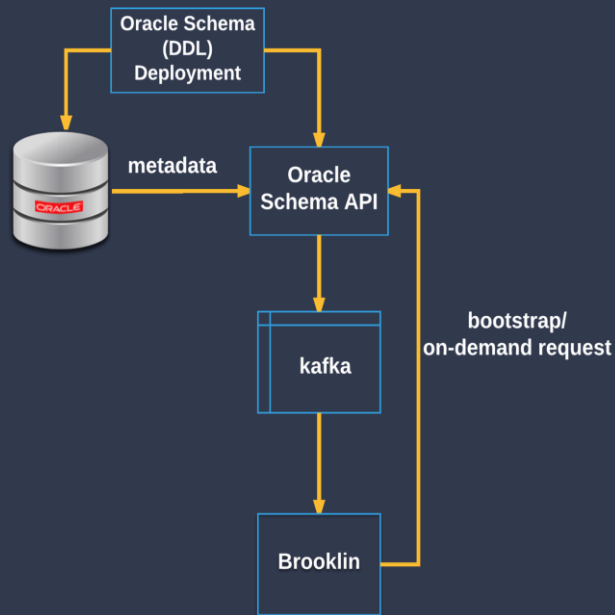
Incremental Data Pipeline

- ❑ OGG Big Data (Kafka) Adapter
- ❑ Pluggable kafka client
- ❑ Table-level granularity for source data
- ❑ Kafka topic per schema; partitioning for parallelism

