

Understanding MLOps with Azure

February 2020



Introductions



Raki Rahman

Solution Architect

slalom | Data & Analytics

<https://rakirahman.com>

Professional Experience

7+ years at Technology Consulting and Software Engineering firms:

- Big Data and Software Engineering on Azure + Open source stack
- Environment Operationalization: DevOps, MLOps, DataOps + ***anything else that ends with "Ops"***
- Delivering and operationalizing Azure environments for large Canadian Enterprise Clientele

Education

B.A.Sc. Aerospace Engineering (Engineering Science) – University of Toronto



Agenda

00:00*

Deployment Kickoff: Azure Environment Spinup for demo

00:02



Big Data Demystified: Hadoop, Spark and the Cloud

00:30

Project Showcase



Bloomberg overview and Project Goal

Azure Solution Architecture

Bitcoin POC use case with Natural Language Processing

00:50

Useful resources

NLP applied - Harry Potter Q&A - BERT on Azure Kubernetes

Big Data Deep Dive - Spark Certification Study Guide

MLOps Deep Dive - MLOps on Azure with Malware Detection

01:00

5 Minute Q/A

01:05

Miflow Deep Dive – Predicting Airbnb listing prices

01:50

Demo Q&A

*Timings are approximate ☺



Demo Prep

Azure Environment Deployment Kickoff



Microsoft Azure



Big Data Demystified

Hadoop, Spark and the Cloud

Recurring questions

Prevalent across multiple Industry Verticals

- **What's so BIG about Big Data** - that I can't solve by just throwing more hardware at my On-Premise SQL Databases?
- Does Big Data processing use some sort of **special computers**?
- If not, how does it **actually leverage** regular Computer Architecture (that you and I are using right now – our laptops) to **process large amounts of data**?
- What's **Distributed Computing** and what do I need to know? Why should I care and how is it different from my On-Premise SQL Databases?
- What's **Hadoop**? What's **Spark**? And why do both exist?
- I already have expensive servers I bought from Dell in my Data Center, who should I **care about the Cloud**?
- What's the **big deal about Azure Databricks**?

"Open this 1 Terabyte XSLX file in Excel"

Motivation Use Case

- You work for a popular online shopping website
- Your boss comes by with a 2TB external Samsung hard drive containing a **1 TB XSLX** file that contains 10 years worth of transaction information
- You're asked to perform a straightforward routine task:
 1. Open the file in Excel
 2. Do a sum across the "[Revenue \(column J\)](#)" to generate the Revenue to date
 3. Send that single number over to your boss via e-mail

Inventory and Revenue sheet.xlsx
1,073,741,824 KB

Inventory & Revenue Sheet

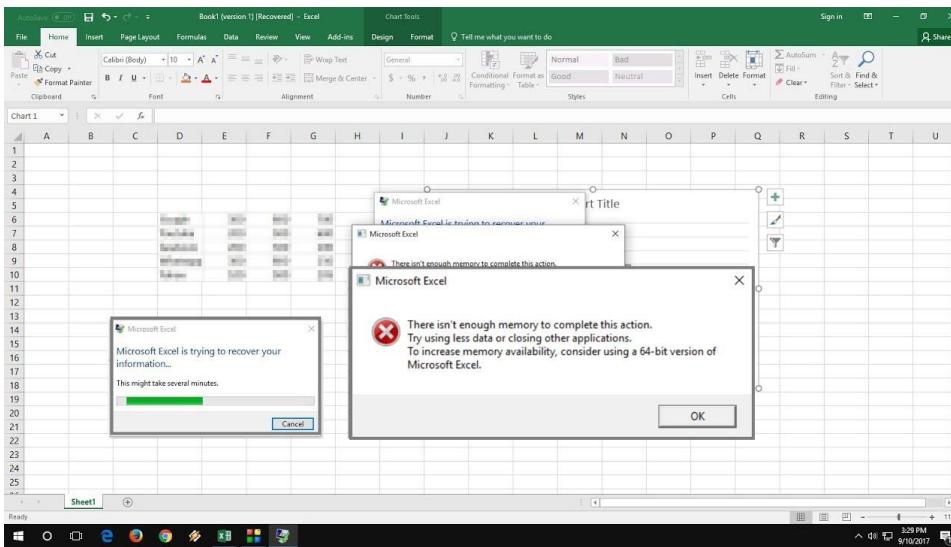
Candles	Quantity Made	Quantity Sold	Current Quantity	Sales Price per Candle	Cost per Candle	Total Sales Revenue	Total Cost
Fragrance 1	30.0	12.0	18.0	\$ 5.00	\$ 2.00	\$ 60.00	\$ 2.00
Fragrance 2	23.0	16.0	7.0	\$ 6.00	\$ 1.00	\$ 96.00	\$ 1.00
Fragrance 3	35.0	14.0	21.0	\$ 8.00	\$ 2.00	\$ 112.00	\$ 2.00
Fragrance 4	26.0	9.0	17.0	\$ 9.00	\$ 3.00	\$ 81.00	\$ 3.00
Fragrance 5	24.0	8.0	16.0	\$ 5.00	\$ 2.00	\$ 40.00	\$ 1.00
Fragrance 6	15.0	5.0	10.0	\$ 4.00	\$ 1.00	\$ 20.00	\$ 1.00
Fragrance 7	25.0	6.0	19.0	\$ 6.00	\$ 2.00	\$ 36.00	\$ 1.00
Fragrance 8	32.0	7.0	25.0	\$ 8.00	\$ 3.00	\$ 56.00	\$ 3.00
Fragrance 9	26.0	11.0	15.0	\$ 7.00	\$ 2.00	\$ 77.00	\$ 2.00
Fragrance 10	29.0	12.0	17.0	\$ 5.00	\$ 1.00	\$ 60.00	\$ 1.00
Fragrance 11	35.0	13.0	22.0	\$ 8.00	\$ 2.00	\$ 104.00	\$ 2.00
Fragrance 12	26.0	16.0	10.0	\$ 9.00	\$ 3.00	\$ 144.00	\$ 3.00
Fragrance 13	24.0	14.0	10.0	\$ 6.00	\$ 2.00	\$ 84.00	\$ 2.00
Fragrance 14	15.0	8.0	7.0	\$ 9.00	\$ 2.00	\$ 72.00	\$ 1.00
Fragrance 15	25.0	5.0	20.0	\$ 9.00	\$ 3.00	\$ 45.00	\$ 1.00
Fragrance 16	32.0	6.0	26.0	\$ 8.00	\$ 2.00	\$ 48.00	\$ 1.00
Fragrance 17	26.0	7.0	19.0	\$ 8.00	\$ 1.00	\$ 56.00	\$ 1.00
Fragrance 18	29.0	11.0	18.0	\$ 8.00	\$ 2.00	\$ 88.00	\$ 2.00
Fragrance 19	32.0	12.0	20.0	\$ 7.00	\$ 3.00	\$ 84.00	\$ 3.00
Fragrance 20	26.0	13.0	13.0	\$ 7.00	\$ 2.00	\$ 91.00	\$ 2.00
Fragrance 21	29.0	16.0	13.0	\$ 8.00	\$ 1.00	\$ 128.00	\$ 1.00

2 Million rows

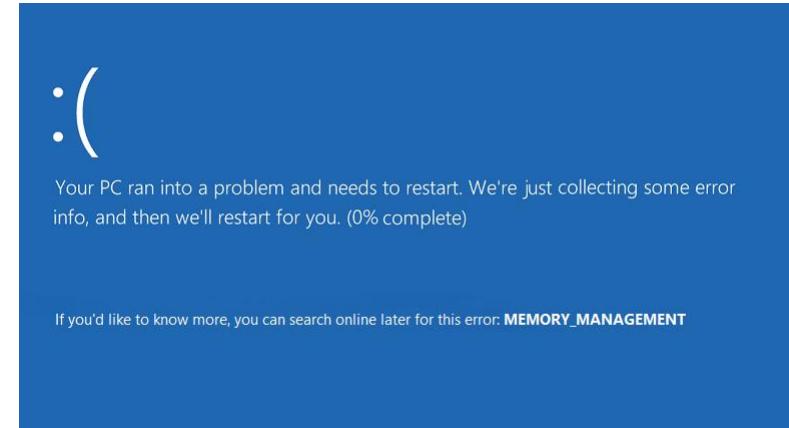
"Open this 1 Terabyte XSLX file in Excel"

