# Data Governance for Regulated Industries Using Hadoop *...and NoSQL*

Justin Makeig, Director – Product Management, MarkLogic
November 2013

# Who am I?

- Product Manager for 6 years at MarkLogic
- Background in consulting and web development
- Passionate about data, infrastructure, and user experience

# What is MarkLogic?

- Enterprise NoSQL since 2001
- Distributed database + search + app platform
- 250+ paying customers, 500+ production applications

**MarkLogic**

# Agenda

- Data governance considerations
- Legacy approaches: Why it's hard
- New generation: Hadoop + Enterprise NoSQL
- Enterprise NoSQL
- Case studies: FATCA, eDiscovery, Dodd-Frank
- Q&A

# Data Governance Considerations

🔒 Security

**MarkLogic**®

# Data Governance Considerations

🔒 Security

👥 Privacy

# Data Governance Considerations

🔒 Security

👥 Privacy

🗄️ Provenance

**MarkLogic**

# Data Governance Considerations

🔒 **Security**     🛡️ Retention

👥 Privacy

🗄️ Provenance

**MarkLogic®**

# Data Governance Considerations

🔒 **Security**

🛡️ Retention

👥 Privacy

⏱️ Continuity

🗄️ Provenance

**MarkLogic**

# Data Governance Considerations

🔒 **Security**          🛡 Retention

👥 Privacy          ⏱ Continuity

🗄 Provenance          ⭐ Compliance

# Why is this difficult?

## And risky?

### And expensive?

#### And behind schedule?

Last Generation

"Unstructured"

ETL · Reference Data

Warehouse

ETL

OLTP

ETL

ETL · Archives

ETL · Data Marts

Documents, Messages · Video · Audio

Metadata

Social

Signals, Logs, Streams

ETL

ETL · Search

# Enter Hadoop



Updates

Queries

?

