

# Embracing Spark as the Scalable Data Analytics Platform for the Enterprise

---

Matthew J. Glickman

[GS.com/Engineering](http://GS.com/Engineering)

Spark Summit East 2015

# How did we get here today?

- The image shows the logo for the Strata+Hadoop WORLD event. The text "Strata+Hadoop" is in a large, white, serif font, with a red plus sign between "Strata" and "Hadoop". Below this, the word "WORLD" is in a smaller, white, sans-serif font, centered between two horizontal red lines. To the right of the logo, the text "Make Data Work" is in red, and "Oct 15-17, 2014 • New York, NY" is in white.

Strata+Hadoop  
WORLD  
Make Data Work  
Oct 15-17, 2014 • New York, NY
- “Mixing Structured Data and Analytics with Spark SQL”  
Michael Armbrust
- “Make Simple things simple and complex things possible”  
Alan Kay (by way of Ali Ghodsi Spark Summit 2014 demo)
- Spark write-up went viral on our internal social media platform

# I'd Seen the Future in Apache Spark...

- Intuitive language bindings to Scala, Java, Python, R
- Combining relational, functional, iterative APIs all into lazy-evaluation data pipelines
- Storage agnostic
- Lambda closures
- Similar abstraction to GS internal platform's very successful tabular dataset framework
- Scala was already becoming a viable GS platform



# What makes Big Data scalable?

- Elasticity in 3 dimensions:
  - Data Storage = \_\_\_\_\_.
  - Compute = \_\_\_\_\_.
  - Users = \_\_\_\_\_.

# Compute Elasticity =



- But first, how many people have...
  - Used a proprietary data analytics framework?
  - Written their own data analytics framework?
  - Wrapped their own framework around Spark?

# Don't wrap Spark!

- Power of Spark is in the API abstractions (e.g., RDD, DataFrame)
- Spark is becoming “Lingua Franca” of Big Data analytics
- Contribute to open source instead of wrapping!

# Everyone is building Data Lakes

- Universal data acquisition makes all big data analytics and reporting easier
- Hadoop provides a scalable storage with HDFS
- How will we scale consumption and curation of all this data?

# There was a dream that was [Spark]...

- Embrace Spark as the elastic data consumption and curation engine to harness the power of the Data Lake
- All Data Lake datasets available as Spark RDD DataFrames
- Achieve data transformation lineage
  - Data Lake manages DAG of all datasets transformation dependencies
  - Spin up CPU-segregated Spark clusters to compute and store curated data back to Data Lake





# What about Integration?

- Embed Spark driver code directly inside JVM applications just like any other Scala library
- Leverage existing SDLC using existing JVM IDE environment instead of spark-submit for easier debugging
- Dynamically deploy code to cluster at run-time with lambda closures and:

```
val sc = new SparkContext(conf)
sc.addJar(JarCreator.createJarFile(JarCreator.
    getClassesFromClassPath(getClass.getPackage.getName)))
```
- Allow multiple applications with different code wavefronts to share Spark Clusters



~GLICKMAN\SPARKCLEAN2\SOURCE - [gsam-examples] - \\SLJEncodanNpOIVzXJEQnM7U291cmNIOkJSQU5DSDp-Z2xpY2tYV4hc3BhcmJibGVhbjJ7U291cmNle30\gsam-examples\src\com\gs\gsam\examples\Spa...

File Edit View Navigate Code Analyze Refactor Build Run Tools VCS Window JSI Help

gsam-examples src com gs gsam examples SparkExample5.scala

```

130 SparkExample5.scala
131
132 def testQuery5(sc: SparkContext) = { sc: SparkContext@8114
133   val sqlContext = new SQLContext(sc) sc: SparkContext@8114
134   sqlContext.setConf("spark.sql.parquet.binaryAsString", "true")
135   import ...
136   val rdd = sqlContext.parquetFile("hdfs://dev162569-001.dc.gs.com:8020/user/hive
137   println(rdd.schemaString)
138   // show(rdd)
139   val rdd2 = rdd.groupBy('bloomberg_id)(Average('cb_implied_volatility) as 'avg.
140   println(rdd2.count())
141   println(rdd2.toDebugString)
142   for (a <- rdd2.take(10).toSeq) {
143     println(a.mkString(","))
144   }
145 }
146
147 def testQuery6(sc: SparkContext) = {...}
148
149 private val clusterURL = "spark://dev162569-001.dc.gs.com:7077"
150 //private val clusterURL = "local";
151
152 def main(args: Array[String]): Unit = {
153   val conf = new SparkConf()
154     .setMaster(clusterURL)
155     .setAppName("glickman")
156     .set("spark.executor.memory", "20g")
157     .set("spark.cores.max", "100")
158   val classes = Seq(...)
159   val jars = classes.map(_.getProtectionDomain().getCodeSource().getLocation()).ge
160   conf.setJars(jars)
161   val sc = new SparkContext(conf)
162   sc.addJar(JsiJarCreator.createJarFile(JsiJarCreator.getClassesFromClassPath(get
163
164   // testQuery2(sc, conn)
165   testQuery5(sc)
166 }
167
168
169

