# FUSING APACHE SPARK AND LUCENE FOR NEAR-REALTIME PREDICTIVE MODEL BUILDING

Debasish Das
Principal Engineer
Verizon

Pramod Lakshmi Narasimha
Principal Engineer
Verizon

**Contributors**
**Platform**: Pankaj Rastogi, Venkat Chunduru, Ponrama Jegan, Masoud Tavazoei
**Algorithm**: Santanu Das, Debasish Das (Dave)
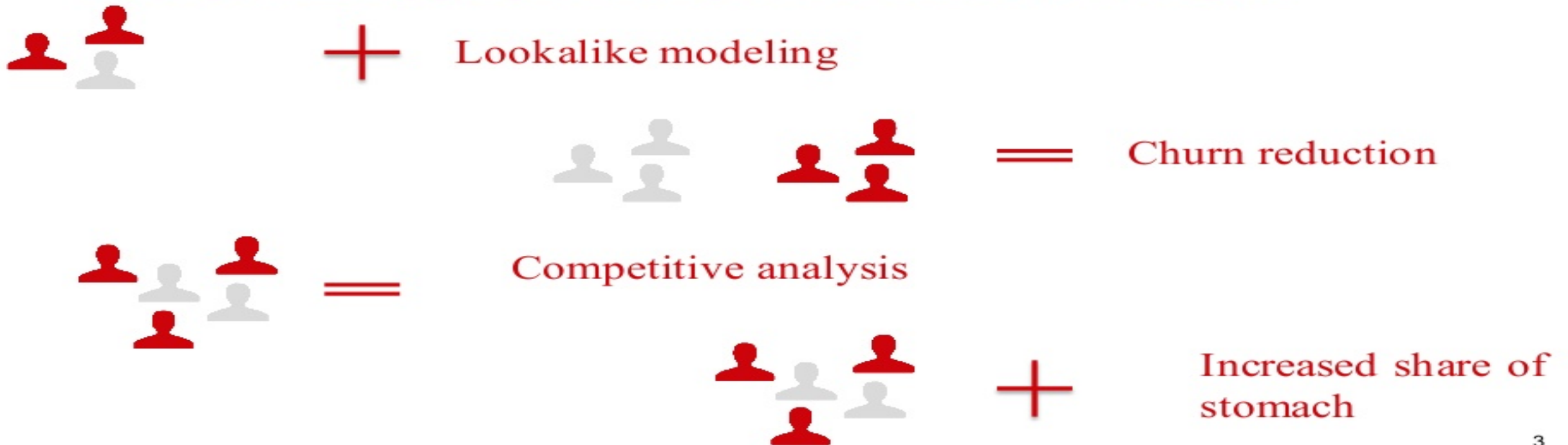**Frontend**: Altaff Shaik, Jon Leonhardt

verizon✓

# Data Overview

- Location data
  - Each srcIp defined as unique row key
  - Provides approximate location of each key
  - Timeseries containing latitude, longitude, error bound, duration, timezone for each key
- Clickstream data
  - Contains clickstream data of each row key
  - Contains startTime, duration, httphost, httpuri, upload/download bytes, httpmethod
  - Compatible with IPFIX/Netflow formats

# Marketing Analytics

- **Anonymous aggregate analysis for customer insights**


Lookalike modeling

Churn reduction

Competitive analysis

Increased share of stomach

# Data Model

- Dense dimension, dense measure
  Schema: srcip, date, hour, tld, zip, tldvisits, zipvisits
  Data: 10.1.13.120, d1, H2, macys.com, 94555, 2, 4

- Sparse dimension, dense measure
  Schema: srcip, date, tld, zip, clickstreamvisits, zipvisits
  Data: 10.1.13.120, d1, {macys.com, kohls.com}, {94555, 94301}, 10, 15

- Sparse dimension, sparse measure
  Schema: srcip, date, tld, zip, tldvisits, zipvisits
  Data: 10.1.13.120, d1, {macys.com, kohls.com}, {94555, 94301}, {macys.com:4, kohls.com:6}, {94555:8, 94301:7}
  Schema: srcip, week, tld, zip, tldvisits, zipvisits
  Data: 10.1.13.120, week1, {macys.com, kohls.com}, {94555, 94301}, {macys.com:4, kohls.com:6}, {94555:8, 94301:7}