Hadoop

Aggregates, Models

Staging

Analytics

Persistence

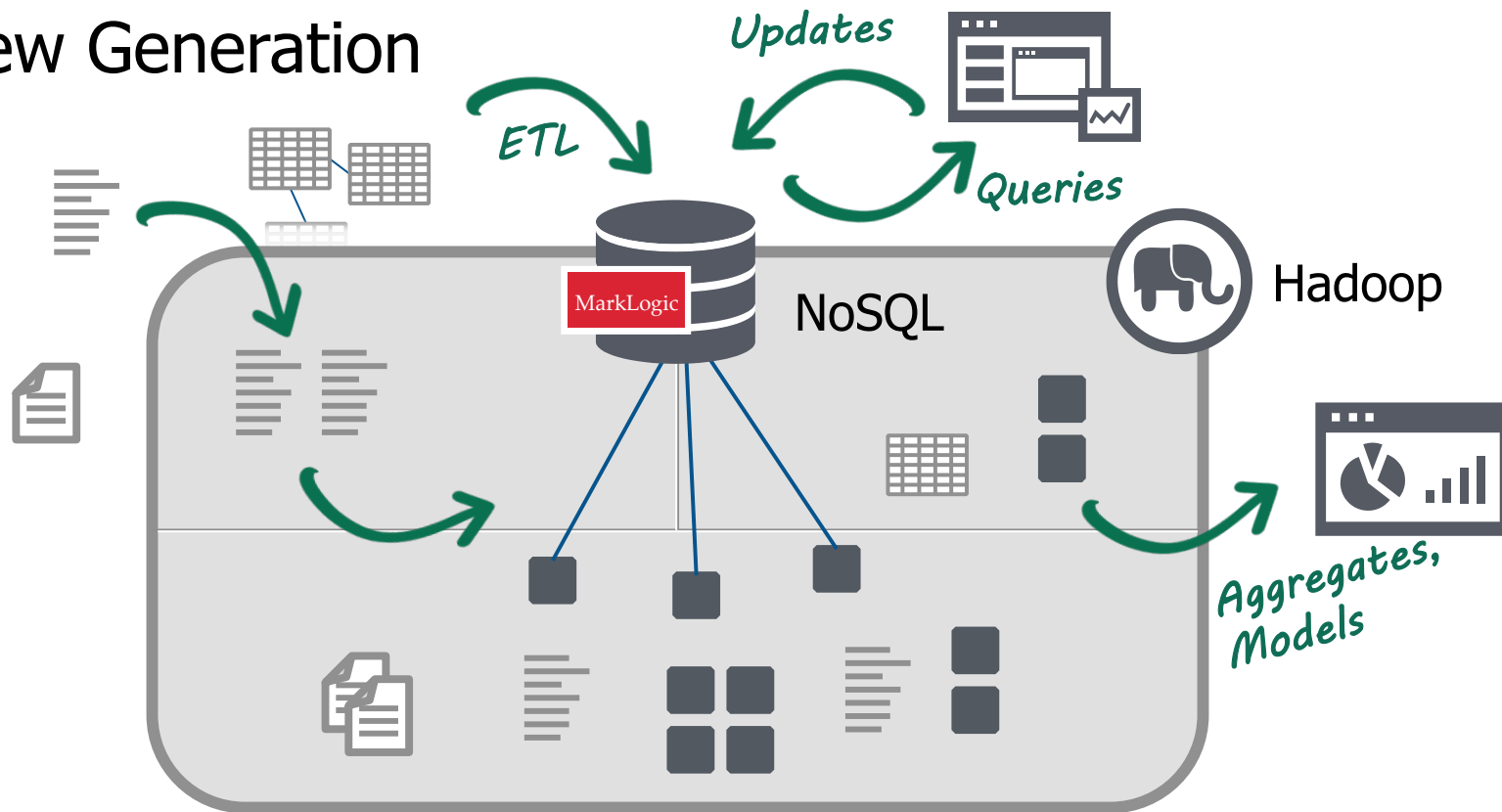**MarkLogic®**

# Why must we choose?

**Legacy RDBMS**

- Indexes
- Transactions
- Security
- Enterprise operations

**"NoSQL"**

- Flexible data model
- Commodity scale out
- Distributed, fault-tolerant
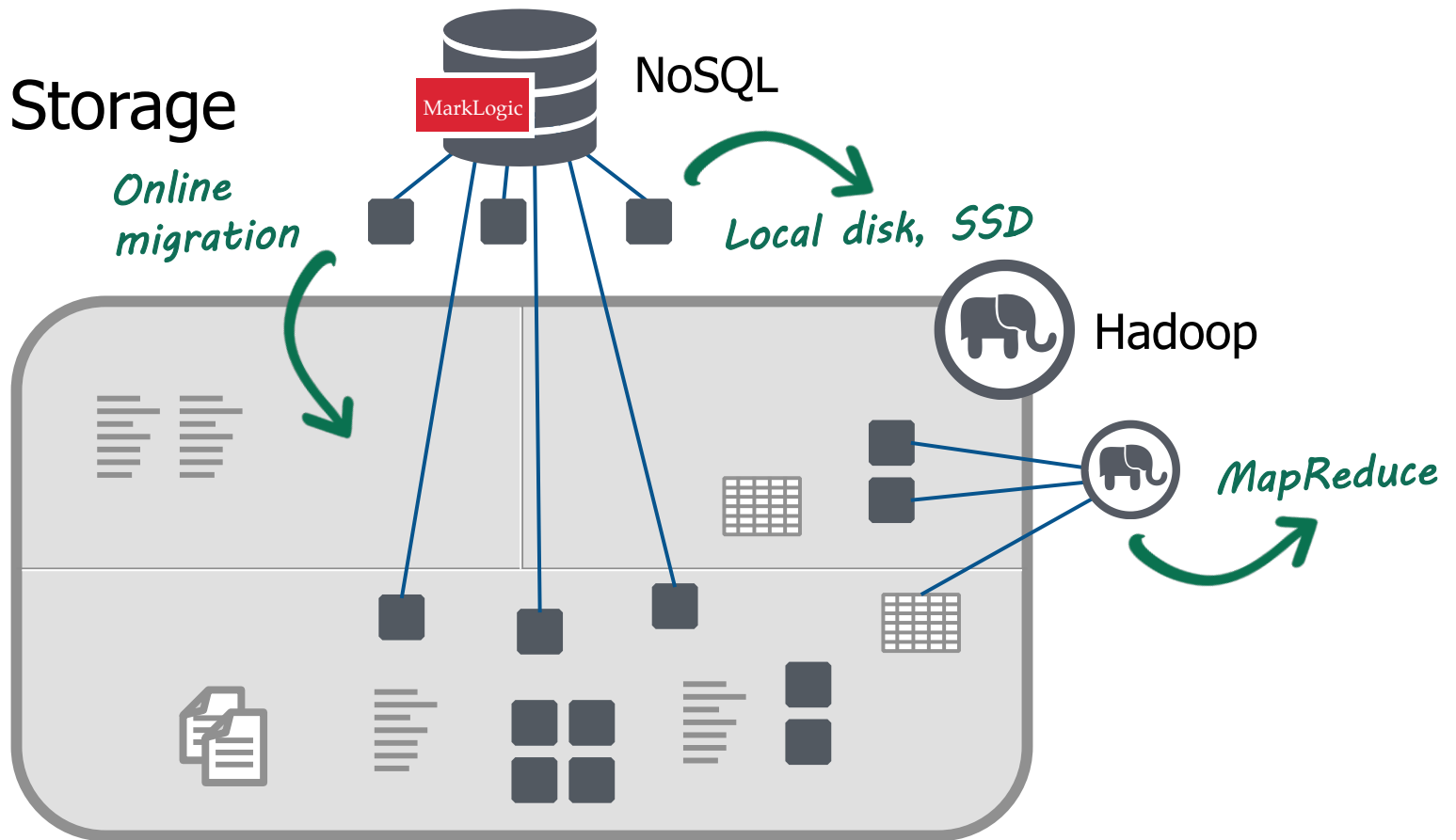- Hadoop sink/source

