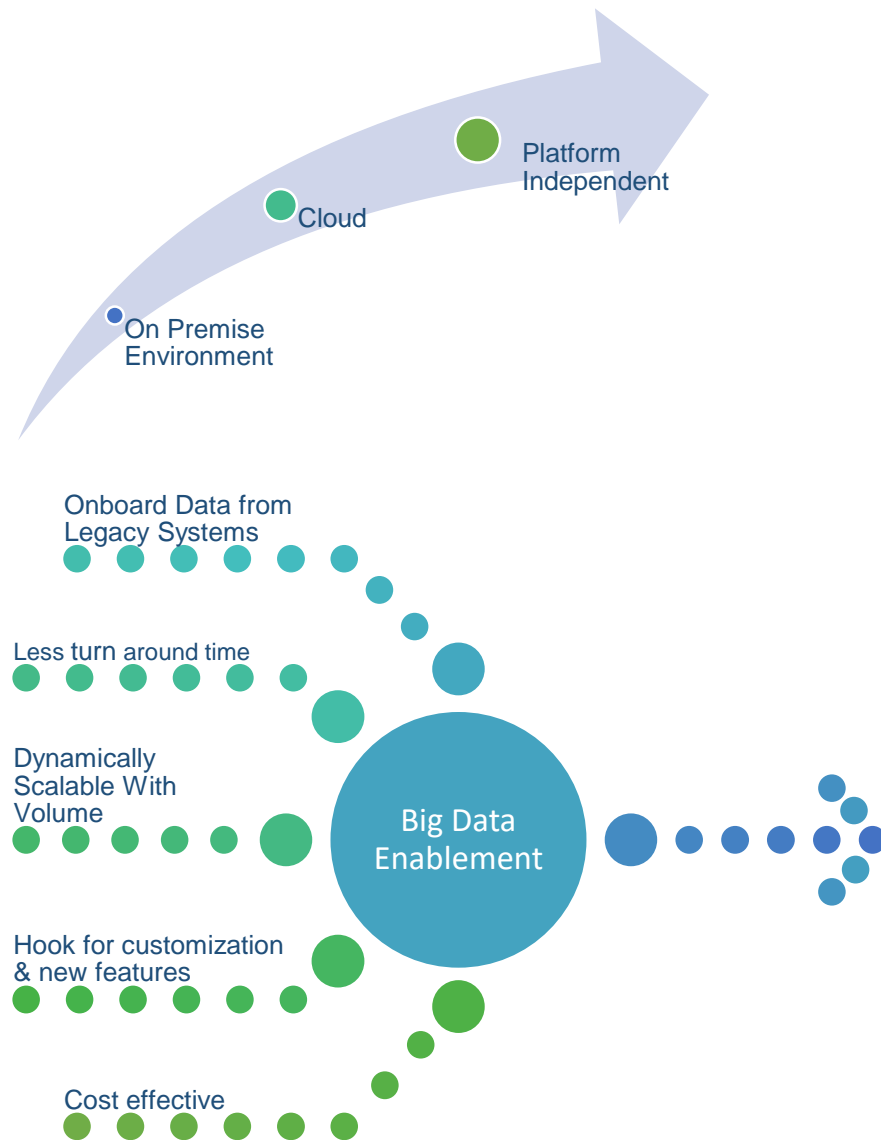


Cloud Enabled Ingestion Framework - Introduction



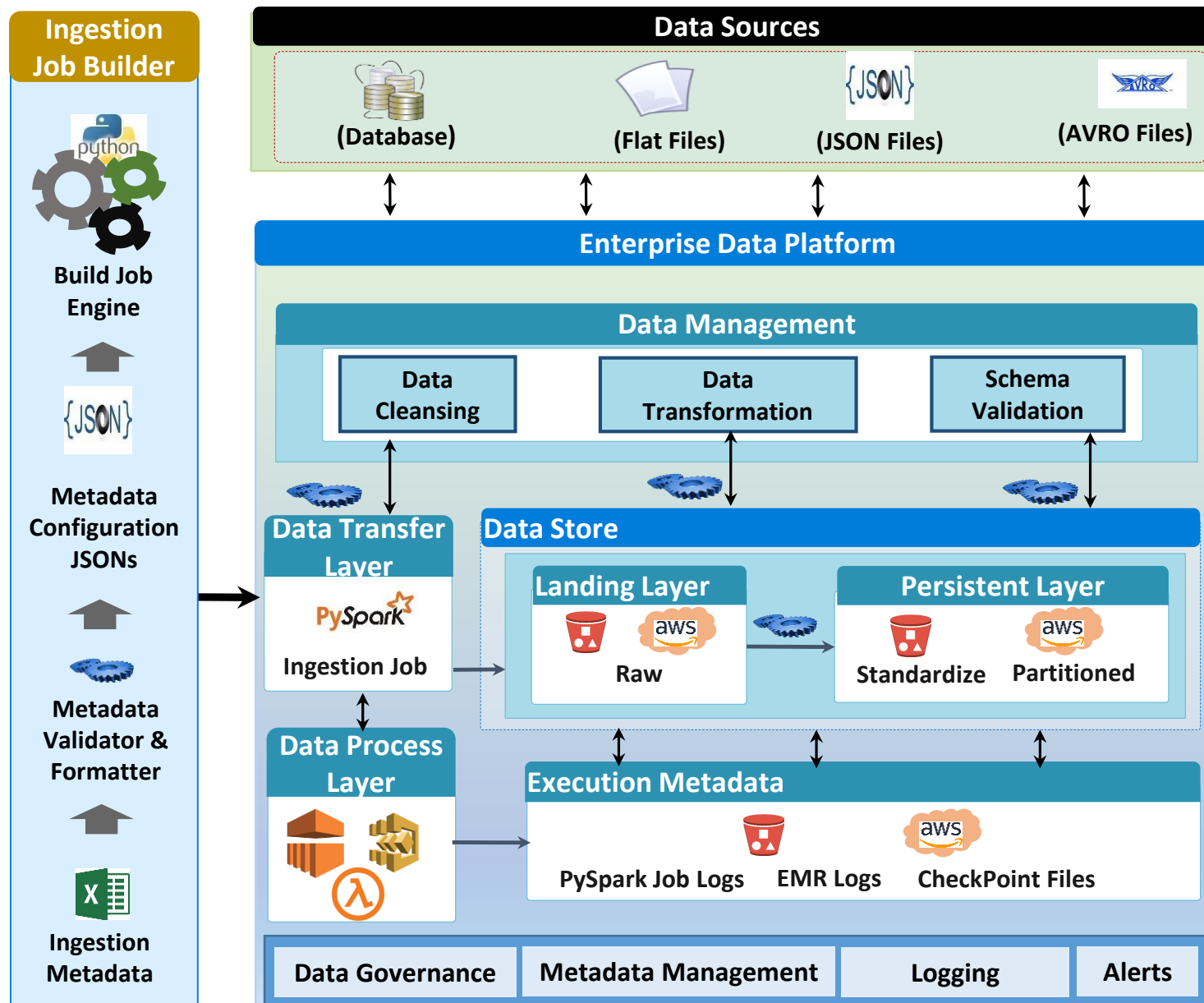
OPPORTUNITY

- As industry is maturing overall w.r.t. Big Data Technologies, most of the customers are planning to shift towards cloud platform. Keeping that in mind, there should be a solution which should be platform independent – On-Premise, Cloud (AWS, Azure, Google), Distribution agnostic keeping the following artifacts tightly coupled:
 - Portability across different platforms
 - Optimized deployment time
 - Optimized Cost w.r.t. development and customization
 - Scalability

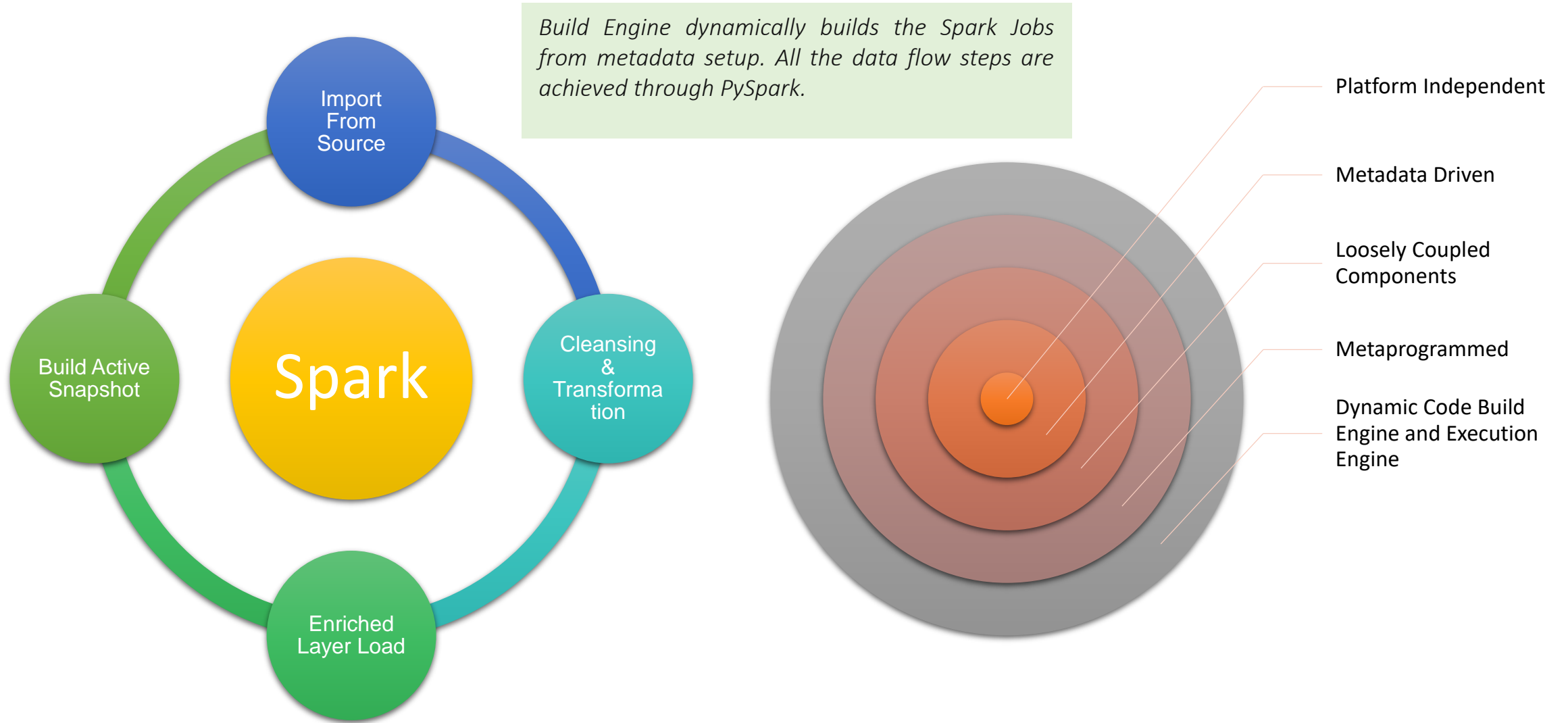
SOLUTION / APPROACH

1. Pure Spark based solution can be ported to any platform