  - Sparse dimension, sparse measure, **last N days**
    Schema: srcip, tld, zip, tldvisits, zipvisits
    Data: 10.1.13.120, {macys.com, kohls.com}, {94555, 94301}, {macys.com:4, kohls.com:6}, {94555:8, 94301:7}
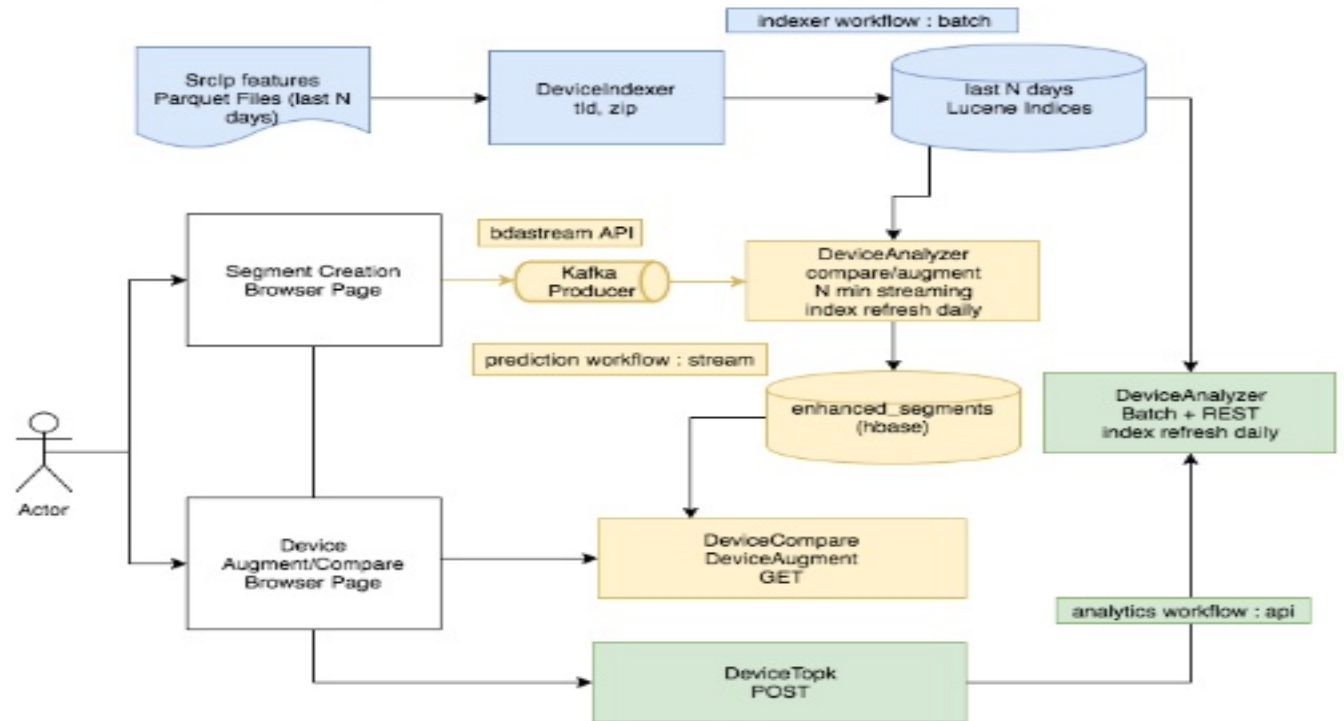
- Competing technologies: PowerDrill, Druid, LinkedIn Pinot, EssBase

4

# Document Dataset Representation

- Example
  Schema: srcip, tld, zip, tldvisits, zipvisits
  Data: 10.1.13.120, {macys.com, kohls.com}, {94555, 94301}, {macys.com:4, kohls.com:6}, {94555:8, 94301:7}

- DataFrame row to Lucene Document mapping

| Store/schema | Row | Document |
|---|---|---|
| srcip | primary key | docId |
| tld<br>zip | String<br>Array[String] | SingleValue/MultiValue<br>Indexed Fields |
| tldvisits<br>zipvisits | Double<br>Map[String, Double] | SparseVector<br>StoredField |