# New Generation

*ETL* *Updates* *Queries*

MarkLogic NoSQL Hadoop

*Aggregates, Models*

# Tiered Storage

**Active**
~$25/GB

**Historical**
~$1/GB

NoSQL

MarkLogic

*Online migration*

*Local disk, SSD*

Hadoop

*MapReduce*

**MarkLogic**

# Enterprise NoSQL

- Flexible data model, comprehensive indexes
    - o Documents: Hierarchy, text, values, tags—schema "on-demand"
    - o Scalars: Aggregates and range filters, including geospatial
    - o Triples: Linked facts and inferencing
    - o Permissions: Users, roles, compartments, and privileges
    - o Queries: Reverse indexes for alerting, matching
- Ad hoc dimensions, lock-free reads
- Real-time transformation
- Strict consistency, security throughout

# Preserving Context with Documents

Before

…movement of materials was observed en route to Abattabad some time after 14:30…

*Inline Enrichment*

After

…movement of materials was observed en route to
```
<place lat="…" long="…" version="2.2.1">
  <original>Abattabad</original>
  <canonical ref="…">Abbottabad</canonical>
  <source>/sources/1234</source>
  <confidence>0.87</confidence>
  …
</place>
```
some time after 14:30…

*Transactional updates*

**MarkLogic**

# Complementary Approaches

## NoSQL

- Online applications
- Delivery
- Decision-making
- Real-time
- Granular updates
- Distributed indexes

## Hadoop

- Offline analytics
- Staging
- Model-building
- Long-haul batch
- Write-once, read-many
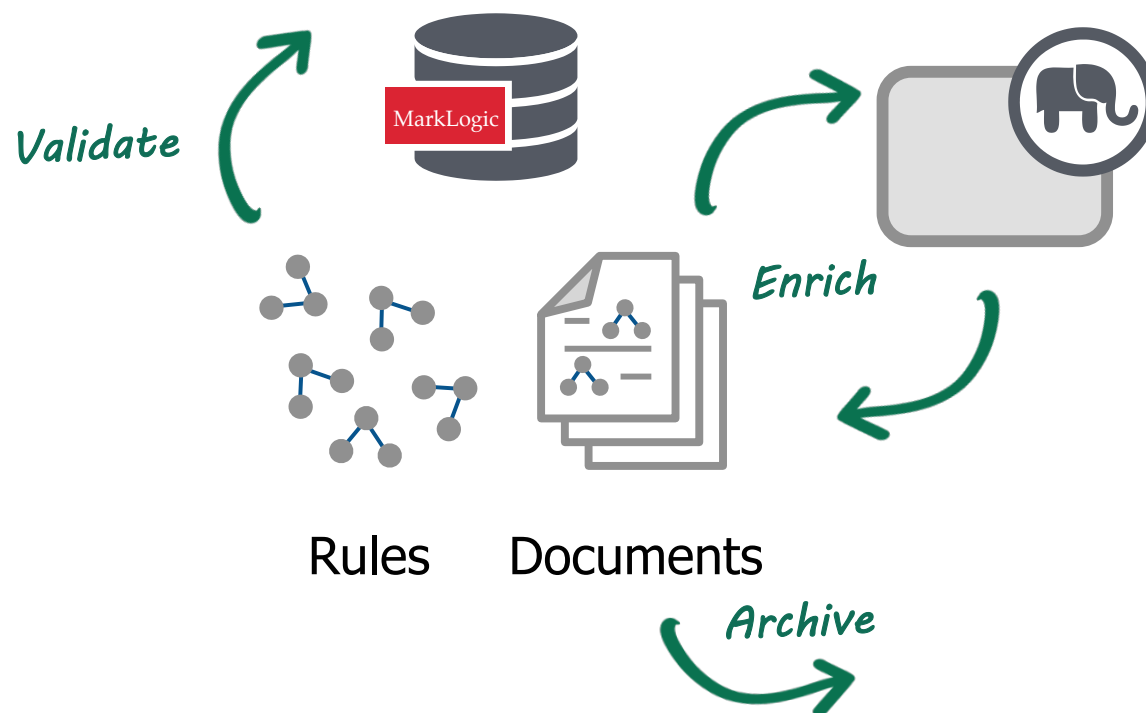- Distributed file system