Motivation Use Case

- You realize your laptop only has 500 GB of space – **so you can't even copy the 1000 GB file to local storage (C:\ Drive)**
- You try to open the file directly from the External drive, and Excel hits you with an "**Out of Memory**" error

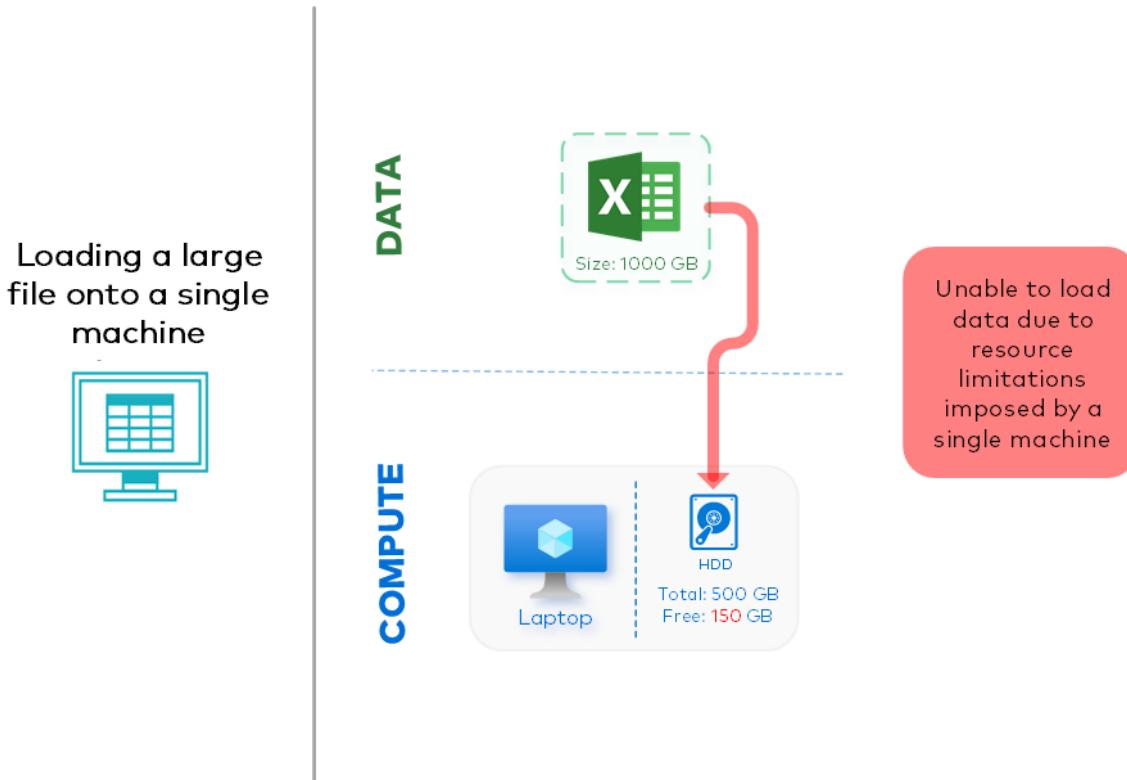


You keep on trying, and get hit with a
Blue Screen of Death



Current Setup – 1 Machine

Motivation Use Case



Loading a large file onto a single machine



DATA

COMPUTE

In conclusion, you realize there's no way given the **resource constraints of 1 Machine** to open this file and perform this simple revenue aggregation.



Brilliant Idea – “Divide and Conquer”

Motivation Use Case

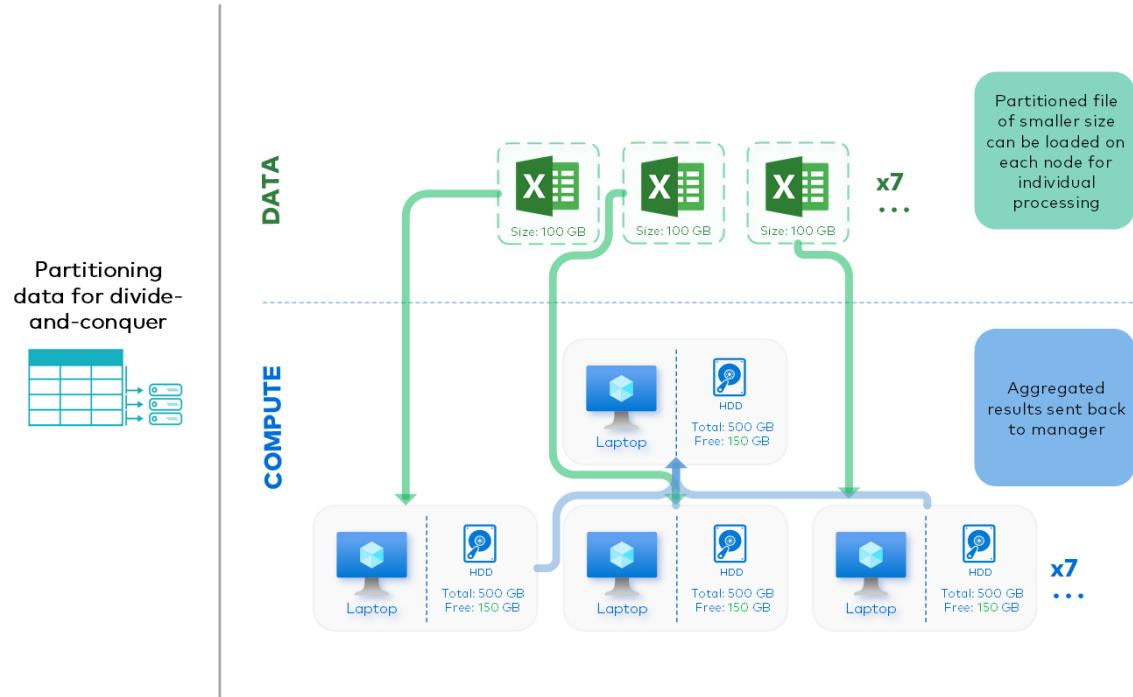
You realize you're able to pass files to your office colleagues' laptops since you're all connected to the office's LAN!

So what if you tried using **10 laptops** to open **10 chunks** of this file? That should work, since most of them have at least 100 GB of space free, so you can:

1. Chunk, or **partition** the 1000 GB file into 10 files of 100 GB each (or say, 1 years worth of transaction per chunk - i.e. **partition by date**)
2. Send each partition to each identical laptop via the office's **LAN** (local area network – **connecting the machines over wire**), and copy it to the laptop's local hard drive so Excel can open it within the machine
3. Open the partitioned file with Excel on each laptop, and do a sum across the price column
4. Take that price aggregation, and save it as a **tiny separate CSV file** (one line containing the price aggregation) on each laptop
5. You send your colleagues an email containing a **network accessible folder** hosted on your laptop, and they each send you all 10 tiny CSVs to that location
6. You open each CSV file, sum up the 10 numbers, and report the full sum back to your boss

Current Setup – 10 Machines

Motivation Use Case

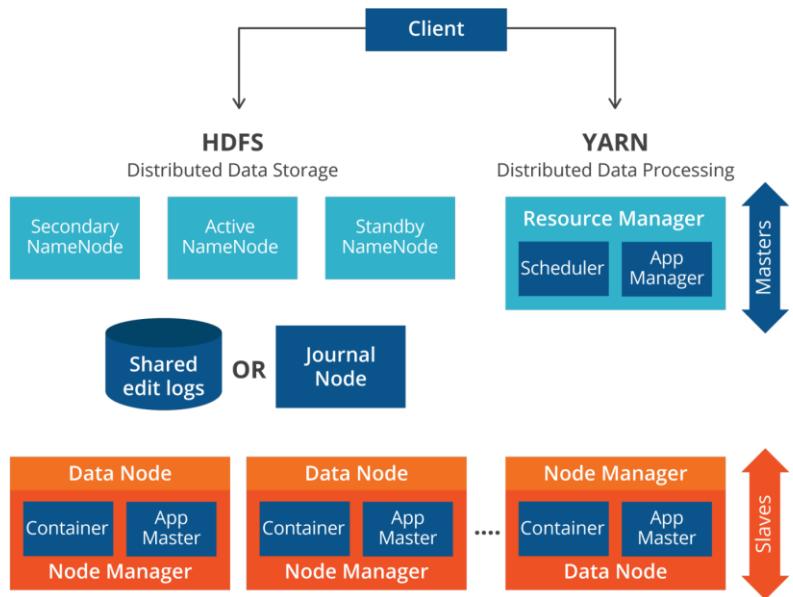


What you just implemented above is called a "**Master/Slave Architecture**", and what is essentially a primitive rendition of Hadoop MapReduce (and as we'll get to – Driver/Executor rendition of Spark).

