# Next Generation Archiving with Hadoop

*A look at archiving, e-discovery, and supervision on the Hadoop platform*

Jordan Volz
Systems Engineer @ **cloudera**®

# What are we talking about?

## Archiving

- Long-term storage of data
  - Compliance (WORM/NENR) vs non-compliance
  - Ingestion + Enrichment
  - Retention
  - Active vs passive
  - Reconciliation
  - Auditing
  - Search

## e-Discovery

- Review of electronic data to assess its relevance in legal proceedings
  - ECA
  - Legal Hold
  - TAR
  - Production
  - Case Management
  - Metadata management
  - EDRM

## Supervision

- Review of electronic communication to detect unethical conduct
  - Risk-based policies
  - Random Sampling
  - Surveillance
  - Auditing
  - Analytics + reporting
  - CO workflow
  - Lexicon management

# Why do we care about this?

**Businesses fail to meet SEC rules on e-mail archiving, risk fines, imprisonment**

**MORE FINES FOR DEUTSCHE BANK FOR RUSSIAN TRADES**

## Wells Fargo agrees to pay $5 Million Penalty for providing an altered document to SEC

**SEC Fines Deutsche Bank (DB) for Failing to Safeguard Material, Nonpublic Info Generated by Research Analysts**

**United States: FINRA Fines Twelve Firms For Recordkeeping Violations**

**Enforcement: Advisor Inflated AUM, Stole From Client, SEC Says**

**SEC fines Citigroup, Morgan Stanley over forex trading program**

FINRA fines Wedbush Securities $1 million on blue sheet failures; more

SPARK SUMMIT 2017

# Shortcomings of Traditional Systems

## Archiving

- Dated architecture
- Scalability
- Can't handle all data types
- Inflexible deployment
- Proprietary technology
- End-to-end vs piecemeal + siloed solutions
- Security/Encryption at Rest

## e-Discovery

- Best tools are not native to archive
- ECA scalability
- Lack of support for advanced media (audio/video)
- Lack of machine learning/advanced analytics
- Demanding SLAs for data export/production – hard to scale up for large cases/demand

## Supervision

- Rule-based
- Too many false positives
- Lack of machine learning/advanced analytics
- Focused on batch processing
- Difficult to customize
- Difficult integrations to archive for surveillance

# Today's Wisdom

Spark is awesome, but for complicated workflows you often need **more than Spark**.

*Luckily, it has great integrations in the Hadoop ecosystem.*

# Strengths of Hadoop/CDH

## Fast

- Fast SQL with **Impala**
- Advanced Analytics + ML via **Spark ML**
- **Kafka** for streaming data ingestion
- In-flight data processing via **Spark Streaming**
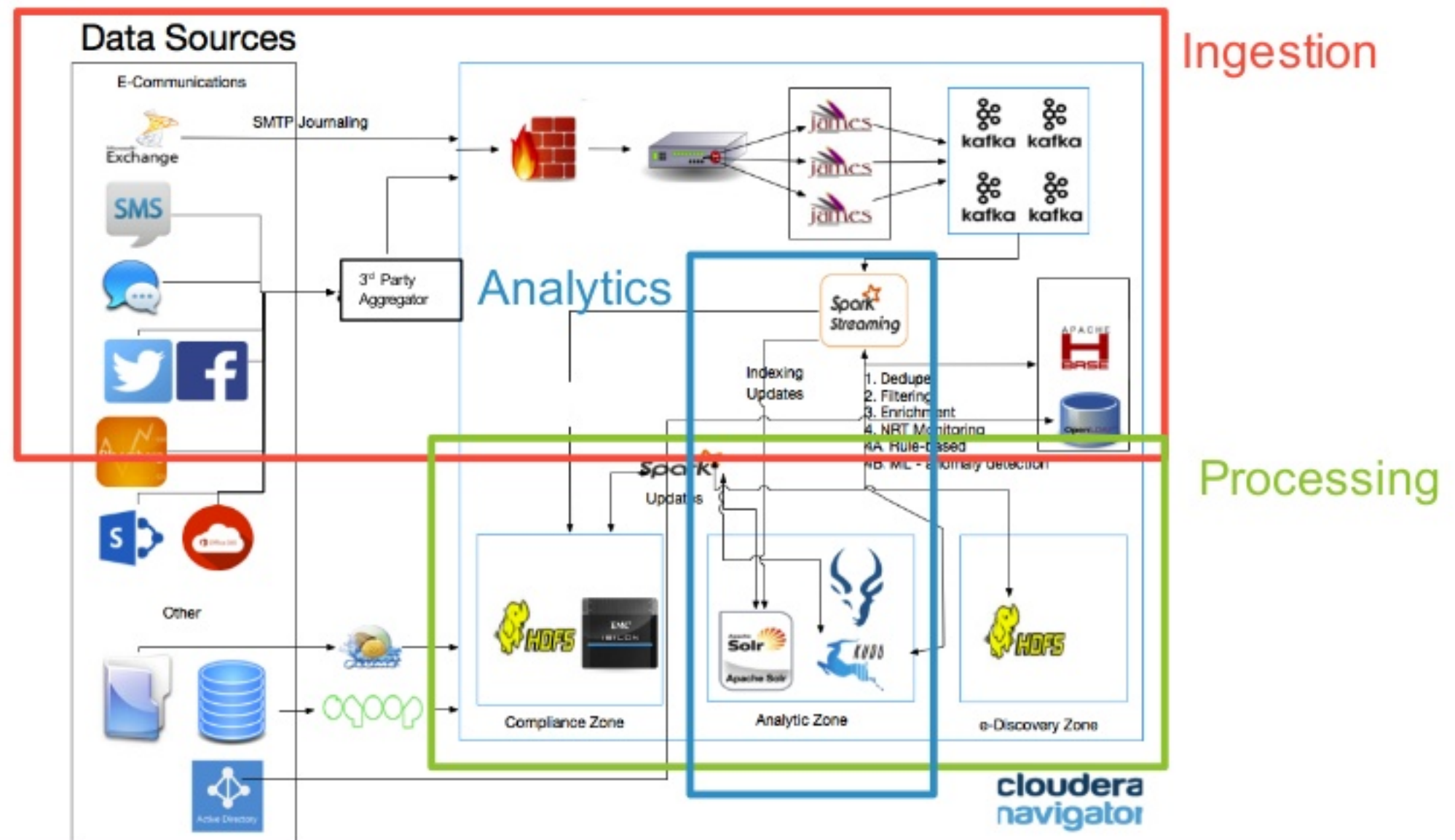- Distributed Search with **Solr**
- **Kudu** for tracking status/updates