# Case Studies

**MarkLogic**

# KPMG: FATCA Compliance for Customer On-Boarding

- Thousands of rules, 1–2M accounts, 30–40M documents
- Encoding, adjusting, and matching rules must scale
- Impossible to pre-define dimensions, relationships
- Vet new accounts and "show your work"
- Real-time decision-making

*6–48 hours to 3 seconds*

**KPMG**

**MarkLogic**

# KPMG: FATCA Compliance for Customer On-Boarding



Validate

MarkLogic

Enrich

Rules     Documents

Archive

# MarkLogic Semantics

- Combine facts from your documents with linked data

- Improve the search experience with facts about the results

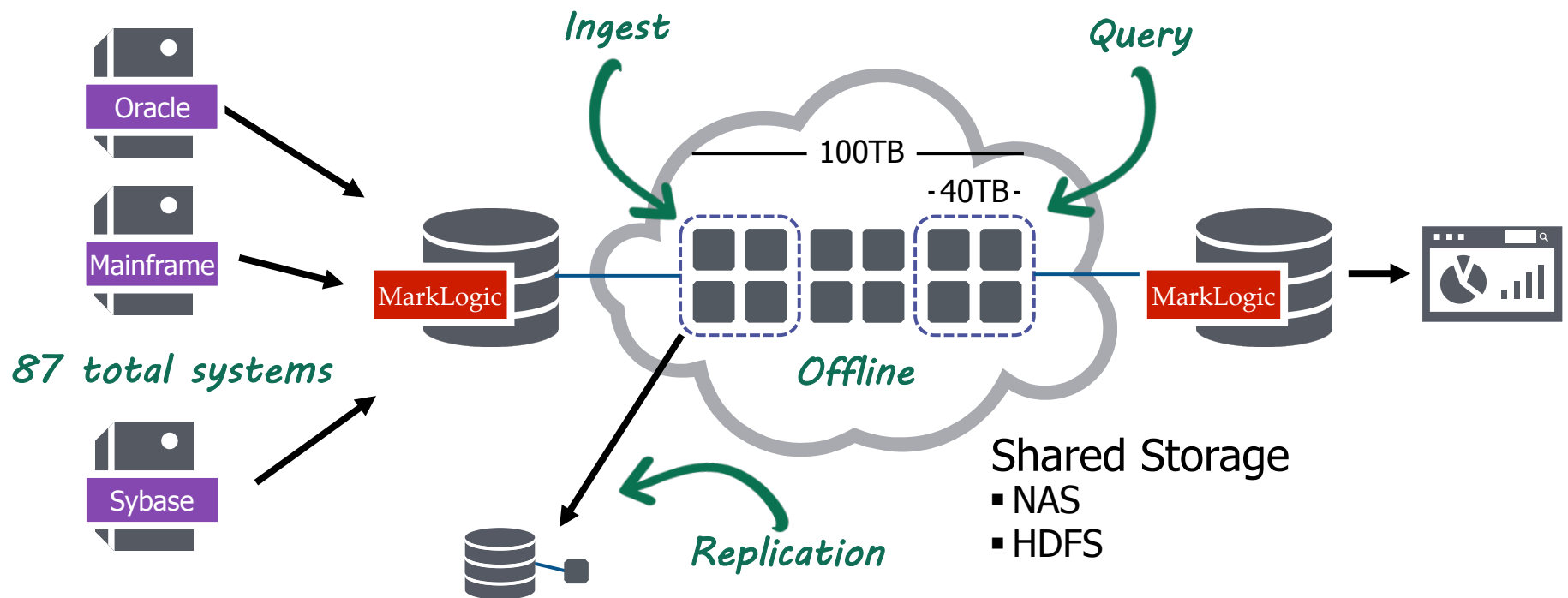- Enhance recall, improve precision with search *and* semantic queries

