



Automating Apache Cassandra Operations with Apache Airflow

Go beyond cron jobs to manage ETL, Data Hygiene, Data Import/Export

Rahul Xavier Singh Anant Corporation

Cassandra Days - London 10/12 | Amsterdam 10/13

We help **platform owners**
reach **beyond** their **potential**
to serve a **global customer**
base that demands

Everything, Now.

☰

We design with our
Playbook, build with our
Framework, and manage
platforms with our **Approach**
so our clients

Think & Grow Big.



Customer Success



Canon

CANON MEDICAL



Deloitte.

Teachstone



intuit



UKG

Our purpose is people



**McKinsey
& Company**



verizon

Export Now



accenture

OUTFRONT
media



AHRI
AIR-CONDITIONING,
HEATING,
& REFRIGERATION INSTITUTE



ONE TO WORLD

Startup Hoyas

AROUND THE RINGS



IBM



U.S. Food Imports

RAND MC NALLY

**Town and
Country Bank**

M
University of Michigan
Health System



Challenge Solutions Subscriptions

Business



Playbook

Platform



[Data] Services Catalog

Technology



Framework

Management



Approach

We offer Professional Services to engineer Solutions and
offer Managed Services to clients where it makes sense, after an
Assessment

Business / Platform Dream

Enterprise Consciousness :

- People
- Processes,
- Information
- Systems

Connected / Synchronized.

Business has been chasing this dream for a while. As technologies improve, this becomes more accessible.