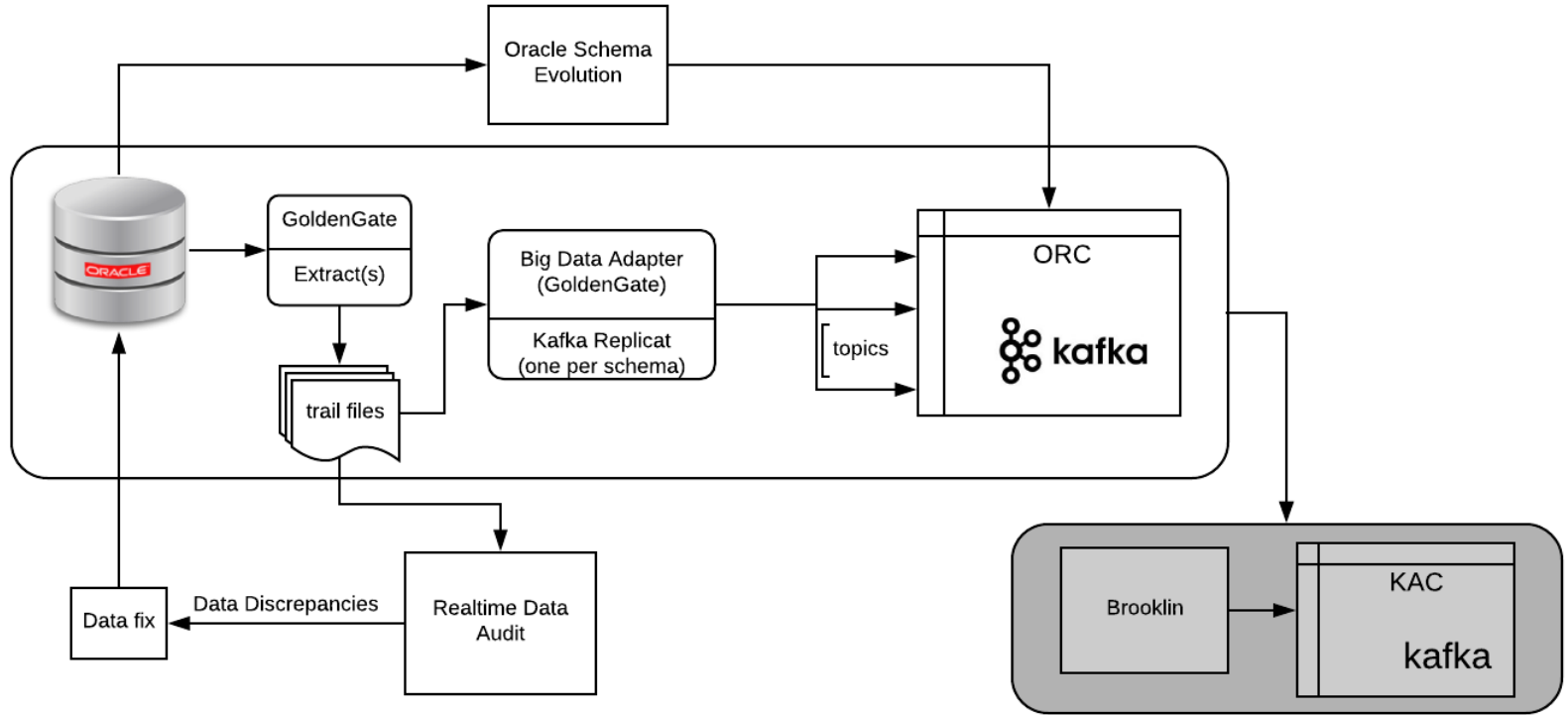


Data and Schema Propagation

- ❑ Generic payload schema; data and schema to different topics
- ❑ Independent framework to propagate schema changes (DDLs, metadata)
- ❑ Integration with release framework; API for on-demand invocation



Big Picture: Data Pipeline



Reliability Infrastructure

- ❑ Robust tooling within active-active, data pipeline components
 - ❑ Custom: Realtime Audit & APIs, OGG Monitor, Schema API
 - ❑ LinkedIn Stack: CFEngine, Ingraphs, Auto Alerts, Iris
- ❑ Redundancy (HA) in data audit and ETL pipelines; integration with schema deployment framework
- ❑ Proactive monitoring and detection of problematic load patterns; auto-recovery in most cases

