# Big Data Technologies And Why They Matter To R Users

Bob Wakefield
Principal
bob@MassStreet.net
Twitter:
@BobLovesData

Mass Street Analytics

# No Deep Dives Today
## It's all 10,000 ft and the speed of heat



Mass Street
Analytics

# Motivation For Today's Presentation

- If I have more data can't I just get a bigger machine?

- How does Spark relate to Hadoop?

Mass Street
Analytics

# Today's Take Away

- The way you access data in the future is going to radically change. Be aware and ready.

# Intended Audience

- People getting started in data

- Academics that work with data

- (If you're already working in industry, you might be bored

# What exactly have you been studying for the past four years?!

Mass Street Analytics

# Big Data Landscape 2016 (Version 2.0)

# We're Experiencing A Modern Tech Renaissance

Hadoop -> SQL On Hadoop -> NoSQL ->

In Memory Databases ->MPP -> Spark 1.0 ->

Spark 2.0 -> Deep Learning -> Blockchain

Mass Street Analytics

# Podcast Plug

- Not So Standard Deviations

- Roger Peng, PhD Biostatistician Johns Hopkins School Of Public Health

- Hillary Parker, PhD Data Scientist Stitch Fix

Mass Street Analytics

# My Workflow



Mass Street Analytics

# Data Retrieval And Storage
# What is Hadoop?

- Can mean a lot of things depending on context.
  - Could mean Hadoop core.
  - Could mean the entire set of Big Data tools.
  - Could mean the MapReduce framework.

Mass Street
Analytics

# Data Retrieval And Storage What is Hadoop?

- Provides distributed fault tolerant data storage

- Provides linear scalability on commodity hardware

- Translation: Take all of your data, throw it across a bunch of cheap machines, and analyze it. Get more data? Add more machines

Mass Street Analytics

# Data Retrieval And Storage Cloud Storage

- **AWS**
  - **S3 is the gold standard**
- **MS Azure**
- **Google Cloud**

Mass Street Analytics

# Cloud Solutions Bob's Thoughts

- **Great for tactical solutions.**
- **If you're a smaller organization, cloud offers a managed solution.**
- **Can be a hassle. You need top flight network engineers.**
- **Moving to the cloud is not something to do lightly. It requires long term thinking.**
- **I have no idea why anybody waste time with on prim servers anymore.**

Mass Street Analytics

# Data Retrieval And Storage
# What Is Hive?

- Hadoop warehouse solution
- SQLesque language called Hive Query Language
- Adds structure to unstructured data
- Provides a window into HDFS
- You can connect through ODBC/JDBC
  - Which means you can work with Hive using standard BI tools

Mass Street Analytics

# Data Retrieval And Storage NoSQL Databases

- Key Value – A simple hash table, primarily used when all access to the database is via primary key.

- Document – The database stores and retrieves documents, which can be XML, JSON, BSON

- Column-family -  Data is stored in columns

- Graph – allows you to store relationships between entities

# Data Retrieval And Storage Issues with NoSQL Databases

- They do NOT work like relational databases
  - The guarantees that you are used to aren't there
  - CAP theory (consistency, availability, partition tolerance)

- Each database has it's own query language
  - That language will frequently look NOTHING like SQL

Mass Street Analytics

# Data Retrieval And Storage Examples of NoSQL Databases

| Document | Key Value | Column Store | Graph |
|---|---|---|---|
| Mongo DB | Apache Accumulo | Apache Accumulo | Neo4J |
| | Redis (in memory) | Druid | |
| | | Cassandra (DataStax) | |
| | | Hbase | |

Mass Street Analytics

# Data Retrieval And Storage
# New SQL Databases

- Distributed versions of databases you're used to

- New concept Hybrid Transactional/Analytical Processing (HTAP)

Mass Street Analytics

# Data Retrieval And Storage New SQL Databases

- In Memory Databases
  - MemSQL
  - VoltDB
  - NuoDB

Mass Street Analytics

# Data Retrieval And Storage MPP Databases

- More suited to analytics
- Examples
  - Greenplumb
  - SQL Server PDW
  - MySQL Cluster CGE

Mass Street Analytics

# Real Time Processing

- **Apache Spark Structured Streaming** - scalable and fault-tolerant stream processing engine built on the Spark SQL engine.

- **Apache Storm** – Storm is a distributed real-time computation system for processing fast, large streams of data adding reliable real-time data processing capabilities to Apache Hadoop 2.x.