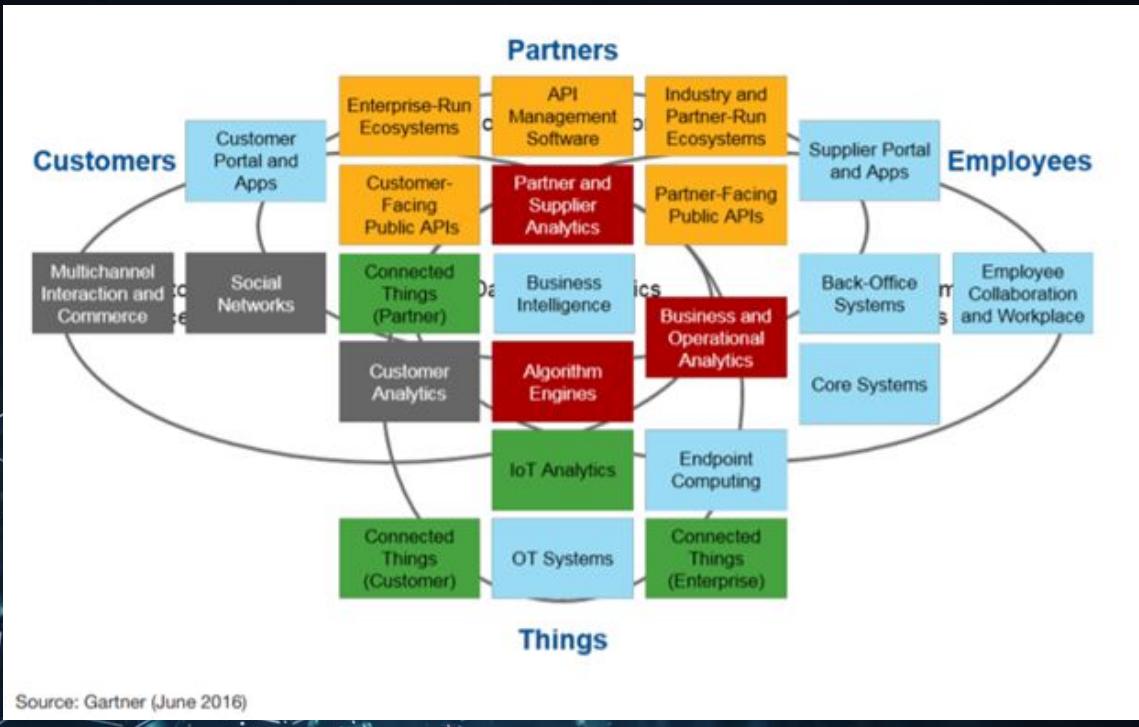


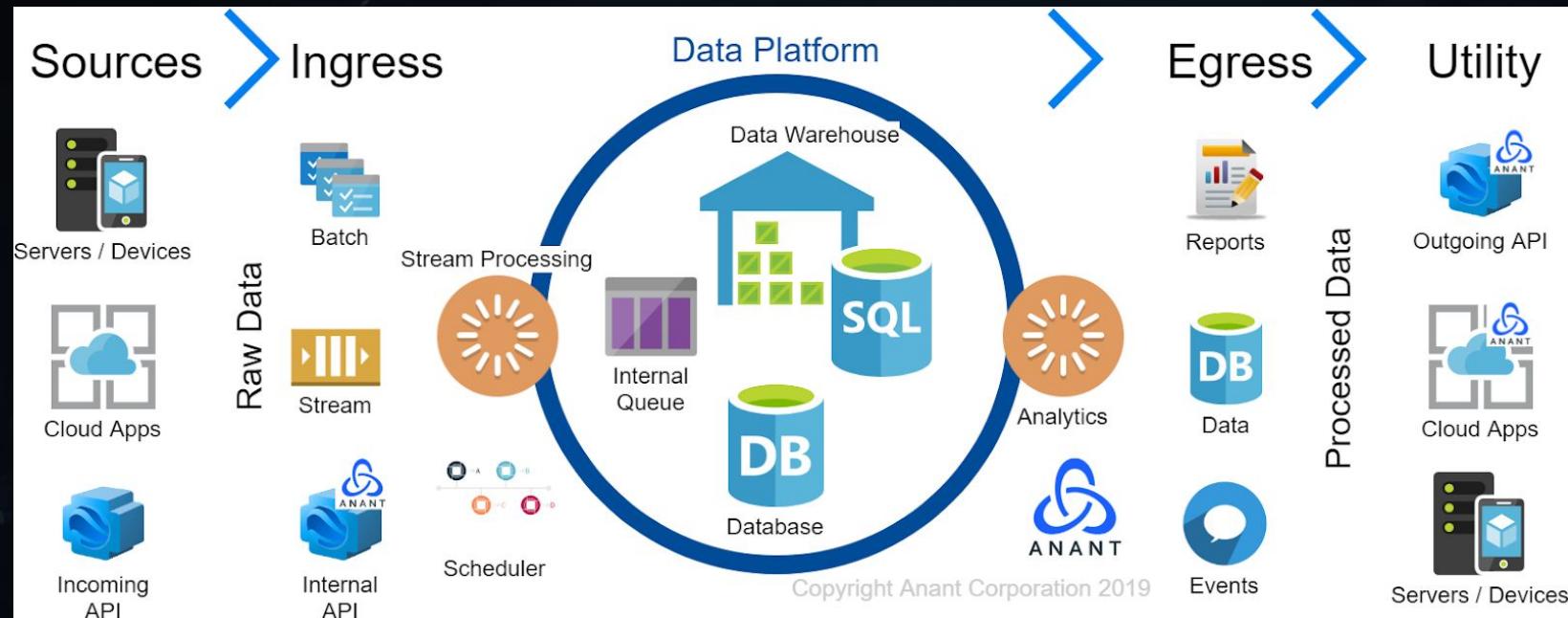
Image Source: Digital Business Technology Platforms, Gartner 2016



Modern Open Data Platform



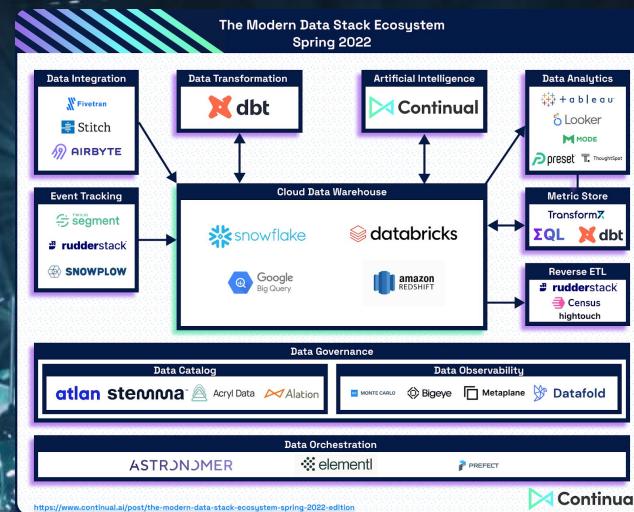
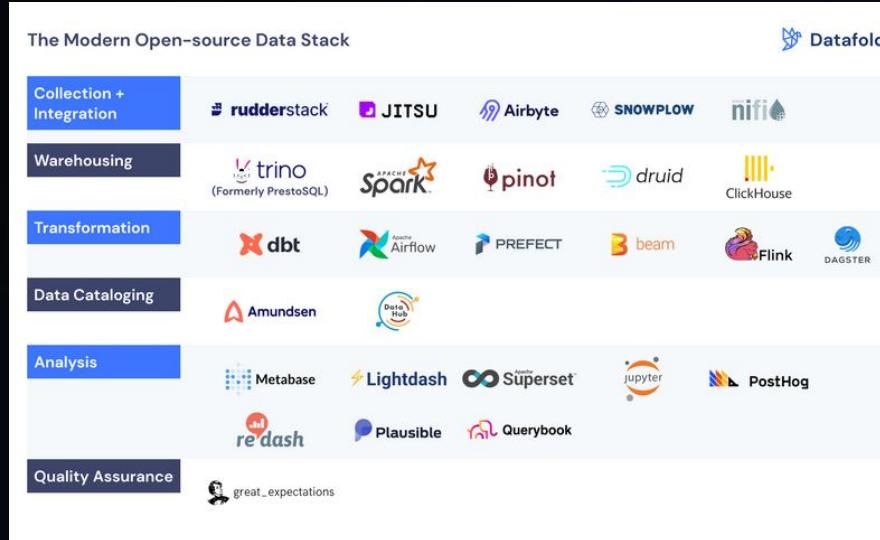
Generic Data Platform Operations