# Tier 1 European Bank: Compliance and Legal Holds

- Accurately respond to discovery as part of litigation
- Hold, review, produce data across current, legacy systems
- Repatriate and reconcile distributed data
- Demonstrate fidelity and audit trail
- Reduce infrastructure and maintenance costs
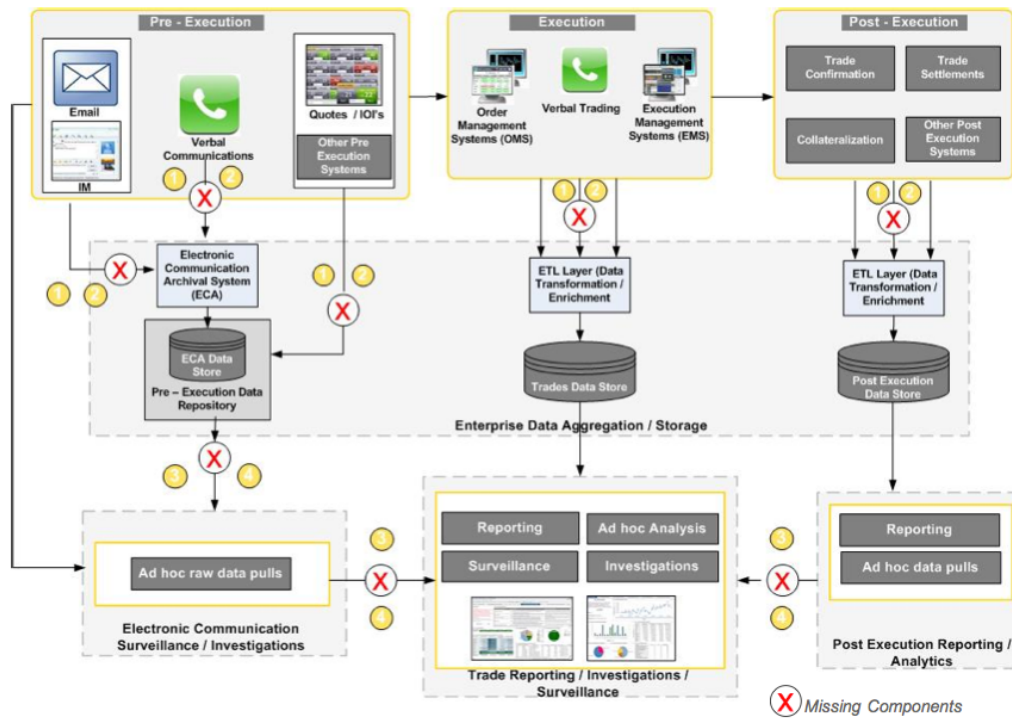
*Estimated $16M savings*