- Distributed collection of srcIp as RDD[Document]
  - ~100M srcip, 1M+ terms (sparse dimensions)
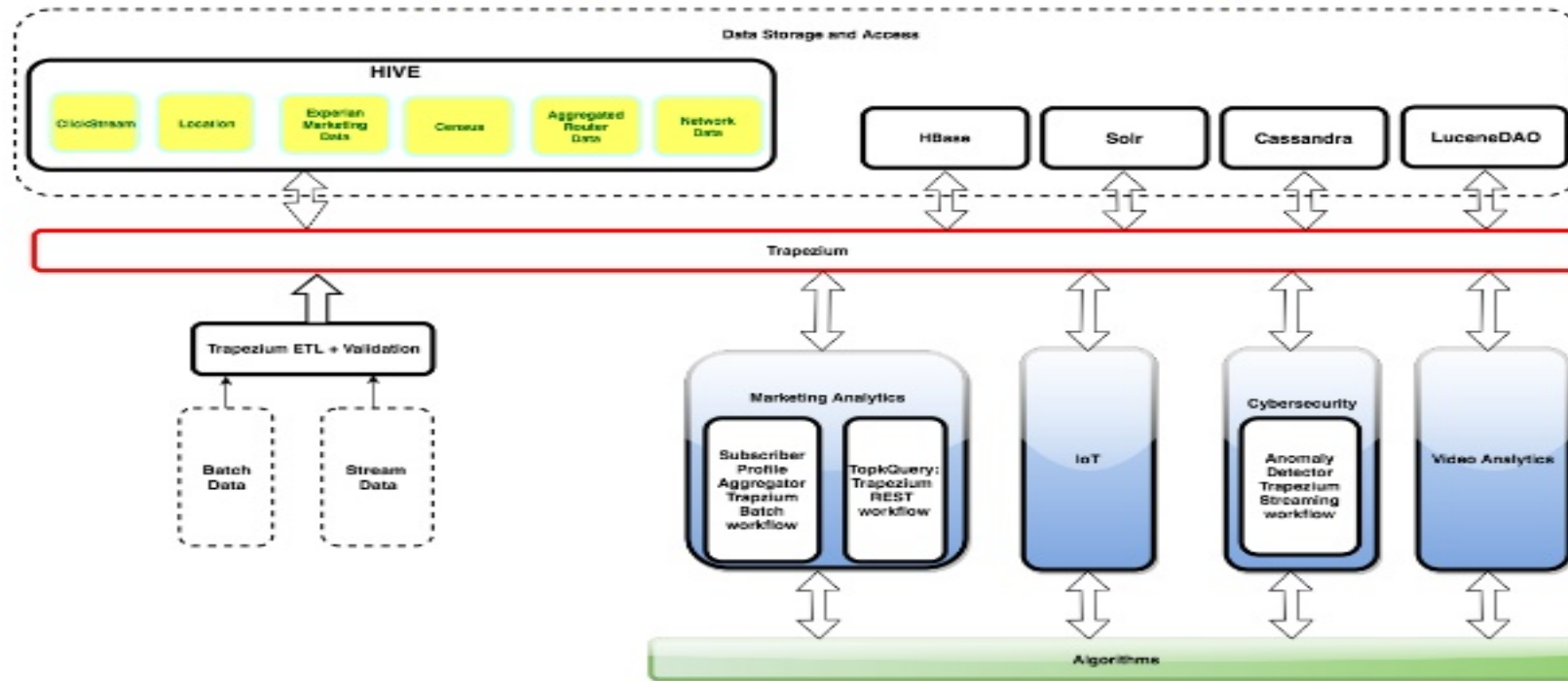
# DeviceAnalyzer