How do you choose from the landscape?

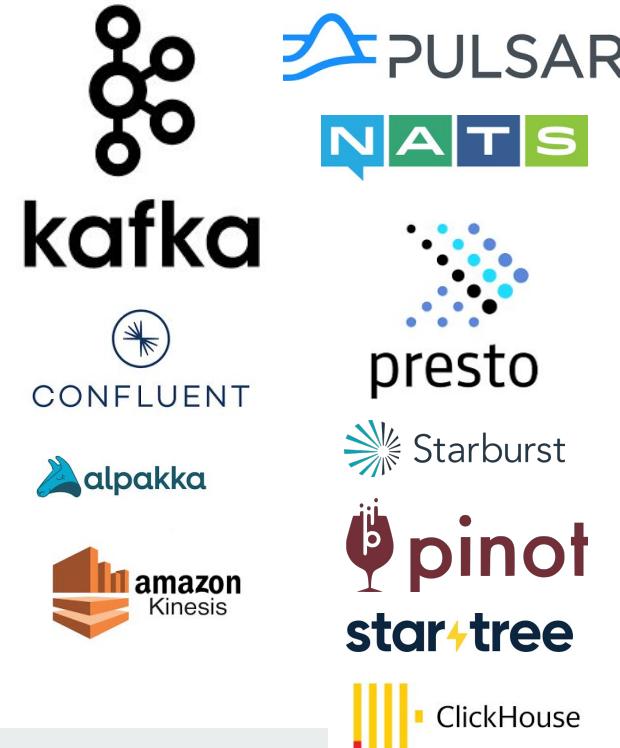


So Many Different “Modern Stacks?”



Lots of “reference” architectures available. They tend not to think about the speed layer since they are focusing on analytics. **Many don’t mention realtime databases... but we can learn from them.**

Distributed Realtime Components



To create globally distributed and real time platforms, we need to use distributed realtime technologies to build your platform. Here are some. Which ones should you choose?

Playbook / Framework



Framework Components



- Major Components
 - Persistent Queues (RAM/BUS)
 - Queue Processing & Compute (CPU)
 - Persistent Storage (DISK/RAM)
 - Reporting Engine (Display)
 - Orchestration Framework (Motherboard)
 - Scheduler (Operating System)
 - Strategies
 - Cloud Native on Google
 - Self-Managed Open Source
 - Self-Managed Commercial Source
 - Managed Commercial Source

Customers want options, so we decided to create a Framework that can scale with whatever Infrastructure and Software strategy they want to use.

Category	Streams / Queues	Processors	Storage	Reporting	Orchestration
Open Source	Apache Kafka Apache Pulsar Apache Nifi RabbitMQ ActiveMQ CLR Threads JVM Threads JS Threads <u>Alpakka/Akka</u> Spark (Scala, Java)	Kafka Connect Kafka Streams Pulsar.io Apache Nifi MySQL Redis Apache Iceberg Delta Lake	Apache Cassandra Apache Solr Elasticsearch MySQL Redis Apache Iceberg Delta Lake	Redash Metabase Zeppelin Jupyter + SQL <u>SparksSQL</u> Presto	Terraform Ansible Puppet Chef Docker Kubernetes Airflow
Google Native	Google Pub/Sub	Google Cloud Tasks Dataflow <u>DataProc</u>	Google Filestore Bigtable Spanner	Google BigQuery	Google CloudRun Build GKE
Commercial	MuleSoft	DSE Analytics	DSE Cassandra	Tableau	Terraform
Commercial	Confluent	Confluent	DSE Search	<u>Microstrategy</u>	Ansible
Open Source	Luna Pulsar	MuleSoft	DSE Graph	+ +	puppet
	Lenses	Talend	Databricks		Chef
		Serverless	Delta Lake		Liquibase
Managed	Confluent	Databricks		DSE	
Open Source	Cloud Aiven Cloud			AlwaysOnSQL Snowflake Starburst Presto Dremio	DSE OpsCenter LCM