Hadoop MapReduce

Master/Slave Architecture

Apache Hadoop 2.0 and YARN



HDFS (hadoop-distributed-file-system) Distributed file system engine designed to run on commodity hardware (your laptop, Samsung External Drive etc). Can store an infinite amount of information because you can throw an infinite of hard drives at the cluster (i.e. ignoring the space-time constraint)

HDFS: You + Your colleagues' laptops - Cheap, easily replaceable, interconnected, fault tolerant storage

Master/Slave

A **master server** manages the file system namespace and regulates access to files by clients. **Slave server(s)** take partitioned data from the Master, process them, and send the results back to the specified directory.

Master: You + Your laptop was the Master that assigned the workloads, gathered the partitioned results back, and aggregated the final value.

Slave: You + Your colleagues' laptops

YARN – Yet Another Resource Negotiator

Resource manager that negotiates resources and workloads across the Cluster to ensure each partition is getting processed for the job to finish timely. This is the engine responsible for intelligent workload assignment and performance tracking.

YARN: You (if a colleague is out sick, you find another, or reassign re-calculated workloads with 9). If one colleague is being specially slow, you yell at them.

Hadoop MapReduce

Birth

Apache Hadoop was released to the public in 2006 under an Open-Source license, made possible thanks to [this](#) scientific paper Google published in 2003 as the "**The Google File System**"

MapReduce is the actual engine (replacing your **pseudo Project/Resource Manager role**). It's responsible for 2 main activities:

- Data Partitioning and Mapping** – sending partitioned data to available slave servers
- Reduction** – bringing the tiny files back, and aggregating the result

Fun Fact: the founder - Doug Cutting who worked at Yahoo! at the time - named the project after his son's toy elephant "Hadoop".



The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung
Google

1. INTRODUCTION

We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven by key observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system design assumptions. We have reexamined traditional choices and explored radically different points in the design space.

The file system has successfully met our storage needs. It is widely deployed within Google as the storage platform for the generation and processing of data used by our service as well as research and development efforts that require large data sets. The largest cluster to date provides hundreds of terabytes of storage across thousands of disks on over a thousand machines, and it is concurrently accessed by hundreds of clients.

In this paper, we present file system interface extensions designed to support distributed applications, discuss many aspects of our design, and report measurements from both micro-benchmarks and real world use.

Categories and Subject Descriptors

D [4]: 3—*Distributed file systems*

General Terms

Design, reliability, performance, measurement

Keywords

Fault tolerance, scalability, data storage, clustered storage

*The authors can be reached at the following addresses:
[\[sanjay, hogoboff, shuntak\]@google.com](mailto:sanjay, hogoboff, shuntak}@google.com).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SOSP '03, October 19–22, 2003, Bolton Landing, New York, USA.
Copyright 2003 ACM 1-58113-757-5/03/0010 ...\$5.00.

First, component failures are the norm rather than the exception. The file system consists of hundreds or even thousands of storage machines built from inexpensive commodity parts and is accessed by a comparable number of client machines. The quantity and quality of the components virtually guarantees that some are not functional at any given time and some will not recover from their current failures. We have seen problems caused by application bugs, operating system bugs, human errors, and the failures of disks, memory, connectors, networking, and power supplies. Therefore, constant monitoring, error detection, fault tolerance, and automatic recovery must be integral to the system.

Second, files are huge by traditional standards. Multi-GB files are common. Each file typically contains many application objects such as web documents. When we are regularly working with fast growing data sets of many TBs comprising billions of objects, it is unwieldy to manage billions of approximately KB-sized files even when the file system could support it. As a result, design assumptions and parameters such as I/O operation and block sizes have to be revisited.

Third, most files are mutated by appending new data rather than overwriting existing data. Random writes within a file are practically non-existent. Once written, the files are only read, and often only sequentially. A variety of data share these characteristics. Some may constitute large repositories that data analysis programs scan through. Some may be data streams continuously generated by running applications. Some may be archival data. Some may be intermediate results produced on one machine and processed on another, whether simultaneously or later in time. Given this access pattern on huge files, appending becomes the focus of performance optimization and atomicity guarantees, while caching data blocks in the client loses its appeal.

Fourth, co-designing the applications and the file system API benefits the overall system by increasing our flexibility.

HDFS

Fun fact about Cloud Storage

HDFS – founded 2003 by Google – is what made Cloud Storage and Distributed computing of today possible. This is how we store data dirt cheap in the Cloud (iCloud, Dropbox, OneDrive etc.)

Companies (Microsoft, Amazon, Google) have distributed, interconnected Data Centers across the world, with millions of Hard Drives. **HDFS makes the storage capacity of these drives ~infinite.**

These companies have an army of folks that take care of these servers (air conditioned, contained environments etc.). **Their job is to make sure any data that's put in there is in good health** (99.999...% SLA).

When you and I put our photos in the OneDrive Cloud, we're basically **renting out a tiny portion of Microsoft's massive HDFS cluster – which is why Cloud Storage is dirt cheap.**



MapReduce – Pain Points

Inherent design-based limitations

MapReduce Summary: a reliable, good-for-its-time, Big Data processing engine, that did the best it could with the technology it originated from.

The Fundamental Problem:

- MapReduce persists the dataset at **each step along the way on the hard disk**
 - **Hard disks were the best we had in 2003 to resiliently store data** (we're talking Computer Science of Data Structures) – MapReduce reflected this limitation
 - **Hard disks are slow.** Specially if you're using the cheap kind (non Solid-State-Drives).

So basically, combine:

1. **Hard disks**, specially mechanical ones (i.e. non SSDs), which are fundamentally slow (thanks to the laws of Physics), but reliable/resilient at storing data
2. **Unoptimized I/O Operations**, to move data across the network (you can think of this as a sophisticated spider-web of ethernet cables) back and forth between Hard disks on different machines, which is also slow

Note that none of this really takes proper advantage of the fact that there's a much faster way to access data in a Computer or Cluster of Computers -

the RAM.

- **RAM – Random Access Memory** – was always used for temporary caching of information, so MapReduce couldn't use RAM for anything related to transforming Data
- The person who invented RAM back in 1968 - Robert Dennard - never meant for it to be a resilient data store (*write cached data quickly, read quickly, forget-about-it*)

RAM is faster

Putting things into Perspective

~100x faster
than HDD based
architecture

Hard drive

MB/s	Read	Write
100	100	

Booting up your laptop from 2001 – 30 seconds

Solid State Drive

MB/s	Read	Write
500	250	

Booting up your laptop from today – 10 seconds

RAM

MB/s	Read	Write
6000	8500	

Playing a video game on your PS4 – ~instant response to action

Apache Spark

Birth

Matei Zaharia was at UC Berkeley when the Netflix Challenge (solving Sparse Matrix for Recommendations Engine) was released in 2009. He tried to solve the problem with Hadoop's MapReduce before quickly realizing the algorithm was being **bottlenecked by the Platform's Read-Write throughput**, rather than execution efficiency.

Over that weekend, he invented a premature rendition of the **RDD – Resilient Distributed Dataset** – which essentially allowed for his Server Cluster to store data resiliently at the RAM layer - UC Berkley RDD [Paper](#) on the right.

Apache Spark was released to the public in 2014 under an Open-Source license, fueled primarily by the power of the RDD, and the capabilities it brought.

- Spark runs workloads **100x faster** than Hadoop's MapReduce.
- Spark has over 1450 contributors, the **most adopted Big Data open-source project**.
- From 2014, Spark has been widely adopted in the industry as the **fastest way to process data** on modern computer paradigm.

Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing

Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, Ion Stoica
University of California, Berkeley

Abstract

We present Resilient Distributed Datasets (RDDs), a distributed memory abstraction that lets programmers perform in-memory computations on large clusters in a fault-tolerant manner. RDDs are motivated by two types of applications that current computing frameworks handle inefficiently: iterative algorithms and interactive data mining tools. In both cases, keeping data in memory can improve performance by an order of magnitude. To achieve fault tolerance efficiently, RDDs provide a restricted form of shared memory, based on coarse-grained transformations rather than fine-grained updates to shared state. However, we show that RDDs are expressive enough to capture a wide class of computations, including recent specialized programming models for iterative jobs, such as Pregel, and new applications that these models do not capture. We have implemented RDDs in a system called Spark, which we evaluate through a variety of user applications and benchmarks.

1 Introduction

Cluster computing frameworks like MapReduce [10] and Dryad [19] have been widely adopted for large-scale data analytics. These systems let users write parallel computations using a set of high-level operators, without having to worry about work distribution and fault tolerance.

Although current frameworks provide numerous abstractions for accessing a cluster's computational resources, they lack abstractions for leveraging distributed memory. This makes them inefficient for an important class of emerging applications: those that reuse intermediate results across multiple computations. Data reuse is common in many *iterative* machine learning and graph algorithms, including PageRank, K-means clustering, and logistic regression. Another compelling use case is *interactive* data mining, where a user runs multiple ad-hoc queries on the same subset of the data. Unfortunately, in most current frameworks, the only way to reuse data between computations (*e.g.*, between two MapReduce jobs) is to write it to an external stable storage system, *e.g.*, a distributed file system. This incurs substantial overheads due to data replication, disk I/O, and serializa-

tion, which can dominate application execution times.

Recognizing this problem, researchers have developed specialized frameworks for some applications that require data reuse. For example, Pregel [22] is a system for iterative graph computations that keeps intermediate data in memory, while HaLoop [7] offers an iterative MapReduce interface. However, these frameworks only support specific computation patterns (*e.g.*, looping a series of MapReduce steps), and perform data sharing implicitly for these patterns. They do not provide abstractions for more general reuse, *e.g.*, to let a user load several datasets into memory and run ad-hoc queries across them.

In this paper, we propose a new abstraction called *resilient distributed datasets (RDDs)* that enables efficient data reuse in a broad range of applications. RDDs are fault-tolerant, parallel data structures that let users explicitly persist intermediate results in memory, control their partitioning to optimize data placement, and manipulate them using a rich set of operators.

The main challenge in designing RDDs is defining a programming interface that can provide fault tolerance *efficiently*. Existing abstractions for in-memory storage on clusters, such as distributed shared memory [24], key-value stores [25], databases, and Piccolo [27], offer an interface based on fine-grained updates to mutable state (*e.g.*, cells in a table). With this interface, the only ways to provide fault tolerance are to replicate the data across machines or to log updates across machines. Both approaches are expensive for data-intensive workloads, as they require copying large amounts of data over the cluster network, whose bandwidth is far lower than that of RAM, and they incur substantial storage overhead.

In contrast to these systems, RDDs provide an interface based on *coarse-grained* transformations (*e.g.*, map, filter and join) that apply the same operation to many data items. This allows them to efficiently provide fault tolerance by logging the transformations used to build a dataset (*its lineage*) rather than the actual data.¹ If a partition of an RDD is lost, the RDD has enough information about how it was derived from other RDDs to recompute

¹Checkpointing the data in some RDDs may be useful when a lineage chain grows large, however, and we discuss how to do it in §5.4.

Why Spark is a big deal

It's not going anywhere soon

From 2014, Spark has been widely adopted in the industry as the fastest way to process data on modern computer paradigm.

Simple concept:

1. RAM is the fastest way to perform data operations on modern computer architecture the world runs on today (Motherboard, CPU, Hard Drive, RAM) – **this is an indisputable truth**

2. Spark's Engine runs on the RDD – which is the brilliant way Matei invented to store data resiliently on the RAM layer – **you can't really reinvent this**
= Spark is the fastest possible way to analyze Data

...Until Quantum Computers become publicly accessible:

- Manufacturing of Quantum computers becomes as streamlined as manufacturing Binary computers (**good luck with 50+ years in the making**)
- Running an ETL job in Quantum stops costing a [few million dollars per run](#)

Quantum Computing can dethrone Spark's Execution capabilities, once we figure out all the logistical issues in operationalizing it for public use ([don't hold your breath](#)).

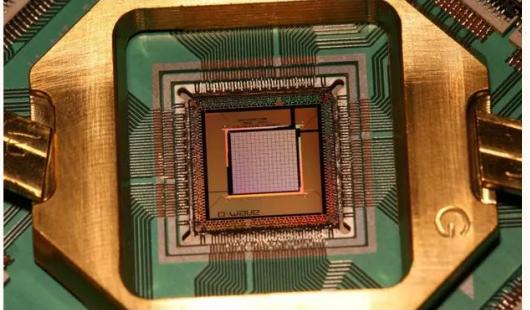
Computing

Sabine Hossenfelder

Fri 2 Aug 2019 07.00 BST



193 323



▲ Quantum computers can solve certain problems much faster than conventional ones. However, such systems are exceptionally fragile, sensitive and expensive. Photograph: D-Wave

The unveiling of the marvel had the media gushing. It was Valentine's Day 1946, and the New York Times broke the story. [The front page](#) spoke of "an amazing machine" and "one of the war's top secrets". By crunching numbers at unprecedented speed, the Electronic Numerical Integrator and Computer, with its 18,000 vacuum tubes, was poised to "revolutionise modern engineering". Eniac would usher in a new epoch of industrial design, some said.

More than 70 years on, another overblown announcement is near. Several companies, notably Google, IBM and the California-based Rigetti, are racing to build a machine that achieves what is grandly termed "quantum supremacy". The feat will mark the moment when a quantum computer, for the first time, outperforms the best conventional computers. Google, the frontrunner, could claim the record this year.

Quantum computers can solve certain problems much faster than conventional ones. Drug development, materials design, weather forecasts, stock trades - all could potentially be helped by quantum computers that perform scores of calculations at once. Physicists have known this since the 1980s, but have only recently made working prototypes.

Operationalizing Spark

Spark + Scalability = Awesome

Now that we've established that Spark is amazing:

How does your organization actually implement it for real-life day-to-day benefit?

Where can you get:

1. **Enough servers in a Data Center** to run your Big Data workloads right In-Memory with Spark, with the ability to scale-up when you process data in bursts, and completely scale-down so you spend \$0 during downtime
2. **A reliable resource negotiator** that takes away the overhead of managing YARN out of your hands and does it all behind the scenes (with room for tweaking, should you choose)
3. **Managed Clusters**, that removes the burden of updating Spark, Python, Scala and other Drivers as new versions are released, as well as dependency and library management
4. **An intuitive, easy, visual interface** to curate your workloads (like Jupyter notebooks) using a multitude of programming languages (like Python, SQL, R, Scala), depending on what your team is comfortable with (I don't know any R and Scala for example, but can get around in Python and SQL – this doesn't mean I should be forced to learn R to do my job)
5. **Secured integration** with your other Enterprise systems/Cloud Services (Like Azure Storage Account, SQL Data Warehouse, Azure Synapse, Tableau and Power BI)
6. **Keep up to date with the industry as it evolves** without you having to invest in architecture uplifts/data processing engine migrations every year

Azure Databricks

Spark + Scalability = Awesome

The Engine



The Platform



The Solution



Azure Databricks

Job execution engine able to efficiently distribute and execute tasks on RAM layer

Scalable RAM + Processing power without management overhead

Managed Apache Spark Platform able to seamlessly integrate within Azure Big Data/ML ecosystem