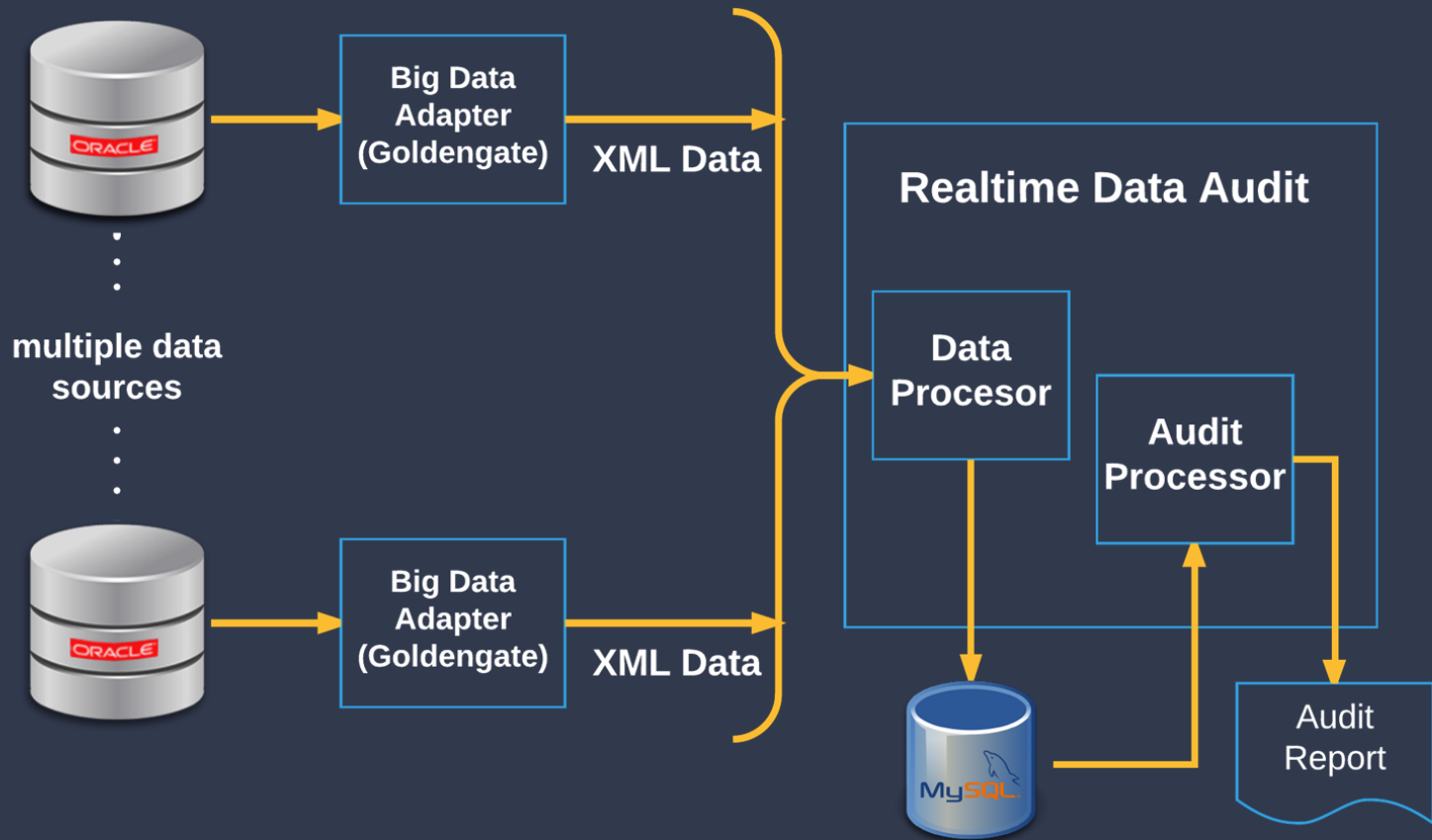


Reliability
/re-ly-a-bi-li-ti/

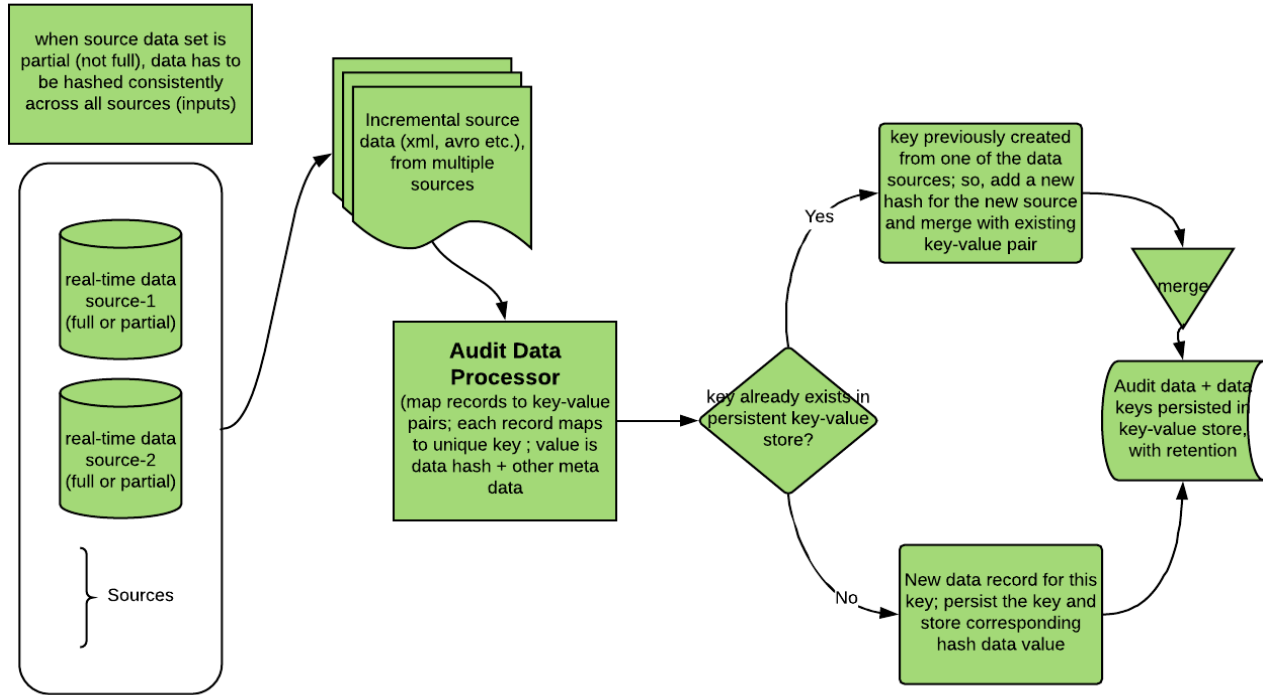
The ability of the software
to function under the given
conditions