- DeviceAnalyzer goals
  - Search and retrieve devices that matched query
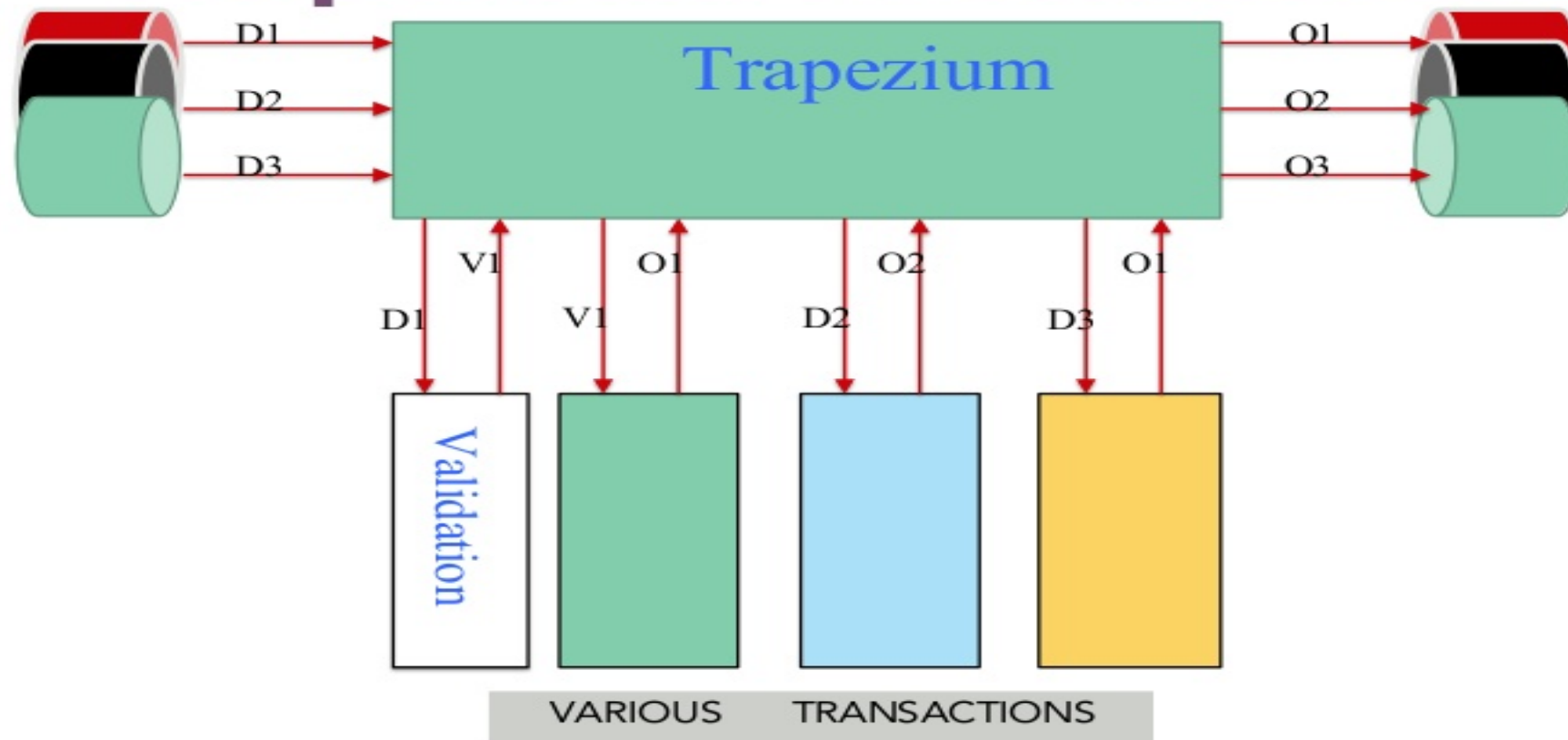  - Generate statistical and predictive models on retrieved devices

# What is Trapezium ?

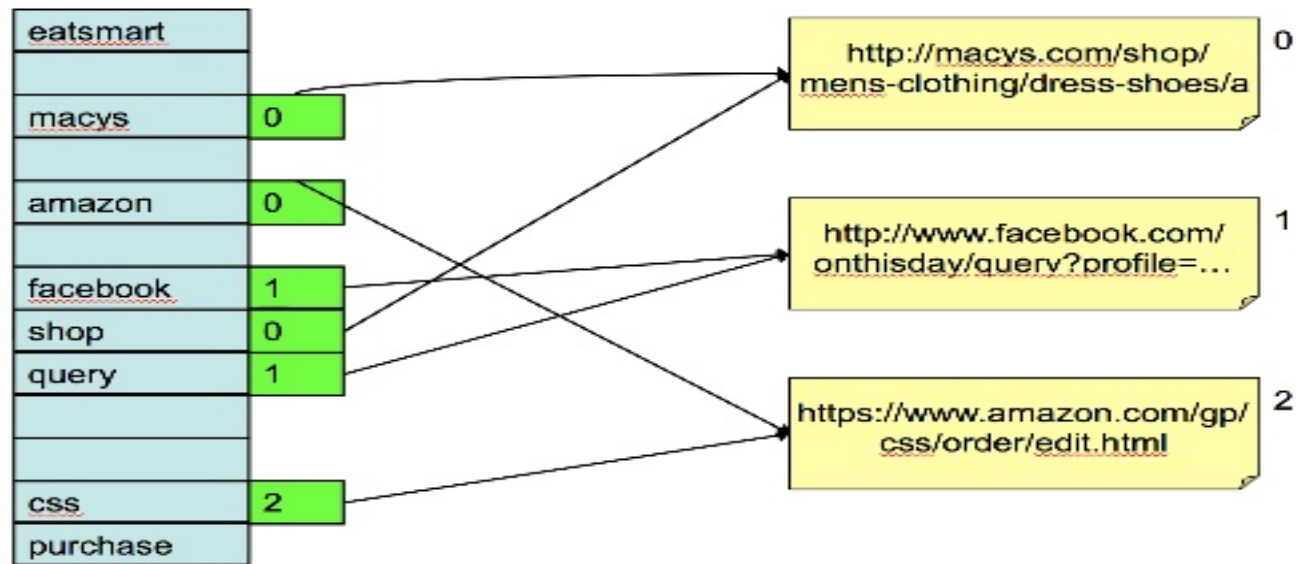DAIS Open Source framework to build batch, streaming and API services