Matei Zaharia

Professor of Computer Science
Stanford University

The Engine



Founded 2009

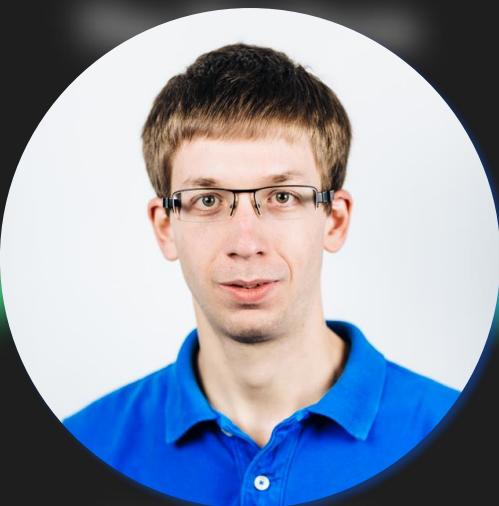
Unbeatable performance



Open Source transparency



The best Big Data processing engine



Founder of:

Apache Spark
Databricks
Delta Lake
MLflow

The Solution



Azure Databricks

Founded 2018

Cloud Scalability



Secret Sauce (Catalyst optimizer
+ Resource Negotiator)



The best platform to implement Spark

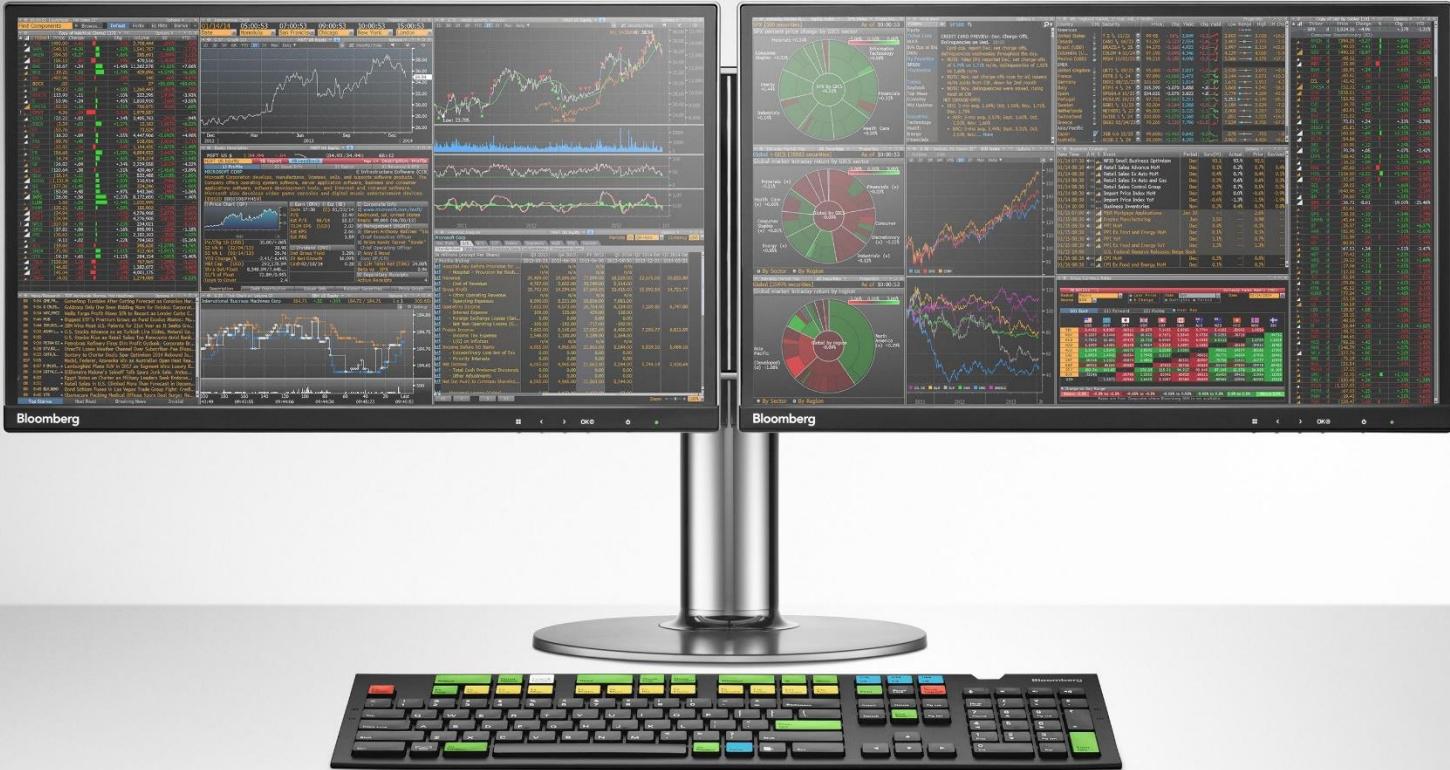


+ Azure Databricks + Bloomberg

Project Showcase

Azure Real-time Streaming from Bloomberg

Bloomberg overview



Bloomberg overview

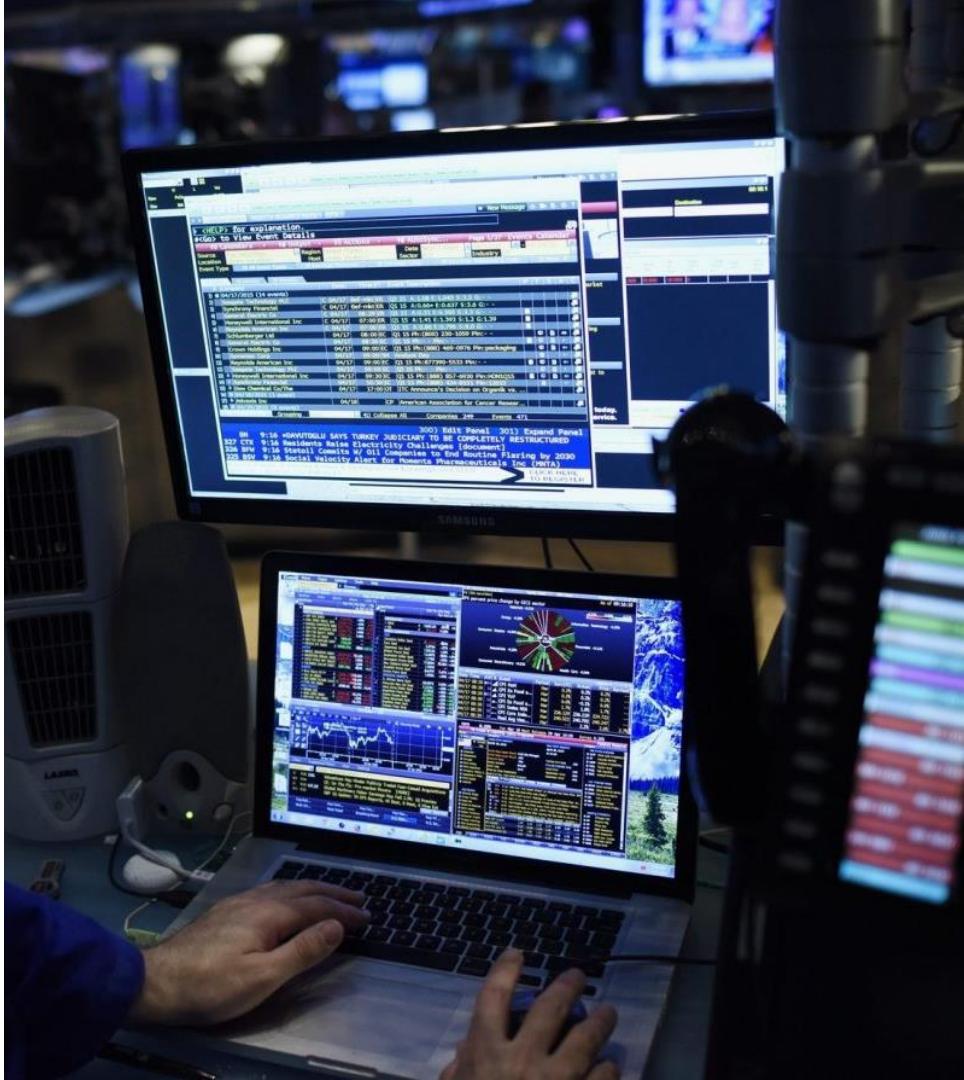
Massive global database of real-time and historical financial data – including trading, currency, supply chain, news, economic forecasts and more.

Bloomberg Terminal

The core revenue generator – implementing a client-side, physical terminal to connect to the Bloomberg Data Center.