```

Evaluate Expression

Expression: show(rdd)

Result: result = undefined

Evaluate Code Fragment Mode Close

ID vlf\_batch\_id\_in vlf\_batch\_id\_out vlf\_digest import\_date import\_source ...

vlf_batch_id_in	has cb model	cb market price	cb model price	cb parity	cb premium to parity	cb
1	1	99.06	99.06	19.81		79.25
1	1	96.52	96.25	30.57		64.95
1	1	97.66	93.55	68.60		29.06
1	1	113.15	123.13	120.68	(7.53)	
1	1	72.50	83.27	18.20		54.31
1	1	95.03	95.09	48.29		46.74
1	1	100.00	100.00	37.96		62.04
1	1	132.21	132.24	120.69		11.52
1	1	100.00	100.00	25.06		74.94
1	1	90.00	98.49	66.66		23.34
1	1	140.05	140.05	140.05	0.00	
1	1	99.93	99.92	62.34		37.60
1	1	128.08	129.24	126.37		1.71
1	1	99.23	99.35	36.15		63.08
1	1	97.01	96.87	25.57		71.44
1	1	99.01	99.36	47.05		51.95
1	1	639.91	639.91	671.97	(32.07)	
1	1	112.51	112.51	66.79		45.72
1	1	420.63	420.63	420.73		(0.10)
1	1	82.25	103.12	60.45		21.80
1	1	126.00	120.99	35.78		90.22
1	1	122.38	122.38	104.82		17.57

Pivot: Filter: View: Full

Debug SparkExample5

Debugger Console

Frames

- jsilink-launcher@2111 in group "jsilink...
- testQuery5@137, SparkExample5\$ (com.gs.gsam.examples)
- main@187, SparkExample5\$ (com.gs.gsam.examples)
- main@1, SparkExample5 (com.gs.gsam.examples)

Variables

- this = (SparkExample5\$@8107)
- sc = (SparkContext@8114)
- sqlContext = (SQLContext@8115)
- rdd = (SchemaRDD@8116) "SchemaRDD[0] at RDD at SchemaRDD.scala:111" n= Query Plan == n= Physical Plan == n=

Watches

No watches

G2 Browser Changes Messages Debug TODO

Event Log AREA Explorer Nimbus Explorer GS Projito Browser

1884 CRLF UTF-8

# What are the Integration Challenges?

- Getting machines provisioned to run Spark
- Should Spark and HDFS be run on the same cluster?
  - CPU segregation versus minimizing I/O
  - Data replication for segregation
- Library version synchronization of all open source libraries between Spark, HDFS as well as driver applications
  - Hadoop vendors only offer a partial solution

# Is Enterprise Data going to Public Cloud?

- Managed cloud data services can provide:
  - Centralized management of clusters
  - Offer choice of Hadoop/Spark versions
  - Same SQL, Spark client APIs
  - Elastic scale and self-service automation
- Why not Public Cloud for Enterprise Data:
  - Regulation
  - Single Cloud Provider Lock-in
  - Data Encryption
  - User entitlement provisioning