Playbook for Modern Open Data Platform



Platform Design

Discovery (Inventory)

- People
- Process
- Information (Objects)
- Systems (Apps)

Evaluate Framework

User Experience

- No-Code/Low Code Apps/Form Builders
- Automatic API Generator/Platform
- Customer App/API Framework

Cloud

- Public
- Private
- Hybrid

DevOps

- Infrastructure as Code
- Systems Automation
- Application CICD

Data

- Data:**Object**
- Data:**Stream**
- Data:**Table**
- Data:**Index**
- Processor:**Batch**
- Processor:**Stream**

DataOps

- ETL/ELT/EtlT
- Reverse ETL
- Orchestration

Execute Approach

Architecture (Design)

- Cloud
- Data
- DevOps
- DataOps

Engineering

- Configuration
- Scripting
- Programming

Operation

- Setup / Deploy
- Monitoring/Alerts
- Administration



Framework



ASTRONOMER



Jenkins



ANSIBLE

Terraform

kubernetes

ANANT



Google Cloud

Azure

amazon web services™



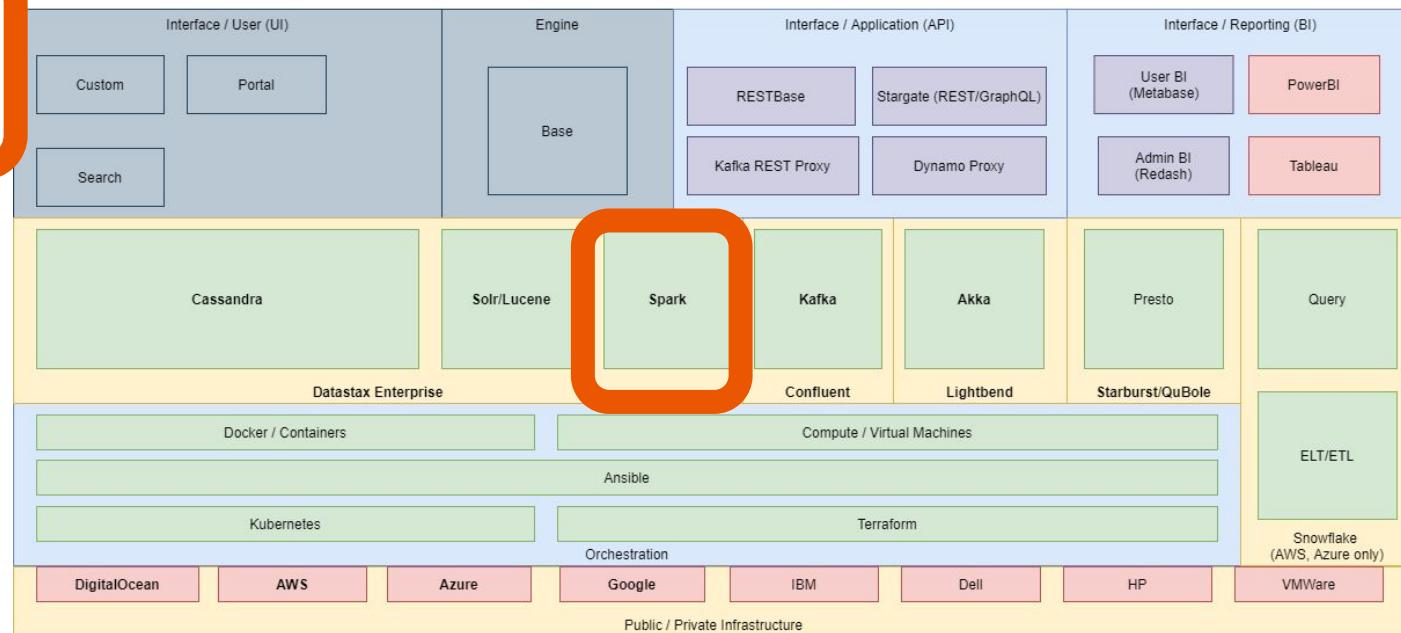
LENSES



databricks



yugabyteDB



presto

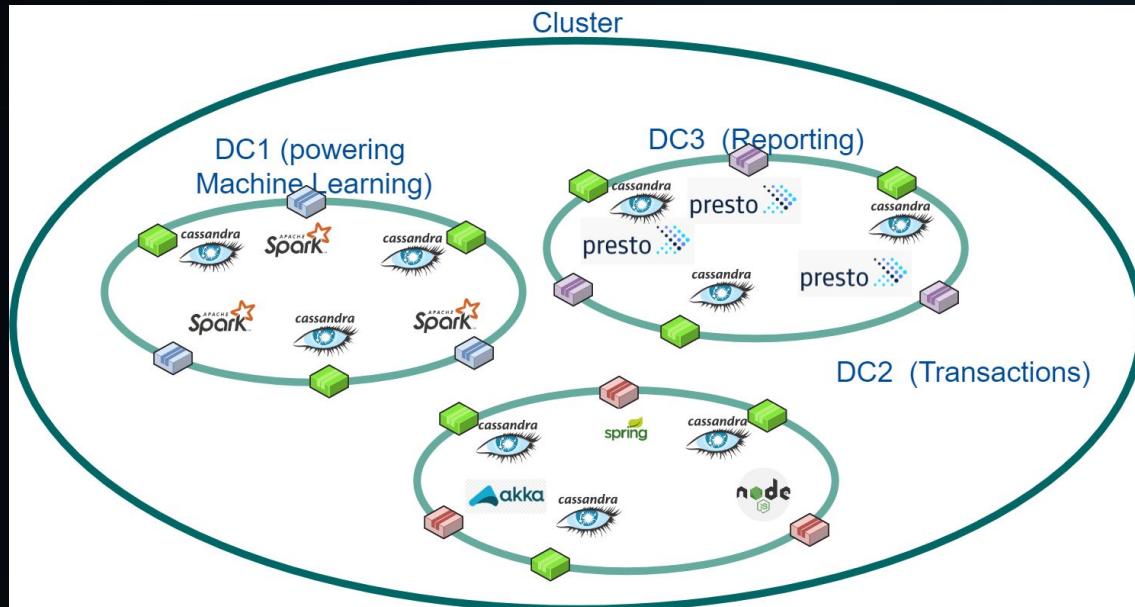


starburst

pinot

star tree

Thinking about Cassandra as a Data Fabric



XDCR: Cross datacenter replication is the ultimate data fabric.

Resilience,
performance,
availability, and scale.

Made widely available
by Cassandra and
Couchbase



Data Modernization / Automation / Integration



great_expectations



Querybook



In addition to vastly scalable tools, there are also modern innovations that can help teams automate and maximize human capital by making data platform management easier.



Playbook / Approach



Approach

Setup	Training	Administration	Customization	Knowledge Management
Dev / Ops / Arch Team	Dev / Ops / Arch Team	Dev / Ops / Arch Team	Dev / Ops / Arch Team	Dev / Ops / Arch Team
Understanding current / future needs	Discovery : Understand skill gaps	Service management	Self Service process	Platform Architecture
Mapping a next viable architecture	Curriculum Planning	Service transition	Testing environments	Platform Operations Runbook
Deciding on Infrastructure	Architecture Training	Incident management	Service transition	Platform Implementation Guide