Bloomberg Market Data Feed (B-PIPE)

B-PIPE enables real-time access to streaming data – allowing client applications access to the Bloomberg Data Center via the BLPAPI (C++, .NET, Java, Python).



Project Goal

Financial Services

Platform Requirements

**Architect a state-of-the-art, automated,
operationalized Cloud Data Analytics Platform**

Project Goal

Financial Services

Platform Requirements



Azure Databricks

Bloomberg

Real-Time Ingestion



- Stream + Batch:
- Ticker data
 - Share price
 - Company news
 - Supply chain updates
 - Historical data

Secured API Integration



ExpressRoute
Integration to
Bloomberg API

Rapid ETL



High-speed Data processing
(Stream + Batch loads)

DevOps



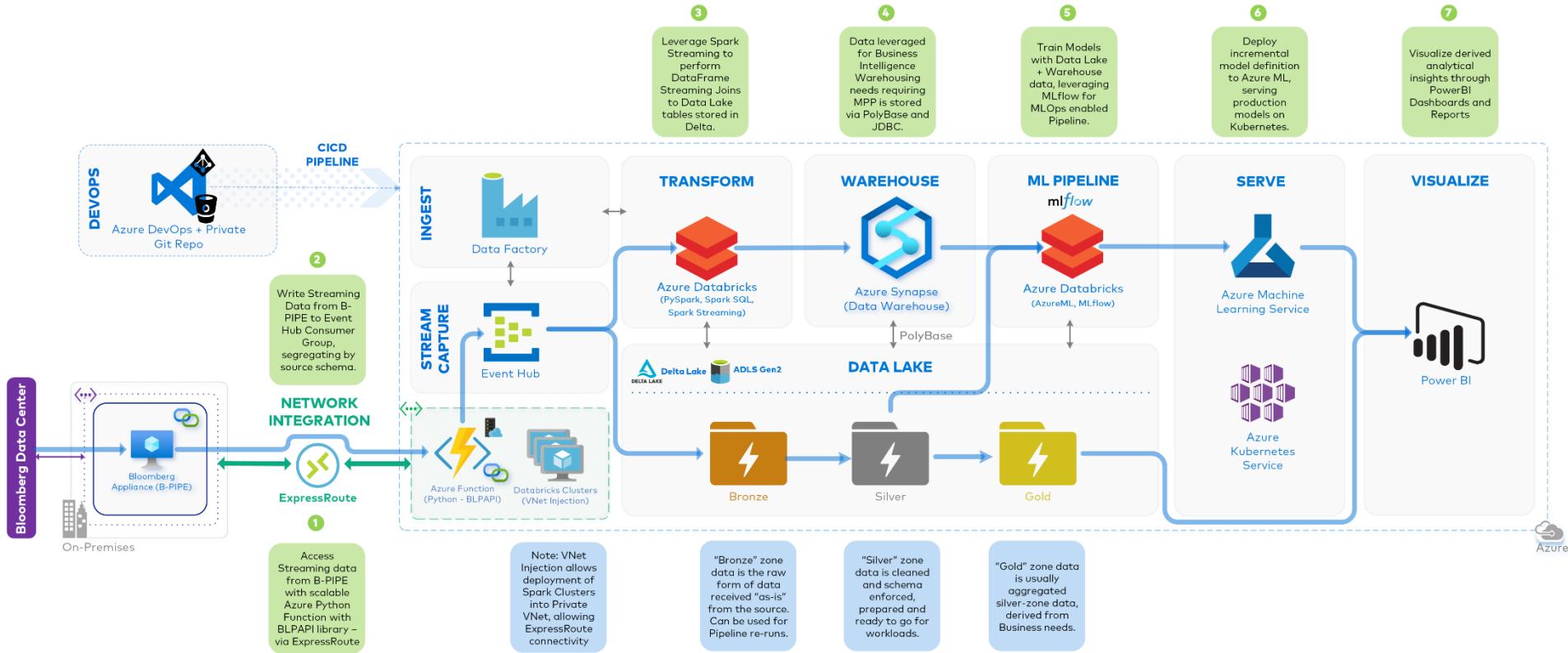
Continuous Integration
Continuous Delivery
(CICD) enabled Data Pipeline

Machine Learning Pipeline



MLOps enabled ML Pipeline capable of real time scoring and incremental model updates

Azure Solution Architecture



ML POC Use Case

Machine readable news with Bloomberg

Proof of Concept Objective

**Use Machine readable news from Bloomberg to explore
ML driven market trading strategy**

Machine readable news

The idea is – Bloomberg provides access to real-time event-driven news feeds in the form of a structured dataset (i.e. consistently formatted, readership stats, sentiment, ticker tagged e.g. STAN LN).

The goal is – train NLP* driven Machine Learning models to filter news and generate Buy/Sell signals.

*NLP: training ML models to process and analyze large amounts of language data.



On August 7th 2012, Standard Chartered stocks took a sharp decline.

STAN LN GBP ↑ AT 1445.50 -4.50 L1445.00 / 1445.50L 1139x852		Message
90 New Search 90 Actions 98 Settings Feedback		Page 1 News
Standard Chartered PLC Sources My Sources Lang My 08/07/12		Sort By Time
Standard Char... X		
Topics My News Securities Search Standard Cha...		
Ticker Lists UKX SXSE INDU +Customize		
In the News Apple Inc Royal Mail PLC Rhoen Kliniku... Access S... Goldman Sac... Deutsche Pos... Assicurazioni... Vivendi SA Banco Monte ... Aker Solution... People Sources		
1) S Korea set to overhaul key lending rate 2) Standard Chartered Probe Said to Require Up to \$700 Million (1) 3) N.Y.'s Standard Chartered Claims Surprised Regulators: NYT Link 4) N.Y.'s Standard Chartered Claims Surprised Regulators: NYT 5) CFS Daily Commet 8 August 2012 6) Standard Chartered New York License Threat Worries Investors (1) 7) GBN: -U.S. Treasury and Fed angered by NY action against Stan Chart. 8) The Mirror: Serious peril looms if US licence lost; ANALYSIS.(News) 9) StanChart Takes 26% Stake in Green Infra: Line Link 10) Stanchart Could Absorb Penalty Equal to Iranian Fees Earned: JPM 11) Standard Chartered's USD clearing business at stake 12) Hong Kong Monetary Authority Reviewing U.S. Order on StanChart 13) India Offshore Work Probed After Standard's 'Sham': Reuters Link 14) Standard Chartered rejects Iran deals claim (Transcript) 15) StanChart stock battered amid claims of illegal Iran deals 16) Stanchart shares tumble amid probe 17) Stanchart ran rogue Iran unit: US regulators 18) US presidential election and fears over Israeli intentions prompt tougher s... 19) Comment: StanChart will not be the last bank to be stung by the US		FTI 08/12 BN 08/12+ BN 08/12 BFW 08/12 CRH 08/12 BN 08/12 RE3 08/12 NS3 08/12 BN 08/12 BFW 08/12 EJ1 08/12 BFW 08/12 BFW 08/12 AUB 08/12 STS 08/12 BTM 08/12 BTM 08/12 FTI 08/12 FTI 08/12+
700 CTX 5:55 Tibet's 1st extradosed bridge opened to traffic 699 BFW 5:55 Italian Equity Movers: Banco Popolare, Generali, Hera, Sorin 698 NS3 5:55 Llanelli Star: Fair brings festive feel 697 NS3 5:55 Llanelli Star: BuyLocal: In the Market		