https://github.com/Verizon/trapezium

7

# Trapezium Architecture

# Lucene Overview

- Scalable, full-text search library
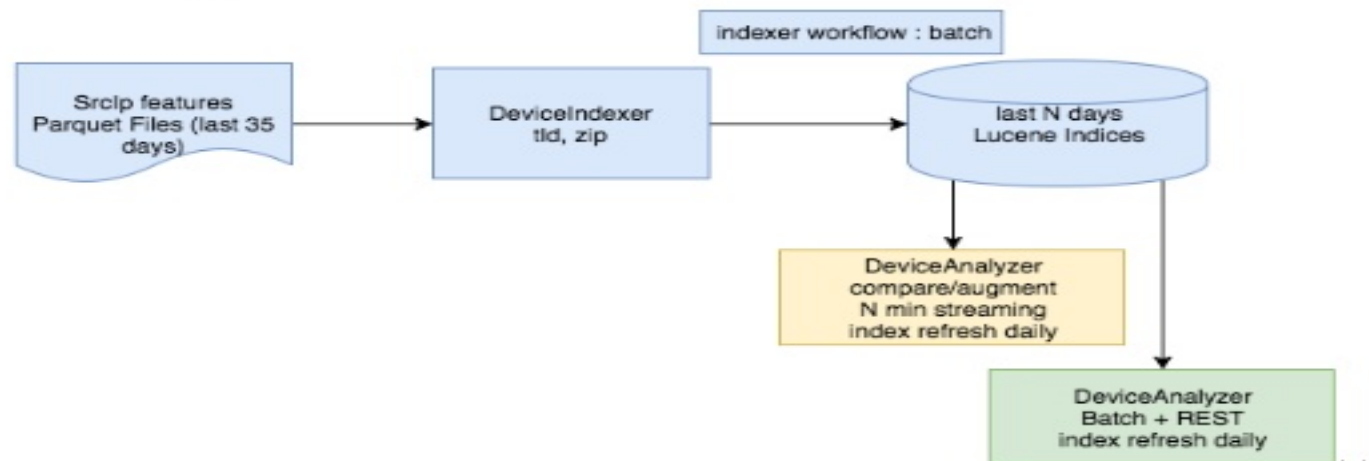- Focus: Indexing + searching documents

# Trapezium LuceneDAO

- SparkSQL and MLlib optimized for full scan, **column indexing not supported**
- Why Spark + Lucene integration
  - Lucene is battle tested Apache Licensed Open Source Project
  - Adds column search capabilities to Spark
  - Adds spark operators (treeAggregate, treeReduce, map) to Lucene
- LuceneDAO features
  - Build distributed lucene shards from Dataframe
  - Save shards to HDFS for QueryProcessor (CloudSolr)
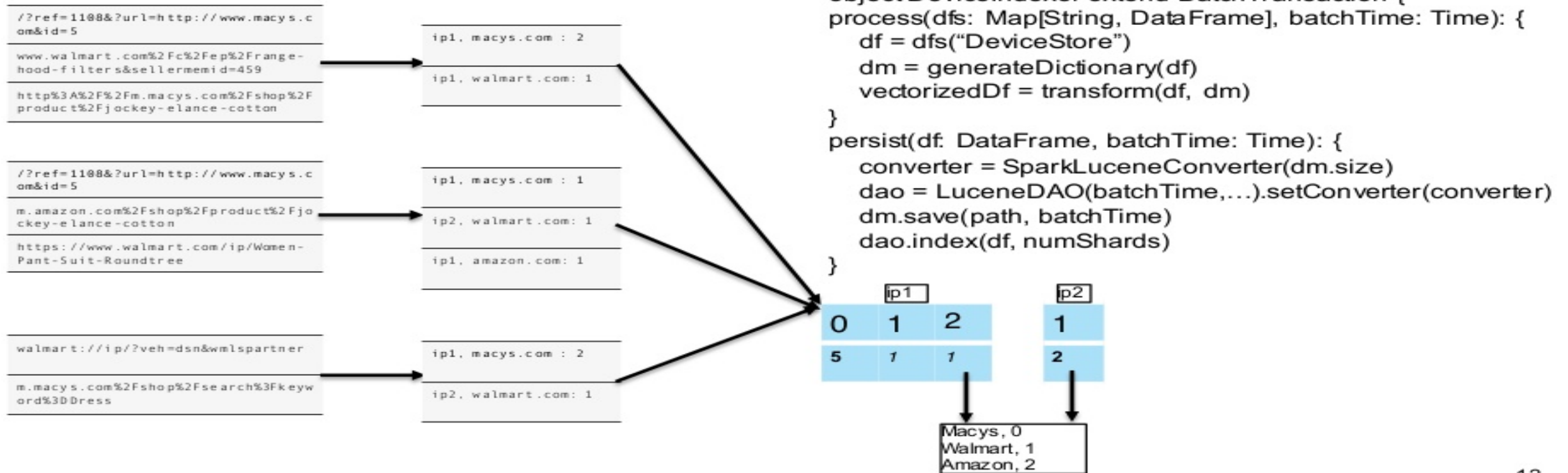  - Access saved shards through LuceneDAO for ML pipelines