# Tier 1 European Bank: Compliance and Legal Holds



**Ingest**

**Query**

Oracle

Mainframe

*87 total systems*

Sybase

MarkLogic

100TB

-40TB-

*Offline*

MarkLogic

Shared Storage
- NAS
- HDFS

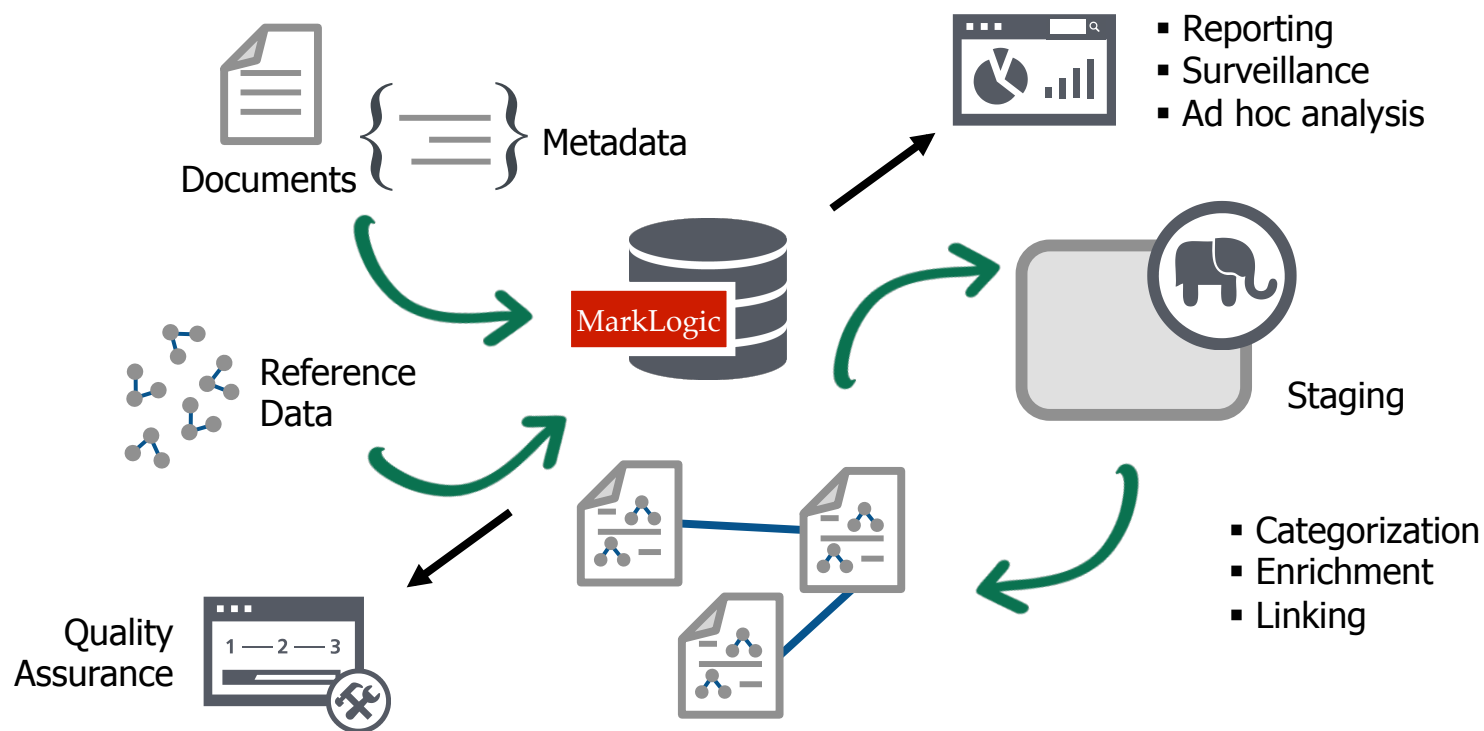*Replication*

# Ernst & Young: Dodd-Frank Compliance

- Trace lineage of order lifecycle for OTC derivatives
- Search, link supporting communications, documents
- Strict reporting and retention rules, response times
- Existing policies, point solutions don't scale