Viewing news from the timestamp shows root cause

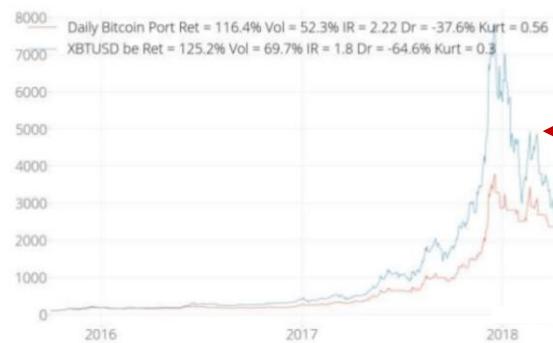
Summarized high level approach

#	To do	Approach	Technology
1	Batch/Stream ingestion pipeline of ticker data (historical/real-time news + stock)	<p>Pull historical/real-time data from B-PIPE into Data Lake and store in "Bronze" Zone.</p> <p>Create Data Ingestion pipeline for Stream and Batch processing.</p>	 Azure Databricks  Delta Lake  Spark Streaming
2	Create Language Data Cleansing Pipeline for news data, and curate into time series format	Leverage Text Pre-Processing tools to perform cleaning operations (stemming, lemmatization, chunking/chinking) to filter noise and structure news data	 Koalas  BeautifulSoup  spaCy
3	Perform tokenization, topic modelling and sentiment analysis	Convert cleaned, structured text into mathematical (vectorized) representation, leveraging NLP specific mathematical techniques (TFIDF, TSNE, LSTM etc.).	 Azure Machine Learning  TensorFlow  Keras
	Create indicators and apply trading rules	Convert into buy/sell signals and/or add other trading factors (e.g. carry)	 Jupyter  Apache Spark  MLlib
4	Score models on test data and deploy best performing model to serving layer	<p>Perform experiments and track parameters/results using Mlflow.</p> <p>Deploy containerized models to Azure Kubernetes using Azure ML version control.</p>	 mlflow  Azure Machine Learning  kubernetes
5	Iterative Model Definition Updates	Connect B-PIPE Streaming pipeline to ML Pipeline for iterative scoring and incremental model updates.	 Azure Databricks  Delta Lake  mlflow  Azure Machine Learning

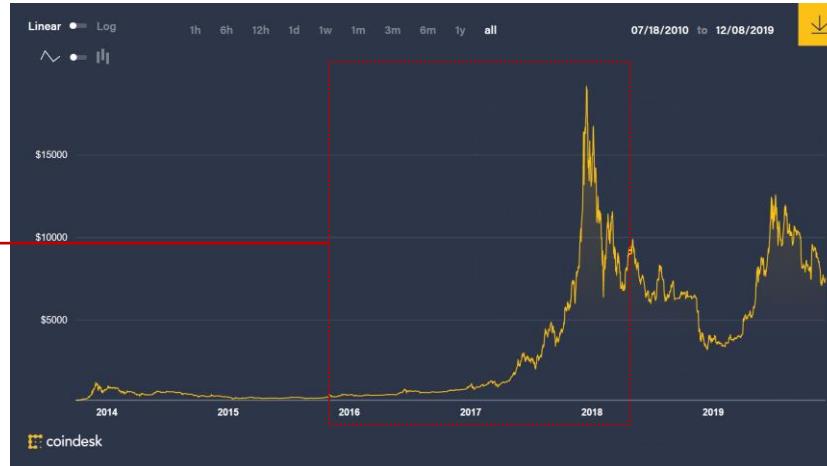
Trading Strategy Output from News Data

```
if __name__ == '__main__':
    model = TradingBitcoinDaily()
    model.construct_strategy()
    model.plot_strategy_group_benchmark_pnl()
```

Daily Bitcoin



Bitcoin (long) actual price VS.
Trading Strategy from NLP



Bitcoin price: 2010 to 2019

Useful Resources

Self-Study/Deep-Dive material

Intelligent Harry Potter Search Engine

Building a BERT-powered Natural Language Processing model, trained on Harry Potter books to answer context-specific questions, served in a containerized web-app.



Answer

Seeker

Passage Context

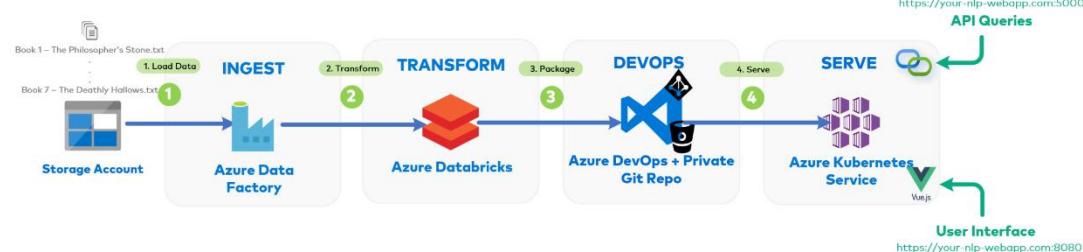
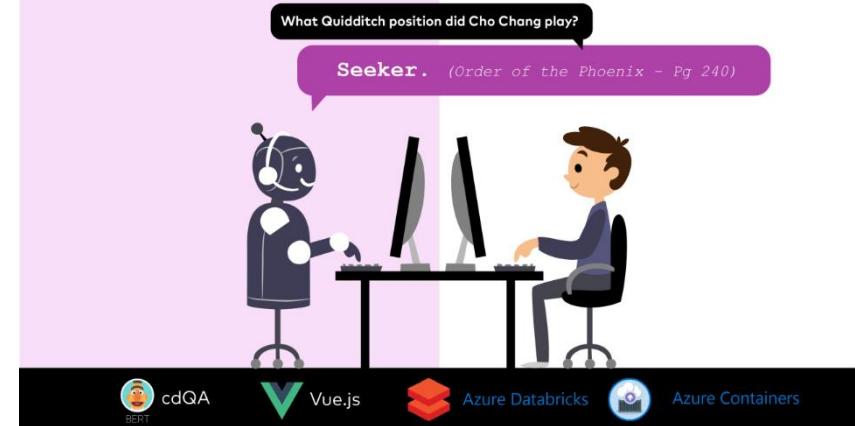
though, Stinksap's not poisonous," he added nervously, as Harry spat a mouthful onto the floor. At that precise moment the door of their compartment slid open. "Oh ... hello, Harry," said a nervous voice. "Um ... bad time?" Harry wiped the lenses of his glasses with his Trevor-free hand. A very pretty girl with long, shiny black hair was standing in the doorway smiling at him: Cho Chang, the Seeker on the Ravenclaw Quidditch team. "Oh ... hi," said Harry blankly. "Um ..." said Cho. "Well ... just thought I'd say hello ... 'bye then." She closed the door again, rather pink in the face, and departed. Harry slumped back in his seat and groaned. He would have liked Cho to discover him sitting with a group of very cool people laughing their heads off at a joke he had just told; he would not have chosen to be sitting with Neville and Loony Lovegood, clutching a toad and dripping in Stinksap. "Never mind," said Ginny bracingly. "Look, we can get rid of all this easily." She pulled out her wand. "Scour gif YY" The Stinksap vanished. "Sorry," said Neville again, in a small voice. Ron and Hermione did not turn up for nearly an hour, by which time the food trolley had already gone by. Harry, Ginny, and Neville had finished their Pumpkin Pasties and were busy swapping Chocolate.

Original Document

The Order of the Phoenix - Page 240

Link

<http://bit.ly/azure-harry-potter>



Spark Certification Study Guide

While studying for the [CRT020: Apache Spark 2.4 with Python 3 Developer Certification](#), I compiled all my study notes into a thorough study-guide required to pass the exam:



[Link](#)

<http://bit.ly/spark-study-guide>

Spark Developer Certification - Comprehensive Study Guide (python)

What is this?

While studying for the [Spark certification](#) exam and going through various resources available online, I thought it'd be worthwhile to put together a comprehensive knowledge dump that covers the entire syllabus end-to-end, serving as a Study Guide for myself and hopefully others.

Note that I used inspiration from the [Spark Code Review guide](#) - but whereas that covers a subset of the coding aspects only, I aimed for this to be more of a comprehensive, one stop resource geared towards passing the exam.