## Easy

- Tested scalability to PBs
- Handles all data types
- Single platform for all solutions
- Easy integrations
- Driven by Open Source innovation
- Cloud or on-premise
- Integrated Data Science Platform (**Cloudera Data Science Workbench**)

## Secure

- End-to-end Security on a common platform
- Full system Data Encryption at Rest (**KMS/KTS**)
- Full system auditing with **Navigator**
- RBAC with **Sentry**
- Multi-tenancy
- **BDR** built-in (geo replication)

# Architectural Overview

# Architectural Overview - Ingestion



**Data Sources**

SMTP Journaling

**Apache James**

**LDAP**

**Hbase**

Identity Enrichment

Deduplication

**Kafka/ Flume**

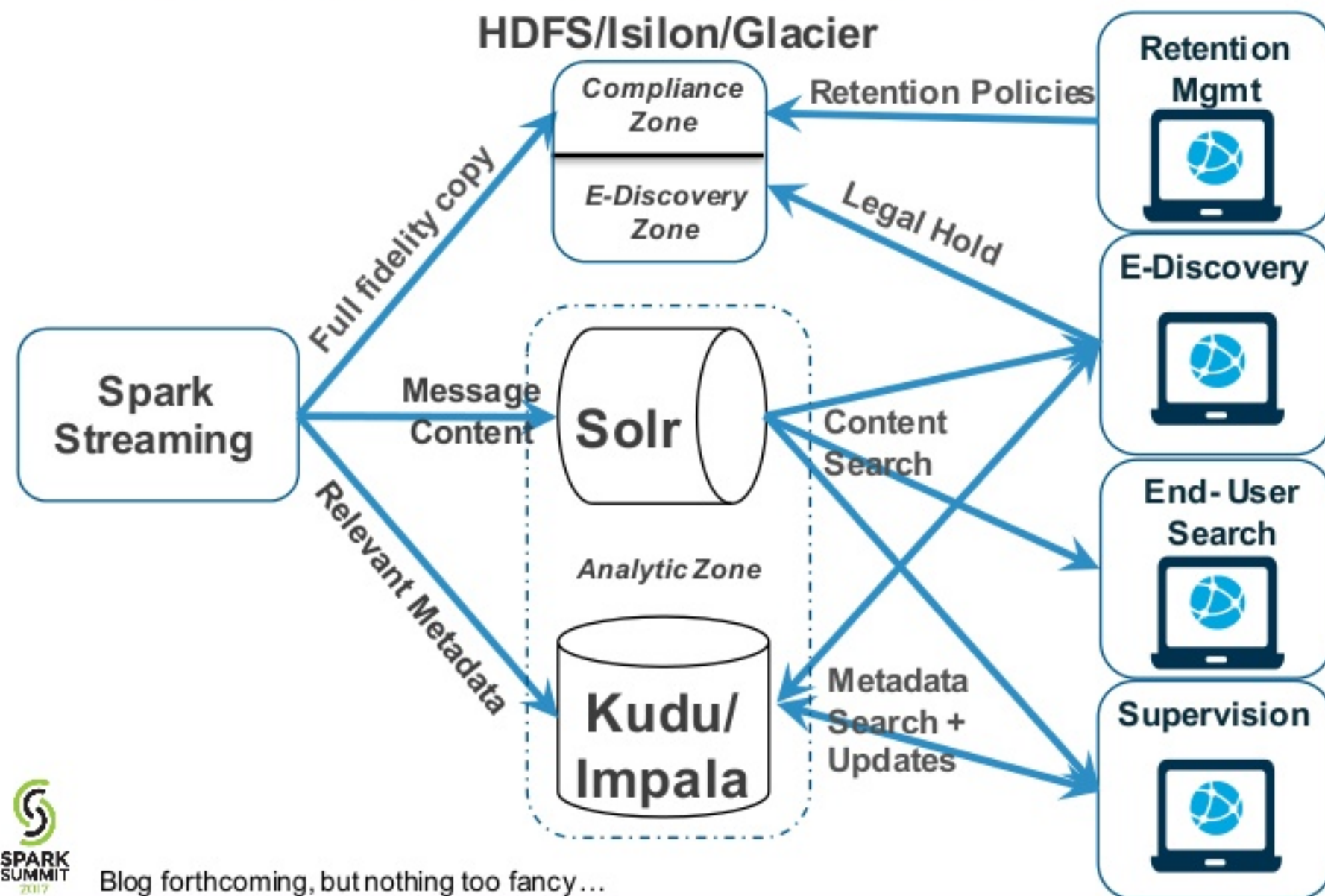**Spark Streaming**

Process

**Archive**

## Uses of Spark

- Deduplication (Hashing)
- Message Enrichment
- Reconciliation
- Message Filtering
- Policy Filtering

For more details, read the blog:
http://blog.cloudera.com/blog/2016/07/how-to-ingest-email-into-apache-hadoop-in-real-time-for-analysis/

SPARK SUMMIT 2017

# Architectural Overview - Processing

**HDFS/Isilon/Glacier**

**Retention Mgmt**

Compliance Zone

E-Discovery Zone

**Retention Policies**

**Legal Hold**

**Spark Streaming**

**Full fidelity copy**

**Message Content**

**Solr**

**Content Search**

*Analytic Zone*

**Relevant Metadata**

**Kudu/ Impala**

**Metadata Search + Updates**

**E-Discovery**

**End-User Search**

**Supervision**

## Uses of Spark

- Message compaction
- Message parsing and indexing
- Metadata extraction and storage
- Metadata updates (optional)
- Content re-indexing (optional)
- Content tokenization (optional)

Blog forthcoming, but nothing too fancy...

# Quick Sample Code

```scala
val messages = KafkaUtils.createDirectStream...

messages.foreachRDD{rdd =>

    val zkHost = "<zookeeper_host>:2181/solr"
    val req = new UpdateRequest
    req.setParam("collection", "<collection_name>")