# Trapezium Batch

```
runMode = "BATCH"
dataSource = "HDFS"
dependentWorkflows={
  workflows=[aggregate]
  frequencyToCheck=100
}
hdfsFileBatch = {
  batchTime = 86400
  timerStartDelay = 1
  batchInfo = [{
    name = "DeviceStore"
    dataDirectory = {saiph-devqa=/aggregates}
    fileFormat = "parquet"
  }]
}
```

```
transactions = [{
  transactionName  = "DeviceIndexer"
  inputData  = [{name = "DeviceStore"}]
  persistDataName  = "indexed"
}]
```
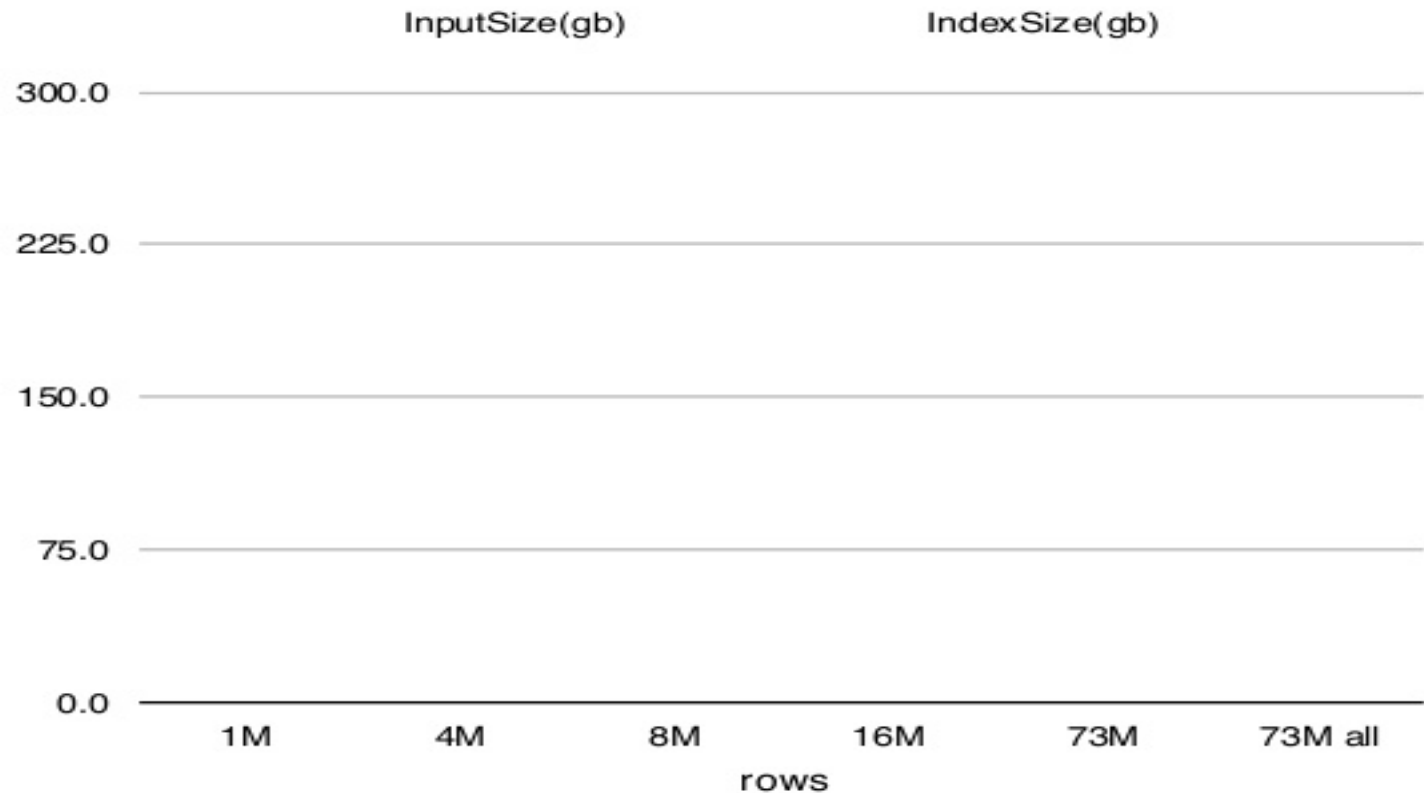
indexer workflow : batch