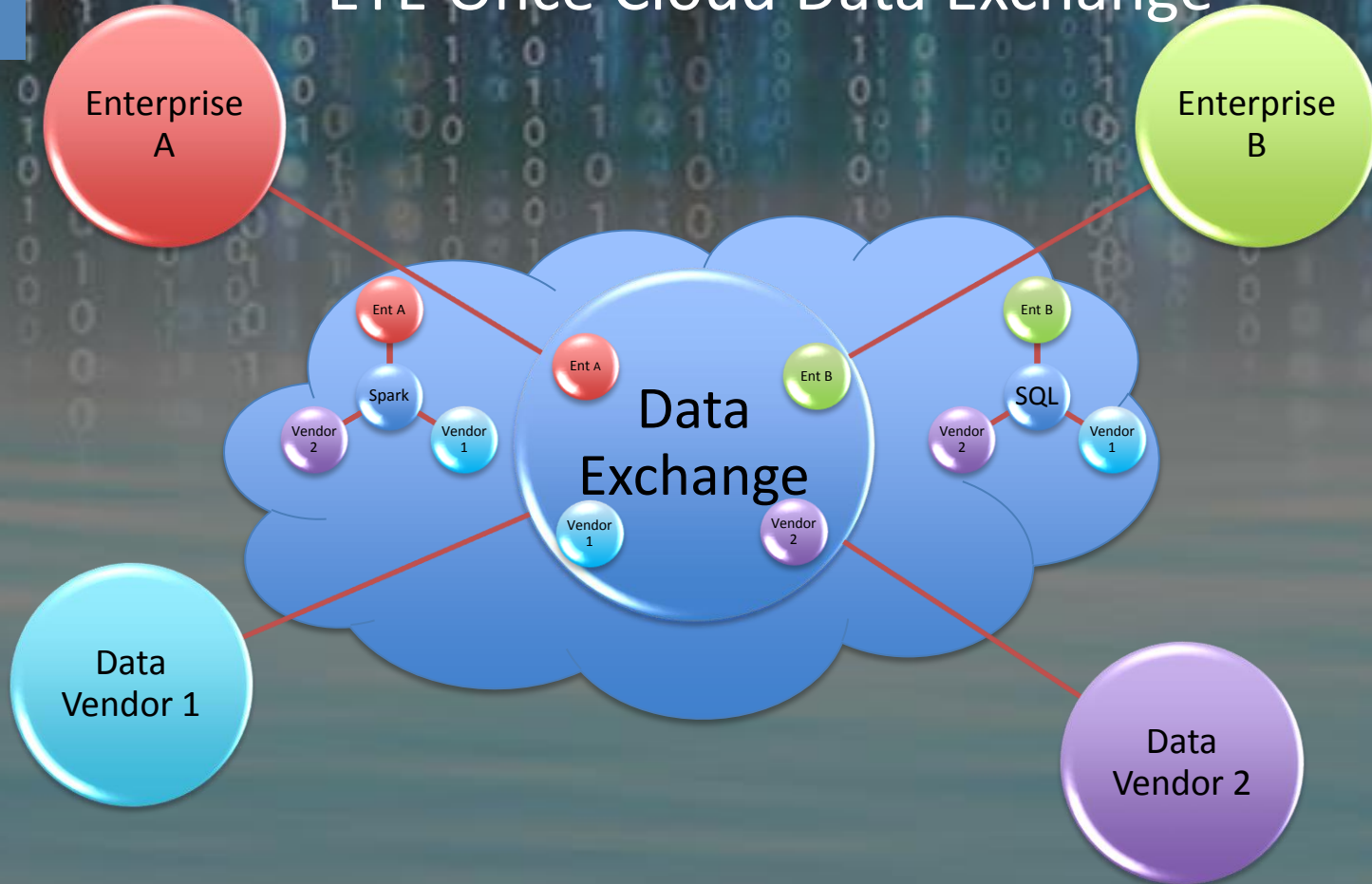
# ETL is still the big problem

- Data needs to be ingested into scalable storage
- Each enterprise will build its own Data Lake
- Moving data is HARD!
- A lot of this data is coming from the same external vendors who are dumping their databases to file feeds for each customer
- Lots of resources at each enterprise are spent reconstructing these vendor databases in Data Lakes alongside internal enterprise data

# Game Changer: Cloud Data Exchanges

- Managed cloud data services have the market disrupting potential to become **Cloud Data Exchanges**:
  - “ETL-once” loading of vendor data by vendors
  - Scalable compute near data with no persistent data movement
  - Vendors can provision access to their data directly to customers
  - No need for each enterprise to ETL the same vendor data
  - Enterprises could then load their data securely alongside vendor data for analytical consumption via standard APIs
- This kind of game changing managed cloud data service will be what really tips the enterprise public cloud data scale

# ETL-Once Cloud Data Exchange



# Takeaways

- Build muscle memory for easier open source contribution
- Think of Spark Client API like ODBC/JDBC
- Embrace don't wrap Spark APIs to prepare for accelerated move of enterprise data to Public Cloud



# we BUILD

Learn more at [GS.com/Engineering](https://www.gs.com/Engineering)

The term 'engineer' referenced in this section is neither a licensed engineer nor an individual offering engineering services to the general public under applicable law.

These materials ("Materials") are confidential and for discussion purposes only. The Materials are based on information that we consider reliable, but Goldman Sachs does not represent that it is accurate, complete and/or up to date, and it should not be relied on as such. The Materials do not constitute advice nor is Goldman Sachs recommending any action based upon them. Opinions expressed may not be those of Goldman Sachs unless otherwise expressly noted. As a condition to Goldman Sachs presenting the Materials to you, you agree to treat the Materials in a confidential manner and not disclose the contents thereof without the permission of Goldman Sachs. © Copyright 2014 The Goldman Sachs Group, Inc. All rights reserved.