    val kuduMaster = "<kudu_host>:7051"
    val kuduContext = new KuduContext(kuduMaster)
    val kuduTableName = "<kudu_table_name>"

    rdd.foreachPartition{partitionOfEmails =>
        //compact emails and save to HDFS --> covered in blog
        //establish solr connection
        val solrServer = new CloudSolrServer(zkHost)
        solrServer.setDefaultCollection(collection)
        solrServer.connect()
        val batch = new util.ArrayList[SolrInputDocument]()
        //establish kudu session
        val kuduClient = kuduContext.syncClient
        val table = kuduClient.openTable(kuduTableName)
        val kuduSession = kuduClient.newSession()
        kuduSession.setFlushMode(FlushMode.AUTO_FLUSH_BACKGROUND)

        val properties=System.getProperties()
        properties.setProperty("mail.smtp.host","cloudera.com")
```

```scala
        partitionOfEmails.foreach{ email =>
            val session=Session.getDefaultInstance(properties)
            val is = new ByteArrayInputStream(email.getBytes())
            val message=new MimeMessage(session,is)
            val id=message.getMessageID()
                ... //parse out email as needed
            val content=processContent(message)


            //solr indexing
            val doc = new SolrInputDocument
            doc.addField("message_id", id)
                ... //add desired fields
            doc.addField("body", content)
            batch.add(doc)

            //kudu upserting
            val operation: Operation = table.newUpsert()
            val row = operation.getRow()
            row.addString("id", id)
                ... //add desired fields
            kuduSession.apply(operation)
        }
        //commit stuff
        req.add(batch)
        solrServer.request(req)
        solrServer.shutdown()
        kuduSession.flush()
        kuduSession.close()
    }
}
```