Srclp features Parquet Files (last 35 days) → DeviceIndexer tld, zip → last N days Lucene Indices

DeviceAnalyzer compare/augment N min streaming index refresh daily

DeviceAnalyzer Batch + REST index refresh daily

11

# DeviceAnalyzer: Indexing

```
/?ref=1108&?url=http://www.macys.c
om&id=5

www.walmart.com%2Fc%2Fep%2Frange-
hood-filters&sellermemid=459

http%3A%2F%2Fm.macys.com%2Fshop%2F
product%2Fjockey-elance-cotton
```

```
ip1, macys.com : 2

ip1, walmart.com: 1
```

```
/?ref=1108&?url=http://www.macys.c
om&id=5

m.amazon.com%2Fshop%2Fproduct%2Fjo
ckey-elance-cotton

https://www.walmart.com/ip/Women-
Pant-Suit-Roundtree
```

```
ip1, macys.com : 1

ip2, walmart.com: 1

ip1, amazon.com: 1
```

```
walmart://ip/?veh=dsn&wmlspartner

m.macys.com%2Fshop%2Fsearch%3Fkeyw
ord%3DDress
```

```
ip1, macys.com : 2

ip2, walmart.com: 1
```

```
object DeviceIndexer extend BatchTransaction {
process(dfs: Map[String, DataFrame], batchTime: Time): {
    df = dfs("DeviceStore")
    dm = generateDictionary(df)
    vectorizedDf = transform(df, dm)
}
persist(df: DataFrame, batchTime: Time): {
    converter = SparkLuceneConverter(dm.size)
    dao = LuceneDAO(batchTime,...).setConverter(converter)
    dm.save(path, batchTime)
    dao.index(df, numShards)
}
```

| ip1 | | | | ip2 | |
|---|---|---|---|---|---|
| 0 | 1 | 2 | | 1 | |
| 5 | 1 | 1 | | 2 | |

```
Macys, 0
Walmart, 1
Amazon, 2
```

# LuceneDAO Index Size

|  | InputSize(gb) | IndexSize(gb) |
|---|---|---|

| rows | InputSize(gb) | IndexSize(gb) |
|---|---|---|
| 1M | 4.0 | 5.1 |
| 4M | 14.4 | 19.0 |
| 8M | 27.9 | 35.7 |
| 16M | 58.8 | 63.2 |
| 73M | 276.5 | 228.0 |
| 73M all | 276.5 | 267.1 |

300.0

225.0

150.0

75.0

0.0

1M    4M    8M    16M    73M    73M all

rows

# LuceneDAO Shuffle Size

Dictionary(mb)     ShuffleWrite(mb)

| rows | ShuffleWrite(mb) | Dictionary(mb) |
|---|---|---|
| 1M | 25 | 22.0 |
| 4M | 56 | 30.0 |
| 8M | 85 | 31.6 |
| 16M | 126 | 32.2 |
| 73M | 334 | 32.4 |
| 73M all | 921 | 146.5 |

# LuceneDAO Index Runtime

20 executors 16 cores
Executor RAM 16 GB
Driver RAM 8g

| rows | Runtime (s) |
|---|---|
| 1M | 135 |
| 4M | 228 |
| 8M | 434 |
| 16M | 571 |
| 73M | 1726 |
| 73M all | 2456 |

Runtime (s)

| | 1M | 4M | 8M | 16M | 73M | 73M all |
|---|---|---|---|---|---|---|

#rows

# Trapezium Api

```
runMode = "BATCH"
dataSource = "HDFS"
httpServer = {
  provider = "akka"
  hostname = "localhost"
  port = 19999
  contextPath = "/"
  endPoints = [{
    path = "analyzer-api"
    className = "TopKEndPoint"
  }]
}
```

# DeviceAnalyzer: Topk

- **Given a query** select * from devices where tld='macys.com' OR 'nordstorm.com' AND (city='SanFrancisco' OR 'Brussels') AND (device='Android') …