Realtime Data Audit

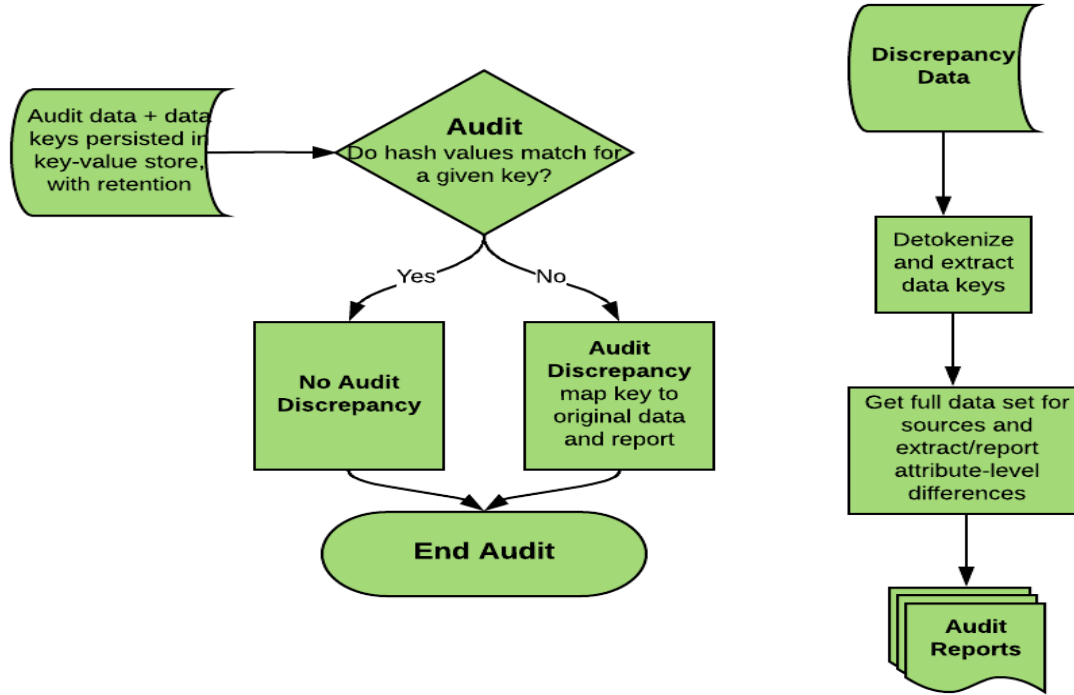
- ❑ Generic framework for data quality in distributed data environments
- ❑ Near-realtime data validation and discrepancy detection across data centers (data sources)
- ❑ Incremental data input from sources -- full (or) partial data set; audit based on last modified timestamp
- ❑ Provides refined discrepancy analysis if full data set is available



Data Processor



Audit Processor





Audit Framework : Features

- Deduplication and Checkpointing; recovery from failures and replay
- Data bucketing window customization to allow more deduplication
- Synchronization with replication and source data watermarks

Active-Active Ecosystem

OGG-Oracle
for multi
datacenter
active-active

OGG Big Data
Adapter for
Incremental
Data and ETL

Realtime Data
Audit
Framework for
Data Quality

Schema and
Metadata
Propagation
for ETL



Thank You!

<https://www.linkedin.com/in/janaonline>

<https://www.linkedin.com/in/ssaisundar>