Install & Stress Test	Implementation Training	Data / Disaster Recovery Plan	Service strategy	Issue resolution knowledge
Orchestration Automation	Operations Training	Continuous monitoring	Continuous testing	Public / Internal knowledge sources
Continuous Integration & Delivery	Ongoing mentoring	Continuous testing	Change management	Continuous knowledge creation



ANANT

Copyright Anant Corporation 2019



Apache **Airflow** +
Apache **Spark** +
Spark Python/Scala/Java/R +
Airflow Python **DAG** =

DataOps for Apache Cassandra
Good enough for rock and roll.



Apache Airflow



- A tool for scheduling and automating workflows and tasks
- Good for automating repeated processes
 - Common ETL tasks
 - Machine learning model training
- Write workflows in Python
 - Tools interactable via Python should work as well
 - Define dependencies for different sections of workflows
 - Workflows are DAG of tasks
- Schedule workflows or execute the processes by hand
 - Cron-like syntax or frequency tags
- Monitor tasks and collect/view logs



Apache Spark



- Unified analytics engine for large-scale data processing.
- Achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.
- Offers over 80 high-level operators that make it easy to build parallel apps. And you can use it interactively from the Scala, Python, R, and SQL shells.
- Powers a stack of libraries including SQL and DataFrames, MLlib for machine learning, GraphX, and Spark Streaming. You can combine these libraries seamlessly in the same application.
- You can run Spark using its standalone cluster mode, on EC2, on Hadoop YARN, on Mesos, or on Kubernetes. Access data in HDFS, Alluxio, Apache Cassandra, Apache HBase, Apache Hive, and hundreds of other data sources.

Big Data Options

Coldish

- S3
- HDFS
- ADLS
- GFS

Warm

- Hive / *
- Data Warehouse
- Data Lakehouse

Hot

- Cassandra*
- Datastax*
- Scylla*
- Yugabyte*
- Mongo
- REDIS
- ...

Hot*

- Astra*
- Scylla Cloud*
- YugaByte Cloud*
- Azure CosmosDB*
- AWS Keyspaces*
- AWS Dynamo
- Google BigTable
- ...



* PSSST. These all use CQL!!!