2. Platform based API oriented solution to Read and Store Data
3. Hook for quickly onboard customer specific Cleansing and Transformation rules
4. Code Build Engine dynamically built the Spark jobs from Metadata keeping restart ability in mind from failure point
5. Metadata based Execution Engine dynamically orchestrates between several feeds within a source (Serial, Parallel, Predecessor dependencies)
6. Framework will be easily deployable across Client environments as Build Engine is Python based solution

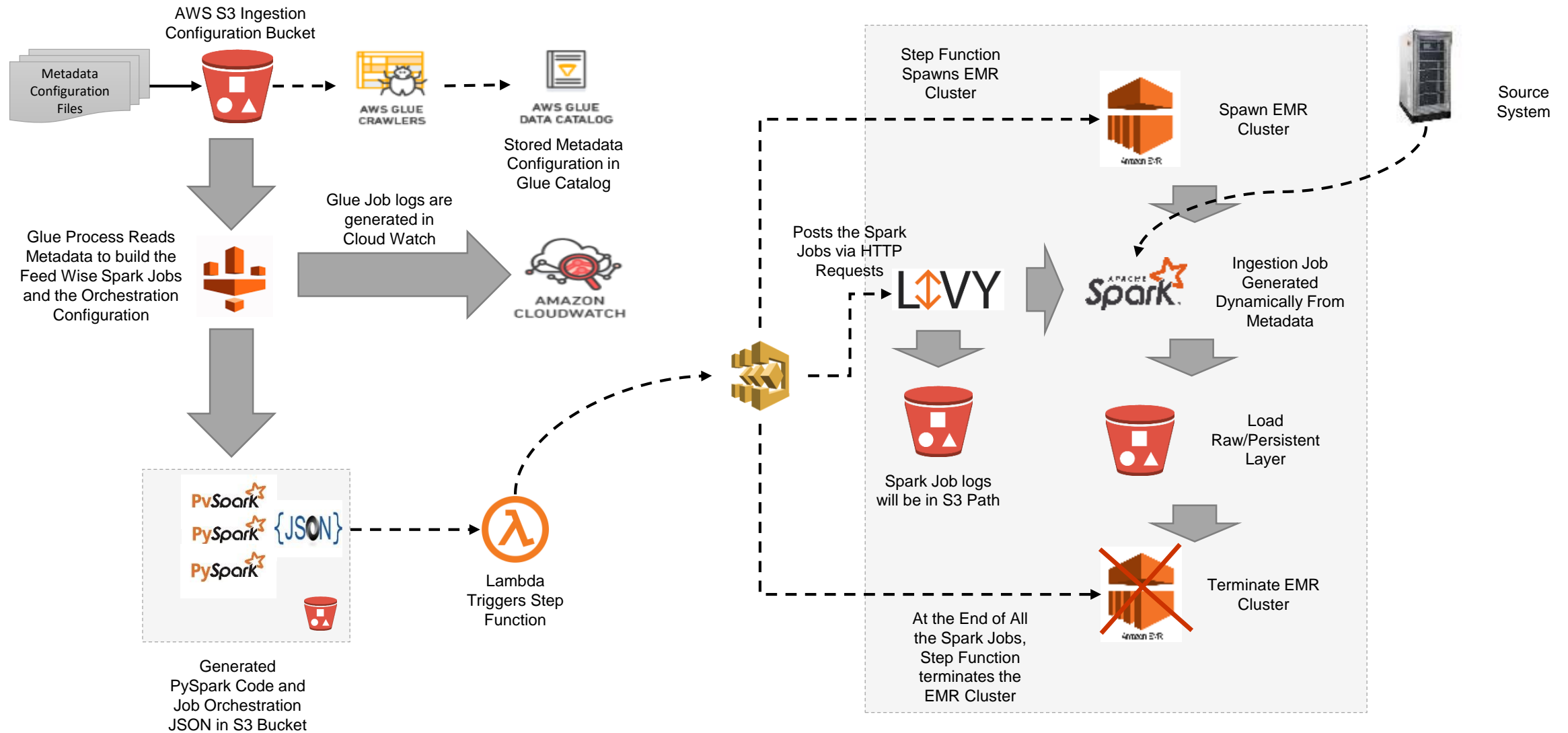
Framework Conceptual Architecture



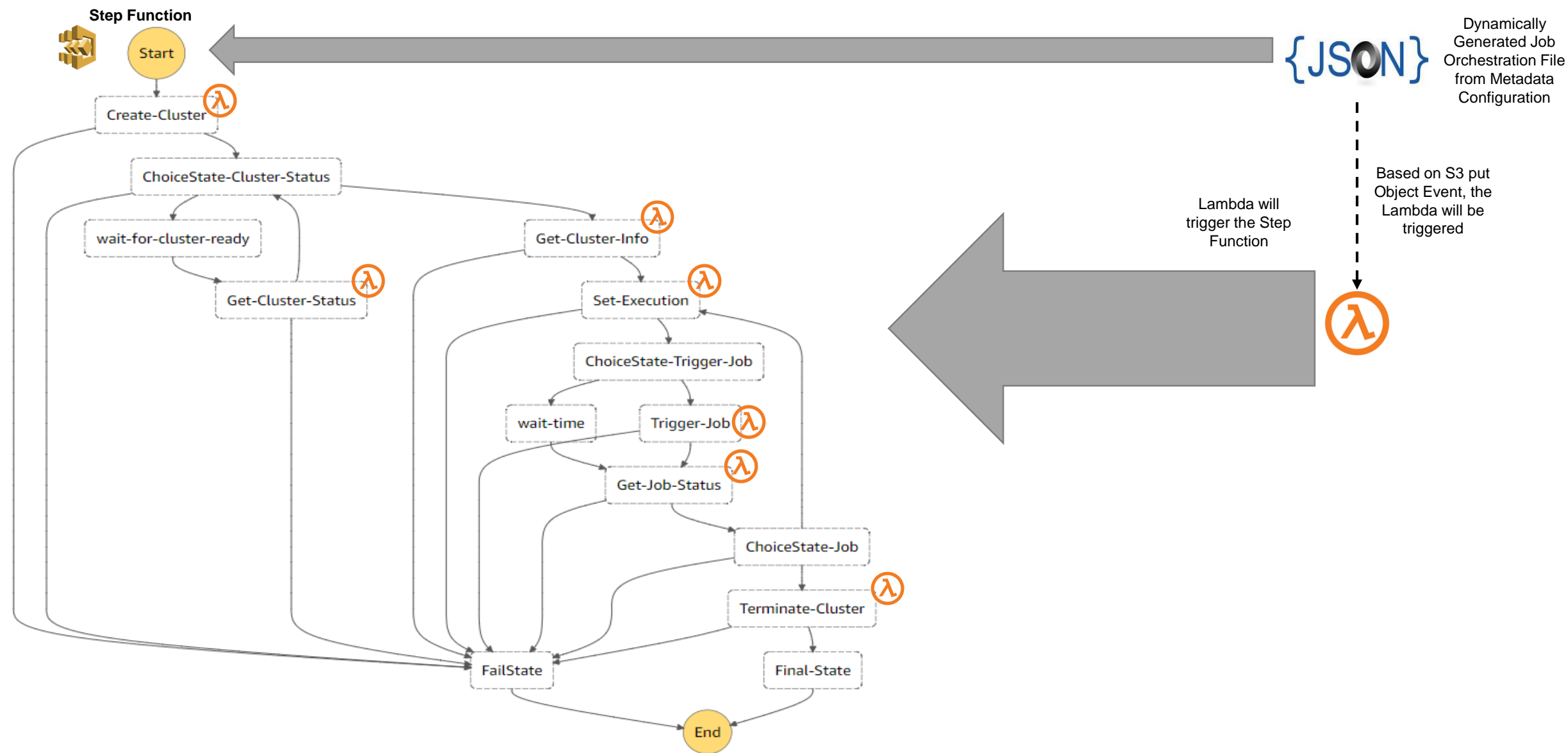
Design Principles



Framework Design - AWS



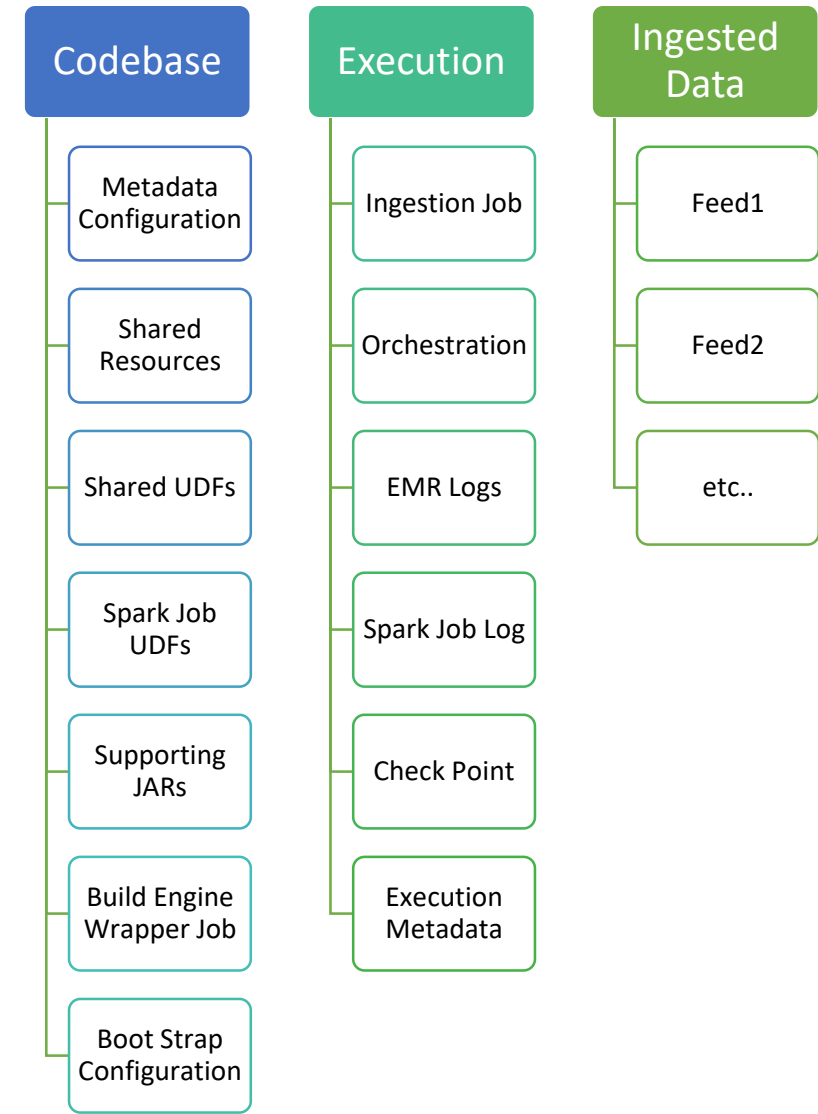