  - **ML: Find topk dimensions highly correlated with selected device**

  - **BI: group by tld order by sum(visits) as tldVisits limit topk**

```
class TopkController(sc: SparkContext) extends
SparkServiceEndPoint(sc) {
override def route : topkRoute
converter = SparkLuceneConverter(dm.size)
batchTime = Trapezium.getSyncTime("indexer")
dao = LuceneDAO(batchTime…)
    .setConverter(converter).load(sc, indexPath)
dict = loadDictionary(sc, indexPath. batchTime)
def topkRoute : {
  post { request => {
    devices = dao.search(request)
    response = getCorrelates(devices, dict, topk)
  }
}
}
```

df[deviceId, vector]

sum, support
mean, median, stddev

# Trapezium Stream

```
runMode = "STREAM"
dataSource = "KAFKA"
kafkaTopicInfo = {
  consumerGroup = "KafkaStreamGroup"
  maxRatePerPartition = 970
  batchTime = "5"
  streamsInfo = [{
    name = "queries"
    topicName = "deviceanalyzer"
}]
}
transactions = [{
  transactionName = DeviceAnalyzer"
  inputStreams = [{name:"queries"}]
  persistStreamName = "deviceanalyzer"
  isPersist = "true"
}]
```



18

# DeviceAnalyzer: Compare

- Given two queries

select * from Devices where tld='[macys.com](macys.com)' OR '[nordstorm.com](nordstorm.com)' AND (city='SanFrancisco') AND (device='Android')

select * from Devices where tld='[macys.com](macys.com)' OR '[nordstorm.com](nordstorm.com)' AND (city='Brussels') AND (device='Android')

- Find the dimensions that discriminate the devices associated with two groups

```
def processStream(streams: Map[String,
DStream[Row]], workflowTime: Time): {
  streams("queries").collect().map{ requests =>
    group1 = dao.search(requests(0))
    group2 = dao.search(requests(1))
    response = runLDA(aud1, aud2, dict)
}
```

- Sparse weighted least squares using Breeze QuadraticMinimizer
- L1 Regularized logistic regression

```
def persistStream(responses: RDD[Row],
batchTime: Time) {
    HBaseDAO.write(responses)
}
```

# DeviceAnalyzer: Augment

- Given a query

select * from Devices where
tld='macys.com' OR 'nordstorm.com'
AND (city='SanFrancisco' OR 'Brussels')
AND (device='Android')…

  - Find devices similar to seed as lookalikes

  - Find dimensions that represent lookalikes

```
object DeviceAnalyzer extends StreamingTransaction {
  converter = SparkLuceneConverter(dm.size)
  batchTime = Trapezium.getSyncTime("indexer")
  dao = LuceneDAO(batchTime…)
      .setConverter(converter).load(sc, indexPath)
  dict = loadDictionary(sc, indexPath, batchTime)
  all = dao.search("*:*")
  def processStream(streams: Map[String, DStream[Row]]) :
  {
    streams("queries").collect().map{ request =>
      audience = dao.search(request)
      response = getLookalikeDimensions(all, audience, dict)
  }
```

- Sparse weighted least squares using Breeze QuadraticMinimizer
- L2 regularized linear regression

20

# FastSummarizer

- Statistical and predictive operators
    - sum: sum over numeric measures
    - support: sum over distinct docID
    - sumSquared: L2 norm
    - gram: Uses BLAS sspr
    - solve: Uses BreezeQuadraticMinimizer to support L1
- Implemented using Array[Float] for shuffle opt
- Scala/Java for Level1 operations
- OpenBLAS for Level3 operations

# Sync API Benchmark

73M rows 1M+ search terms

1 measure on 250K sparse dimensions

20 executors 8 cores

32 GB driver RAM 16 GB executor RAM

akka-http cores: 24 default

topk

| qps | runtime(s) |
|-----|-----------|
| 1 | 1.389 |
| 5 | 1.663 |
| 10 | 3.214 |
| 20 | 5.992 |
| 40 | 12.174 |

runtime(s)

# Async API Benchmark

73M rows, 1M+ search terms

1 measure on 250K sparse dimensions

20 executors 8 cores

32 GB driver RAM 16 GB executor RAM

forkjoinpool = 40

Kafka Fetch + compare/augment + HBase Persist

| predictions | | |
|---|---|---|
| qps | compare(s) | augment(s) |
| 1 | 9 | 16 |
| 5 | 13 | 36 |
| 10 | 23 | 70 |
| 20 | 42 | 142 |

compare          augment

160

120

80

40

0

        1          5          10          20

qps

# topk tld + apps

# Augment: Auto Enthusiastic

# Augment Model Performance



**Decile Chart**

■ **Palomar Prediction Audience**
■ Randomly Selected Audience

Targeted Audience Accuracy Score

# Compare: Leisure vs Business Travellers

# THANK YOU.
## Q&A

Join us and make machines intelligent
Data & Artificial Intelligence Systems
499 Hamilton Ave, Palo Alto
California

**verizon**✓