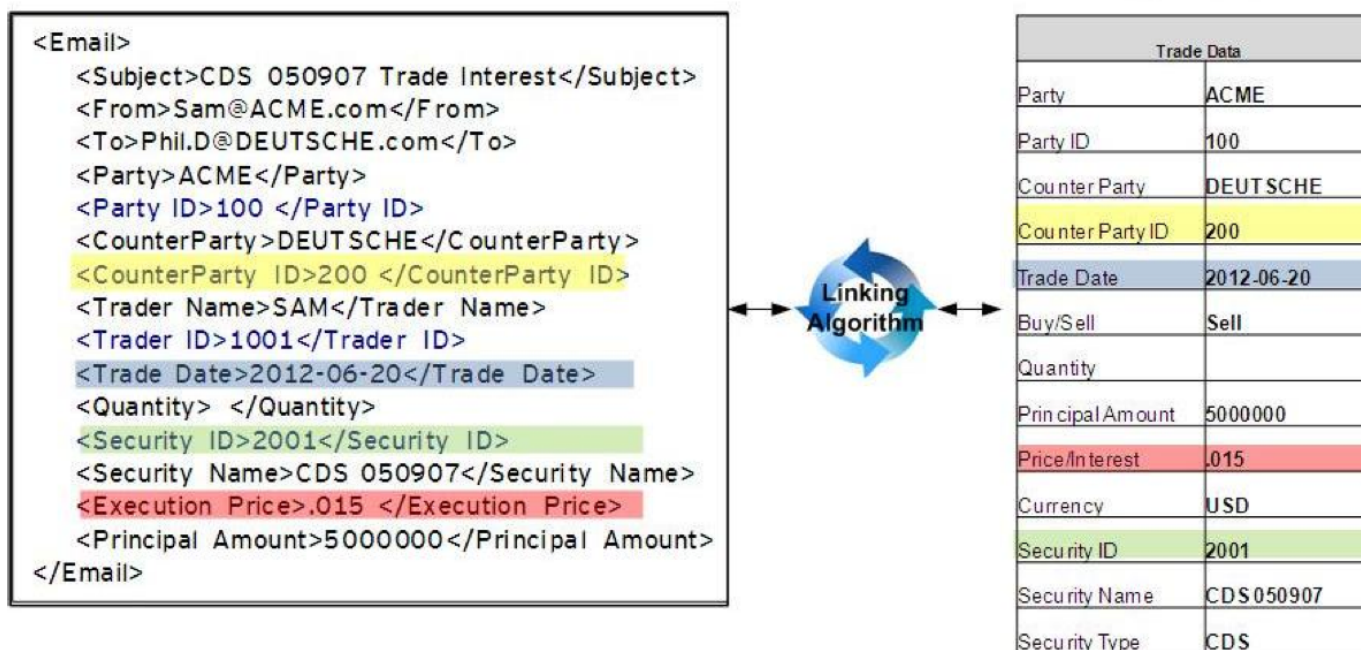**EY**
Building a better
working world

# Current State



- Missing key relationships between pre-/post-trade data
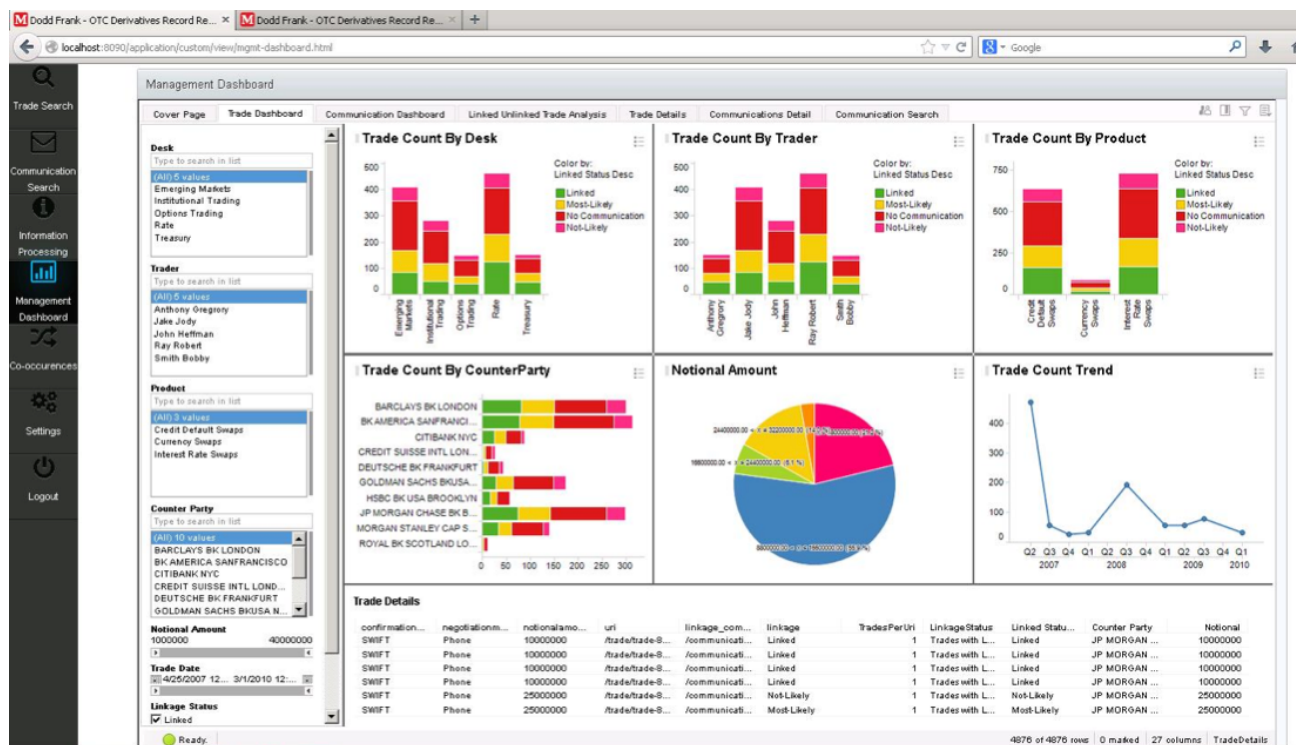- No way to query across silos
- Segregated reporting and surveillance

# Ernst & Young: Dodd-Frank Compliance

Documents

Metadata

- Reporting
- Surveillance
- Ad hoc analysis

MarkLogic

Reference Data

Staging

- Categorization
- Enrichment
- Linking

Quality Assurance

1 — 2 — 3

# Enrichment and Linking



Trade Record

# Management Dashboard

# What now?

# Take-Aways

- New and more data is both an opportunity and a threat
- Last generation of data management is not sufficient
- More copies, representations, transformations increase risk
- Index once and reuse across workloads, lifecycle
  - NoSQL: indexing and updates for interactive apps
  - Hadoop: staging, persistence, and analytics