Common Global Settings

Parsing/Processing

Connection Settings

Commits

# Architectural Overview – Analysis + ML

Electronic Communication

Solr

Spark ML

Reference Data (Transactional Data)

Kudu

Clustering Classification

Unified Surveillance

Impala

E-Discovery

Supervision

## Uses of Spark

- Machine Learning – Clustering for real-time anomaly detection in e-comm
- Model Permanence + Updates
- TAR for e-discovery (automated or recommended tagging + classification)
- Access to unified surveillance measures

Let's look at a quick example

# Surveillance Examples

- With only two tables, ecomm metadata and transactional data, but we can still start to answer some pretty interesting questions like…
  - What communication is occurring after a large transaction?
  - What trades are made after a flagged message?
  - Etc…

# Sample Schemas

## Table: ecomms_msgs

Version:

State:

### Schema

| Column | ID | Type |
|---|---|---|
| id | 0 | string |
| sentdate | 1 | int64 |
| receiveddate | 2 | string |
| archivedate | 3 | doub |
| subject | 4 | string |
| sender | 5 | string |
| recipient | 6 | string |
| sendergroup | 7 | string |
| recipientgroup | 8 | string |
| carboncopy | 9 | string |
| blindcarboncopy | 10 | string |
| legalhold | 11 | bool |
| legalholdpolicy | 12 | int32 |
| flagged | 13 | bool NULLABLE |
| flagType | 14 | int32 NULLABLE |
| supervisionpolicy | 15 | int32 NULLABLE |
| retentiontype | 16 | int32 NULLABLE |
| retentionperiod | 17 | int32 NULLABLE |

## Table: ecomms.trades_kudu (0db56ed82cef490088eb6beaf586e5f6)

| Version: | 0 |
|---|---|
| State: | Running |

### Schema

| Column | ID | Type | Encoding | Compression | Re |
|---|---|---|---|---|---|
| msgseqnum | 0 | int32 NOT NULL | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| trade_time | 1 | int64 NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| trader | 2 | string NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| price | 3 | double NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| total_price | 4 | double NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| msgtype | 5 | int32 NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| sourceseqnum | 6 | int32 NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| symbol1 | 7 | string NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| volume | 8 | int64 NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| exchangeid | 9 | string NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| securitytype | 10 | string NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |
| linkid | 11 | string NULLABLE | AUTO_ENCODING | DEFAULT_COMPRESSION | - |

# Large Transaction Example

# Flagged Messages Example

Impala · ↺ Add a name… · Add a description…

```
1  select trader, symbol1, total_price, trade_time,id, recipient, subject
2      from
3          ecomms.msgs, ecomms.trades_kudu
4      where flagged=true
5          and sender=trader
6          and sentdate between trade_time and trade_time + 600
```

1.43s  default ⌄

| | trader | symbol1 | total_price | trade_time | id | recipient |
|---|---|---|---|---|---|---|
| 1 | sally.beck@enron.com | HPE | 1410 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 2 | sally.beck@enron.com | BIG | 396.1 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 3 | sally.beck@enron.com | KMT | 2708 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 4 | sally.beck@enron.com | AXP | 7100 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 5 | sally.beck@enron.com | G | 5140 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 6 | sally.beck@enron.com | UPS | 20594 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 7 | sally.beck@enron.com | S | 1412 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 8 | sally.beck@enron.com | UNP | 7816 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 9 | sally.beck@enron.com | SU | 2617 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |
| 10 | sally.beck@enron.com | FSL | 2596.24 | 1449480641 | <1000004222.20000101120000@Email4222-3_4_2001 10_21_00 PM.eml> | john.forney@enron.comgeir.solberg@enron.comjohn.hodg |

# Retrieving Content Example

# Putting it all together…

# Some Obstacles

- Small File Problem
- Compliance Archiving (WORM storage)
- Large Files with Kafka
- UIs / workflow management

# Finding a Partner

- BI/Analytics/Reporting/Visualizations
- Compliance Storage
- Machine Learning/Data Science
- ETL/Ingest
- NLP/Sentiment Analysis

# Questions? Thank You.

jordan.volz@cloudera.com

linkedin.com/in/jordanvolz

SPARK SUMMIT 2017