Cleaning Big Data: Same \$hit Different Day

Data Cleaning as part of Data Engineering

- Step 1: Remove duplicate or irrelevant observations
- Step 2: Fix structural errors
- Step 3: Filter unwanted outliers
- Step 4: Handle missing data
- Step 5: Validate and QA

Data Cleaning after the Fact

- Enforce a custom data retention policy (TTL)
- Enforce GDPR / Right to be Forgotten
- Move application, customer, user from one system to another
- Remove x “versions” or “revisions” of data
- Remove test data from a stress test



<https://www.tableau.com/learn/articles/what-is-data-cleaning>

Cleaning Big Data: In SQL

```
9  
10  DELETE FROM [dbo].[SalesData]  
11  WHERE CustomerId IN  
12  (SELECT TOP 10 CustomerId  
13  FROM [dbo].[SalesData]  
14  ORDER BY OrderDate ASC);  
15  GO  
16  
17
```

100 %

Messages

(13 rows affected)

Data Cleaning in SQL

Find what do you want to delete.
Delete it.



Cleaning Big Data: In Spark SQL

SQL



```
> DELETE FROM events WHERE date < '2017-01-01'

> DELETE FROM all_events
  WHERE session_time < (SELECT min(session_time) FROM good_events)

> DELETE FROM orders AS t1
  WHERE EXISTS (SELECT oid FROM returned_orders WHERE t1.oid = oid)

> DELETE FROM events
  WHERE category NOT IN (SELECT category FROM events2 WHERE date > '2001-01-01')
```

Deletes the rows that match a predicate. When no predicate is provided, deletes all rows.

This statement is only supported for Delta Lake tables.

Data Cleaning in SPARK SQL

- What do you want to delete.
Delete it.

WARN: Doesn't work with all data in
Spark SQL. Only if connector supports
Table Delete

<https://docs.databricks.com/spark/latest/spark-sql/language-manual/delta-delete-from.html>
<https://spark.apache.org/docs/latest/api/java/org/apache/spark/sql/connector/catalog/SupportsDelete.html>

Cleaning Big Data: Cleaning data in Spark / SQL

Data Cleaning in Spark for Cassandra

- What do you want to delete.
- Delete it.

```
1 def deleteCategory(keyspace: String, table: String, filter: String): Long = {
2   import com.datastax.spark.connector.cql.CassandraConnector
3   val dseConnector = CassandraConnector(sc)
4   val documentDF = spark
5     .read
6     .format("org.apache.spark.sql.cassandra")
7     .options(Map( "table" -> table, "keyspace" -> keyspace))
8     .load
9   documentDF.createOrReplaceTempView("documents")
10  val deleteDocDF = spark.sql(s"SELECT * FROM documents WHERE $filter")
11  val deletedCount = deleteDocDF.count()
12  deleteDocDF.foreachPartition(partition => {
13    dseConnector.withSessionDo(session =>
14      partition.foreach(log =>
15        val delete = s"DELETE FROM \"$keyspace\".\"$table\" WHERE <primary key>=$log.getString(0);"
16        session.execute(delete)
17      )
18    )
19  })
20  deletedCount
}
```



<https://stackoverflow.com/questions/28563809/delete-from-cassandra-table-in-spark>



Cleaning Big Data: Deduping in Spark SQL

Deduping Data in Spark for Cassandra

- What do you want to dedupe.
- Do some deduping.
- What you want to delete.
- Delete it.