# DO MORE WITH HADOOP

## SECURE
Minimize duplication, costly ETL, reduce risk

## REAL-TIME
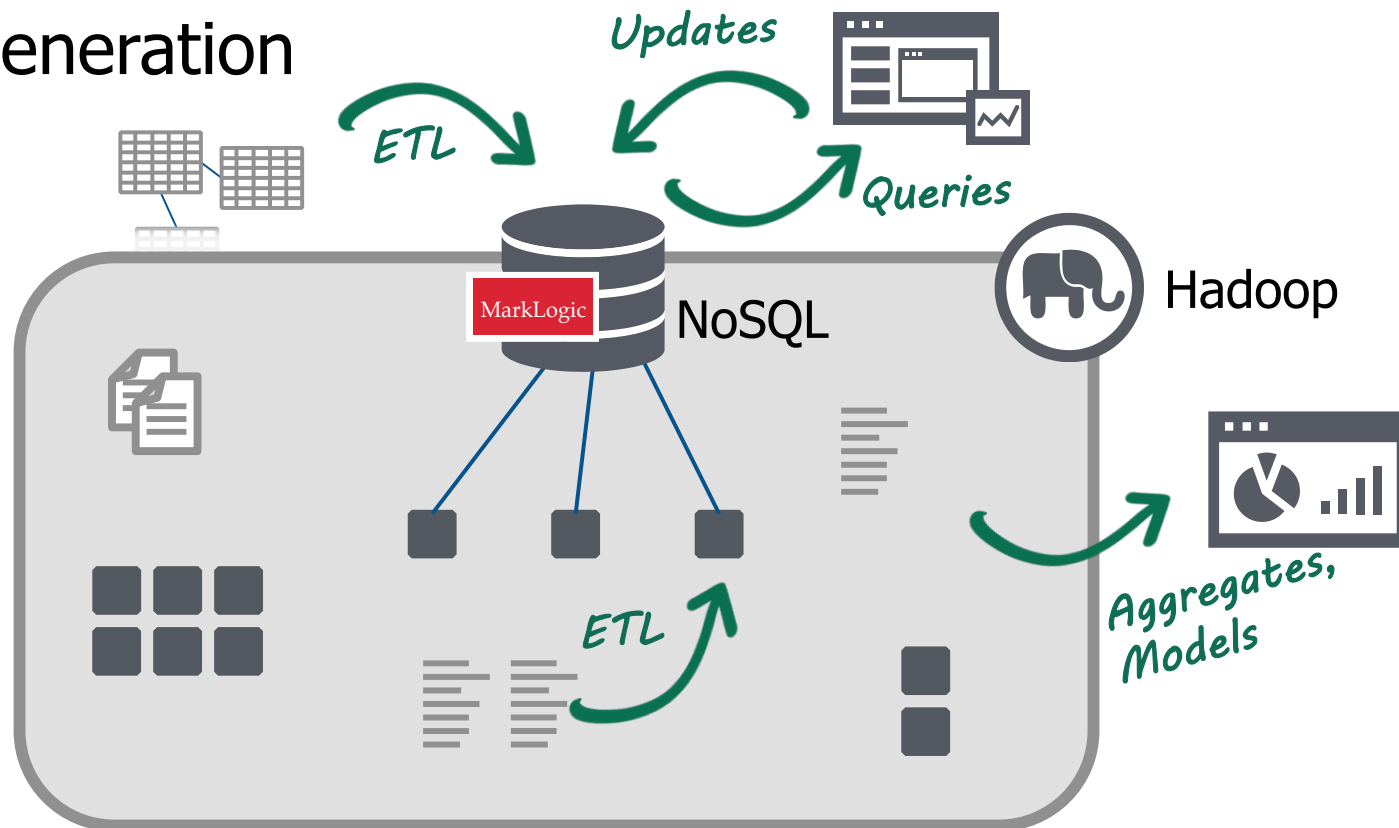Enterprise-class database for real-time search, delivery & analytics

## RUN APPLICATIONS
Run mission critical applications directly on HDFS

# New Generation

Updates

ETL

Queries

MarkLogic

NoSQL

Hadoop

ETL

Aggregates,
Models

# Preserving Context with Documents

- Hierarchy
- Relationships
- Semantics
- Security

Marshal

Shred