- **Apache Flume** – Flume allows you to efficiently aggregate and move large amounts of log data from many different sources to Hadoop.

# Real Time Processing

- Other real time processing frameworks
  - Apache Flink
  - Apache Samza

- They all have benefits and drawbacks
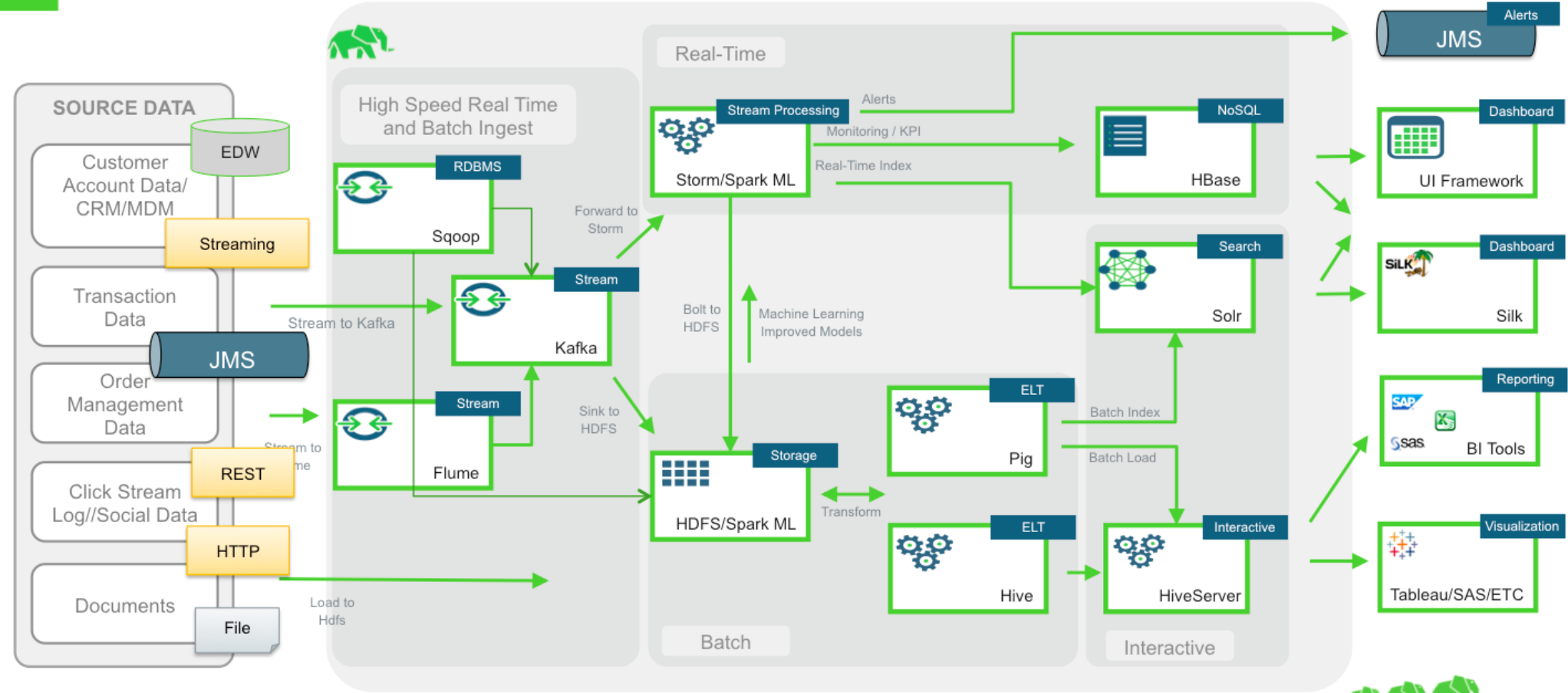
Mass Street
Analytics

# What Is Spark

- Distributed in memory processing framework
- Rapidly replacing MapReduce as a means to crunch data.
- Many ways to interact with Spark
  - R, Java, Scala, Python
  - Sparkr, Sparklyr, H2O with Sparkling Water
- Several APIs
  - RDDs, Dataset/Dataframe, Spark Streaming, Spark SQL, Spark Streaming, Strucutred Streaming

Mass Street Analytics

# How You Work Now

# What You Might Have To Deal With

# Q&A