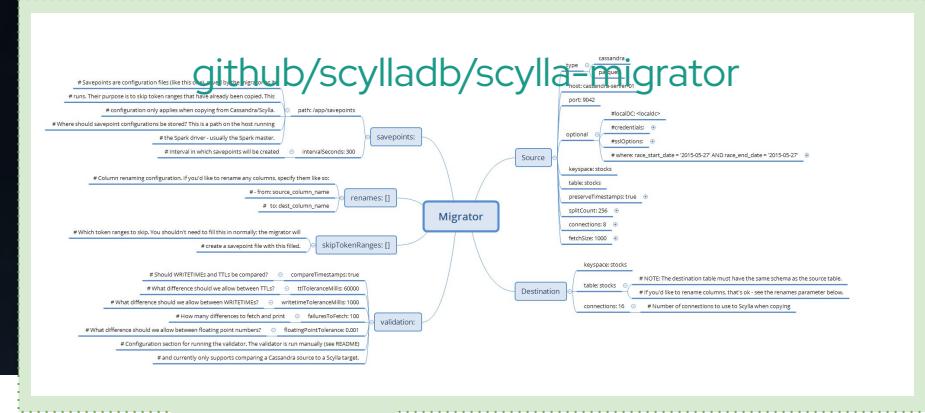
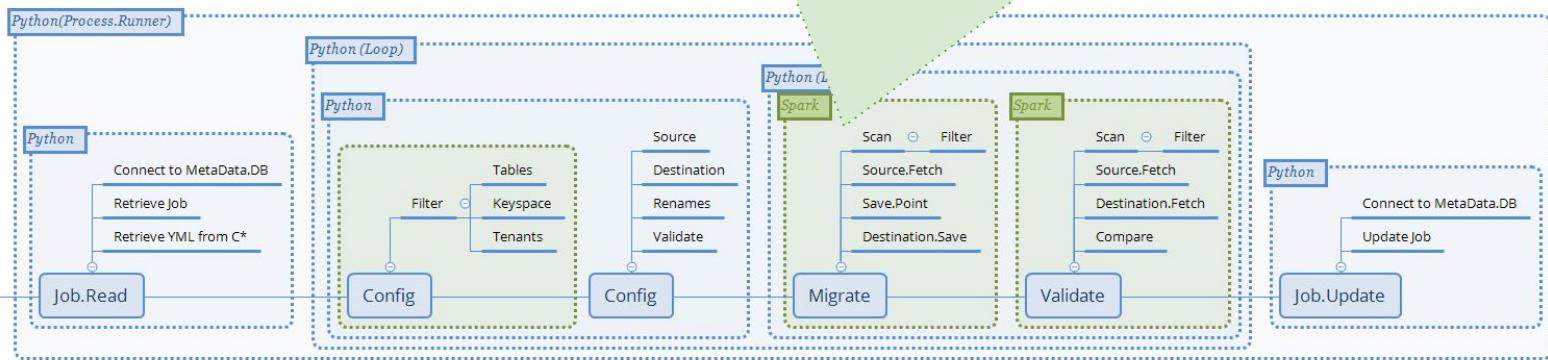
```
1  def dedupeTable(keyspace: String, table: String, columns: Array[String], filter: String): Long = {
2      import com.datastax.spark.connector.cql.CassandraConnector
3      val dseConnector = CassandraConnector(sc)
4      val documentDF = spark
5          .read
6          .format("org.apache.spark.sql.cassandra")
7          .options(Map( "table" -> table, "keyspace" -> keyspace))
8          .load
9      //TODO this category filter can be done using SPARK SQL, and you should be able to just pass it a SQL "where clause"
10     documentDF.createOrReplaceTempView("documents")
11     val docDF = spark.sql("SELECT * FROM documents "+filter)
12     val documentDF_temp = docDF.withColumn("comparable", concat_ws(", ", (columns.map(c => col(c)): _*)))
13     val countsDF = documentDF_temp.rdd.map(rec => (rec(9),1)).reduceByKey(_+_).map(rec => (rec._1.toString, rec._2)).toDF("Compared", "Count")
14     val joinedDF = documentDF_temp.join(countsDF, documentDF_temp("comparable") === countsDF("Compared"),
15         "left_outer").select(documentDF_temp("doc_id"), documentDF_temp("comparable"), countsDF("Count"))
16     val result = joinedDF.where("Count > 1")
17     val originals = result.groupBy("comparable").agg(first("doc_id").as("doc_id"), min("count").as("count"))
18     val to_Delete = (result.select("doc_id")).except(originals.select("doc_id"))
19     val deletedCount = result.count()
20     to_Delete.foreachPartition(partition =>
21         dseConnector.withSessionDo(session =>
22             partition.foreach{ log =>
23                 val delete = s"DELETE FROM "+keyspace+"."+table+" where doc_id="+log.getString(0) +";"
24                 session.execute(delete)
25             })
26     )
27     deletedCount
```