Awesome Resources/References used throughout this guide

References

- [Spark Code Review used for inspiration](#): <https://databricks-prod-cloudfront.cloud.databricks.com/public/402/eecb02393dea0a7f14173bcfc524954772526824/5218887830675/0/1123846766950497/latest.html>
- [Syllabus](#): <https://academy.databricks.com/exam/crt020/python>
- [Data Scientists Guide to Apache Spark](#): <https://dive.google.com/open?id=17KMSlwgMyQ8uv/tbowznlLB64Oer9T>
- [JVM Overview](#): <https://www.jaworld.com/article/3272244/what-is-the-jvm-introducing-the-java-virtual-machine.html>
- [Spark Runtime Architecture Overview](#): <https://reco.rent.manning.com/running-spark-an-overview-of-sparks-runtime-architecture>
- [Spark Application Overview](#): https://docs.claudera.com/documentation/enterprise/5-6-x/topics/cdh_ig_spark_apps.html
- [Spark Architecture Overview](#): <https://queriot.com/en/series/spark-architecture-overview-clusters-jobs-stages-tasks-etc>
- [Mastering Apache Spark](#): <https://packtbooks.gitbooks.io/mastering-apache-spark/spark-scheduler/Stage.html>
- [Manually create DFs](#): <https://medium.com/@mpowern/manually-creating-spark-dataframes-b14da6906393>
- [PySpark SQL docs](#): <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html>
- [Introduction to DataFrames](#): <https://databricks.com/doc/facts/datasets/introduction-to-dataframes-python.html>
- [PySpark UDFs](#): <https://spark.apache.org/docs/latest/api/python/pyspark.html>
- [ORC File](#): <https://wikispark.org/confluence/display/livedLanguageManual-ORC>
- [SQL Server Stored Procedures from Databricks](#): <https://databricks.net/blog/2018/10/12/executing-sql-server-stored-procedures-on-databricks-pyspark>
- [Repartition vs Coalesce](#): <https://stackoverlow.flow/questions/3161097/sql-repartition-vs-coalesce>
- [Partitioning by Columns](#): <https://packtbooks.gitbooks.io/mastering-spark-sql/spark-sql-dynamic-partition-inserts.html>
- [Bucketing](#): <https://packtbooks.gitbooks.io/mastering-spark-sql/spark-bucketing.html>
- [PySpark GroupBy and Aggregate Functions](#): <https://hendra-henriawan.github.io/pyspark-groupby-and-aggregate-functions.html>
- [Spark Quickstart - 1](#): https://spark.apache.org/docs/latest/quick_start.html
- [Spark Caching - 1](#): <https://stackoverlow.com/questions/44958869/where-does-df-cache-is-stored>
- [Spark Caching - 2](#): <https://stackoverlow.com/questions/45558869/where-does-df-cache-is-stored>
- [Spark Caching - 3](#): <https://changhoilee.com/2018/04/24/pyspark-dataframe-basics/>
- [Spark Caching - 4](#): <https://spark.apache.org/docs/latest/rdd-programming-guide.html#rdd-persistence>
- [Spark Caching - 5](#): https://www.tutorialspoint.com/pyspark/pyspark_storagelevel.htm
- [Spark Caching - 6](#): <https://spark.apache.org/docs/2.1.2/python/modules/pyspark/storagelevel.html>
- [Spark SQL Function Examples](#): <https://spark.apache.org/docs/2.3.0/api/python/index.html>
- [Spark Built-in Higher Order Functions Examples](#): <https://docs.databricks.com/static/notebooks/apache-spark-2-4-functions.html>
- [Spark SQL Timestamp conversion](#): <https://docs.databricks.com/static/notebooks/mstamp-conversion.html>
- [RegEx Tutorial](#): <https://medium.com/factory-mind/reg-ex-tutorial-a-simple-cheatsheet-by-examples-ef4dc1c2d85>
- [Rank VS Dense Rank](#): <https://stackoverflow.com/questions/44968912/difference-in-dense-rank-and-row-number-in-spark>
- [SparkSQL Windows](#): <https://databricks.com/blog/2015/07/15/introducing-window-functions-in-spark-sql.html>
- [Spark Certification Study Guide](#): <https://github.com/vivek-bombatkar/Databricks-Apache-Spark-2X-Certified-Developer>

Resources

PySpark Cheatsheet

Python For Data Science Cheat Sheet

PySpark - SP basics
Learn Python for data science interactively at www.DataCamp.com



PySpark & Spark SQL

Spark SQL is Apache's module for working with structured data.

Initializing SparkSession

A SparkSession can be used to create DataFrame, register DataFrame as tables, execute SQL queries, cache tables and read parquet files.

```
>>> spark = SparkSession.builder.appName("SparkSQL").getOrCreate()
```

Creating DataFrames

From RDDs

```
>>> from pyspark import types import *  
Infer Schema
```

```
>>> spark.read.options(header='true').csv('people.csv')  
>>> lines = sc.textFile("people.txt").map(lambda l: l.split(",")*)
```

```
>>> people = lines.map(lambda p: Row(name=p[0], age=int(p[1]),  
    gender=p[2]))
```

Specify Schema

```
>>> schema = StructType([StructField("name", StringType(), True),  
    StructField("age", IntegerType(), True),  
    StructField("gender", StringType(), True)])
```

Register

```
>>> df = people.createDataFrame(people, schema)  
>>> df.registerTempTable("people")
```

Duplicate Values

```
>>> df.dropDuplicates()
```

Queries

```
>>> from pyspark import SQLContext
```

```
>>> df.select("*").show()
```

```
>>> df.select("name", "age").show()
```

```
>>> df.select("name", "age").distinct().show()
```

```
>>> df.select("name", "age").groupby("name").count().show()
```

```
>>> df.select("name", "age").groupby("name").count().sort("count", ascending=False).show()
```

Understanding MLOps with Azure Databricks

"Because the focus has always been on the model's prediction capabilities and not how the model actually gets there, we've poured thousands of hours and developed hundreds of models, and now managing this beast of an environment is a never-ending battle every single day when you actually start using this pipeline to predict mission-critical features."

Link

<http://bit.ly/mlops-databricks>



Understanding MLOps with Azure Databricks



Operationalized ML Pipelines with Delta Lake, MLflow and Docker Containers in Azure Databricks

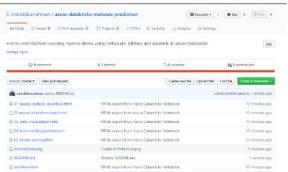
Understanding MLOps with Azure Databricks

Published on November 25, 2019 | [Edit article](#) | [View stats](#)



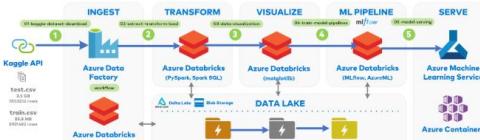
Ki Rahman
Data, MLOps & DevOps | Microsoft Certified
Azure Expert | Machine Learning in ... See More

Github repo for Project



All the notebooks referenced are available in my Public Github repo [here](#) in HTML format, that can be directly imported into your Databricks environment (or viewed on a browser).

Architecture Diagram



The architecture above is labelled with the corresponding notebooks that correspond to each high level flow.

MLFlow end-to-end demo

Our Demo for today



An end-to-end demo showcasing Mlflow **Tracking**, **Projects** and **Models**, using an Airbnb dataset.

We will finish by deploying a Random Forest model into Azure ML, and querying served model with Postman against the deployed Docker Container – mimicking a Production App.

Link

<http://bit.ly/mlops-mmai>

mlops-queens-university

Big Data and MLOps content presented at Queen's University MMAI Workshop

This repo contains the content I presented at the Queen's University MMAI Workshop ([link](#)).

Primary components

We're going to demo the three core components that make up mlflow:

mlflow
TRACKING

Record and query experiments: code, data, config, results

databricks.com/mlflow

mlflow
PROJECTS

Packaging format for reproducible runs on any platform

mlflow.org

mlflow
MODELS

General model format that supports diverse deployment tools

github.com/mlflow

twitter.com/MLflow

Note: you don't have to use all three, each feature can be used independently.

Tracking

This allows us to log all aspects of the ML process - like different hyperparameters we tried, evaluation metrics, as well as the code we ran - alongside other arbitrary artifacts such as test data.

This also provides a *leaderboard-style UI* that makes it easy to see which model performed the best.

Projects

These are all about **reproducibility** and **sharing**. They combine GIT, the environment/model framework, either conda or docker and the specification that makes the code re-runnable.

Models

An abstraction that allows us to **create/export models** from any open source framework via the *Tracking* and *Projects* abstractions. We can also export them to a standard format that can be deployed to any number of systems. Since most deployment systems use some sort of container based solution (e.g. AzureML or Sagemaker), models make easy deployments to these systems - or we can deploy directly to Kubernetes or Azure Container Registry.

Agenda

In this notebook we will demonstrate the following topics:

Before the Demo

Questions?



Raki Rahman

Solution Architect

slalom | Data & Analytics

<https://rakirahman.com>