Airflow DAG to Migrate Cassandra Data



Spark-Migrator
-Process

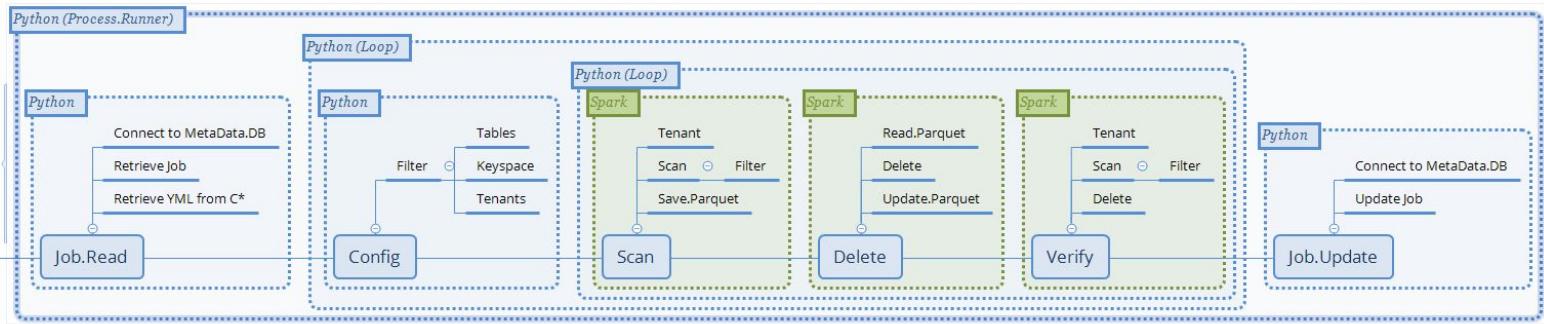


Airflow can help us take any data process, compiled, interpreted etc, coordinate the steps as a DAG ("Directed Acyclic Graph") and then to make it even more awesome, parametrize it either via the Airflow GUI or our own table somewhere.

Airflow DAG to Clean Cassandra Data



Spark-Cleanup
-Delete-Process



Since we write abstracted code, we can replace the “Migrator” process with a Delete, Dedupe, Validate. Whatever.

Airflow allows us to reuse conventions we set in a team for large scale operations, and most importantly... make it easy for people to run Data Operations like this without being Cassandra , Spark, Python experts.

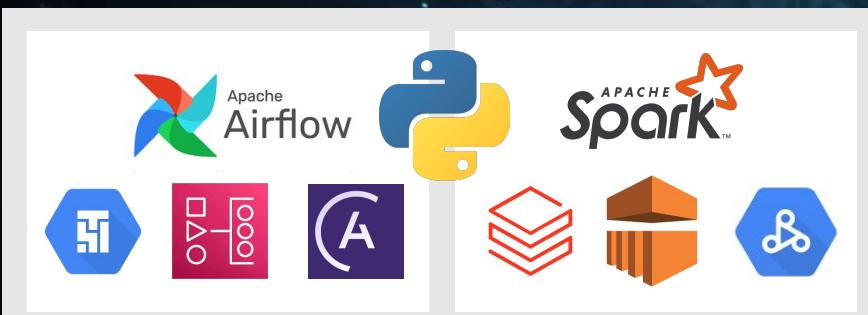
Considerations for Spark/Airflow Solution

Considerations for Spark

- Figure out if you want to **manage it / run it**.
- Not all Spark code is created equally.
- Not all Spark languages run the same.
- Compiled Jobs with input parameters can work better in the long run, less room for code drift.
- Don't let people do adhoc delete operations until and unless it's absolutely necessary.
- Who has access to it?

Considerations for Airflow

- Figure out if you are going to manage it / run it.
- Figure out for whom you are going to run it (Platform, Environment, Stack, App, Customer?)
- Not all DAGs just work. Sometimes they need tweaking across Environments, Stacks, Apps, Customers.
- The same DAG may fail over time. Need to watch execution times.
- Who has access to it?



Key Takeaways for Cassandra Data Operations

Don't reinvent the wheel.

Use Apache Spark

Use a Scheduler (Apache Airflow w/ Python)

- You can look, but Apache Spark is basically it. Look no further.
- Learn Spark, Python / Scala is fine. Just start using Apache Spark.
- Airflow, Jenkins, Luigi, Prefect, any scheduler can work, but Airflow has been proven for this.
- Airflow works with more than just Apache Cassandra, Apache Spark.. There are numerous Connections and Operators.





Demo

- <https://github.com/Anant/example-cassandra-etl-with-airflow-and-spark>
- Set up Astra
- Set up Airflow and Spark
- Connect Airflow and Spark
- Trigger DAG with PySpark jobs via Airflow UI
- Confirm data in Astra



Thank you and Dream Big.

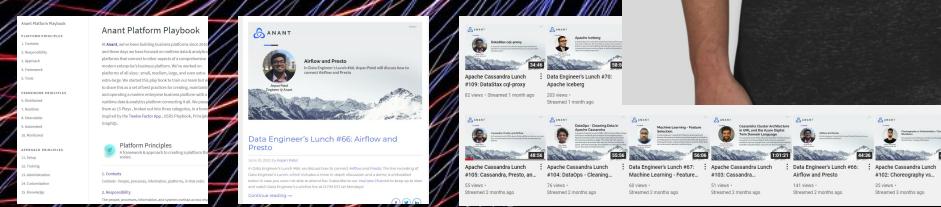


Hire us

- **Design** Workshops
- **Innovation** Sprints
- **Service Catalog**

Anant.us

- Read our Playbook
- Join our Mailing List
- Read up on Data Platforms
- Watch our Videos
- Download Examples



www.anant.us | solutions@anant.us | (855) 262-6826
3 Washington Circle, NW | #301 | Washington, DC 20037