Orchestration Design - AWS



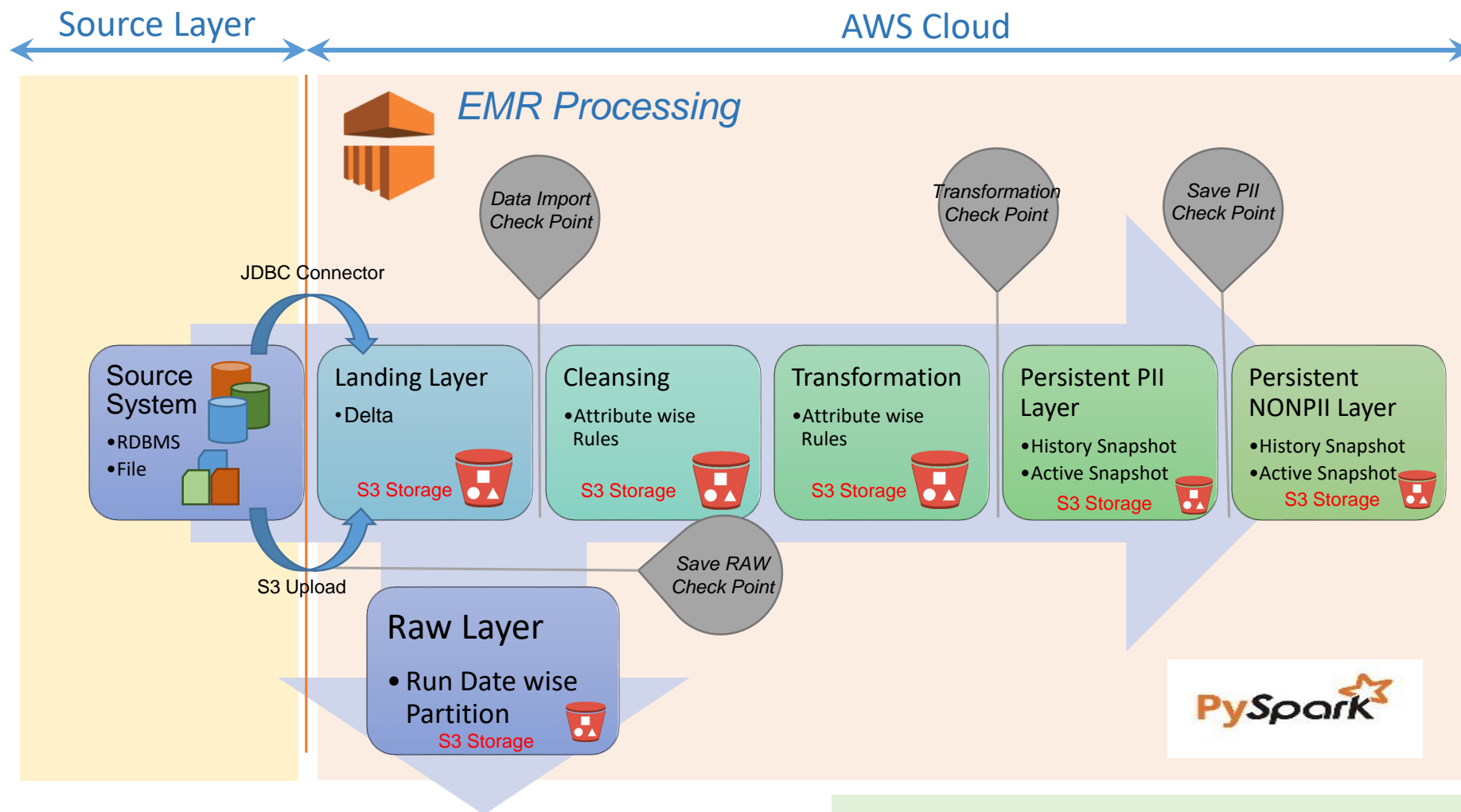
Directory Structure

Codebase Folder	Purpose
<i>Metadata Configuration</i>	<i>Metadata configuration for Source, Feed, Feed Attributes</i>
<i>Shared Resources</i>	<i>Common configuration setup files</i>
<i>Shared UDFs</i>	<i>UDFs to Dynamically generate the PySpark code per Metadata Configuration</i>
<i>Spark JOB UDFs</i>	<i>UDFs for Cleansing and Transformation rules</i>
<i>Supporting JARs</i>	<i>JAR files required for dynamically generated PySpark Job</i>
<i>Build Engine Wrapper</i>	<i>Wrapper script which uses the Shared UDFs to dynamically generate the PySpark Ingestion Job</i>
<i>Boot Strap Configuration</i>	<i>Boot Strap configuration script to spawn EMR Cluster</i>

Execution Folder	Purpose
<i>Ingestion Job</i>	<i>Generated PySpark Job script</i>
<i>Orchestration</i>	<i>Generated Feed wise execution orchestration configuration if ingestion process is executed at source level</i>
<i>EMR Logs</i>	<i>EMR Cluster logs for each time a cluster is spinning up</i>
<i>Spark Job Logs</i>	<i>Feed wise PySpark job logs</i>

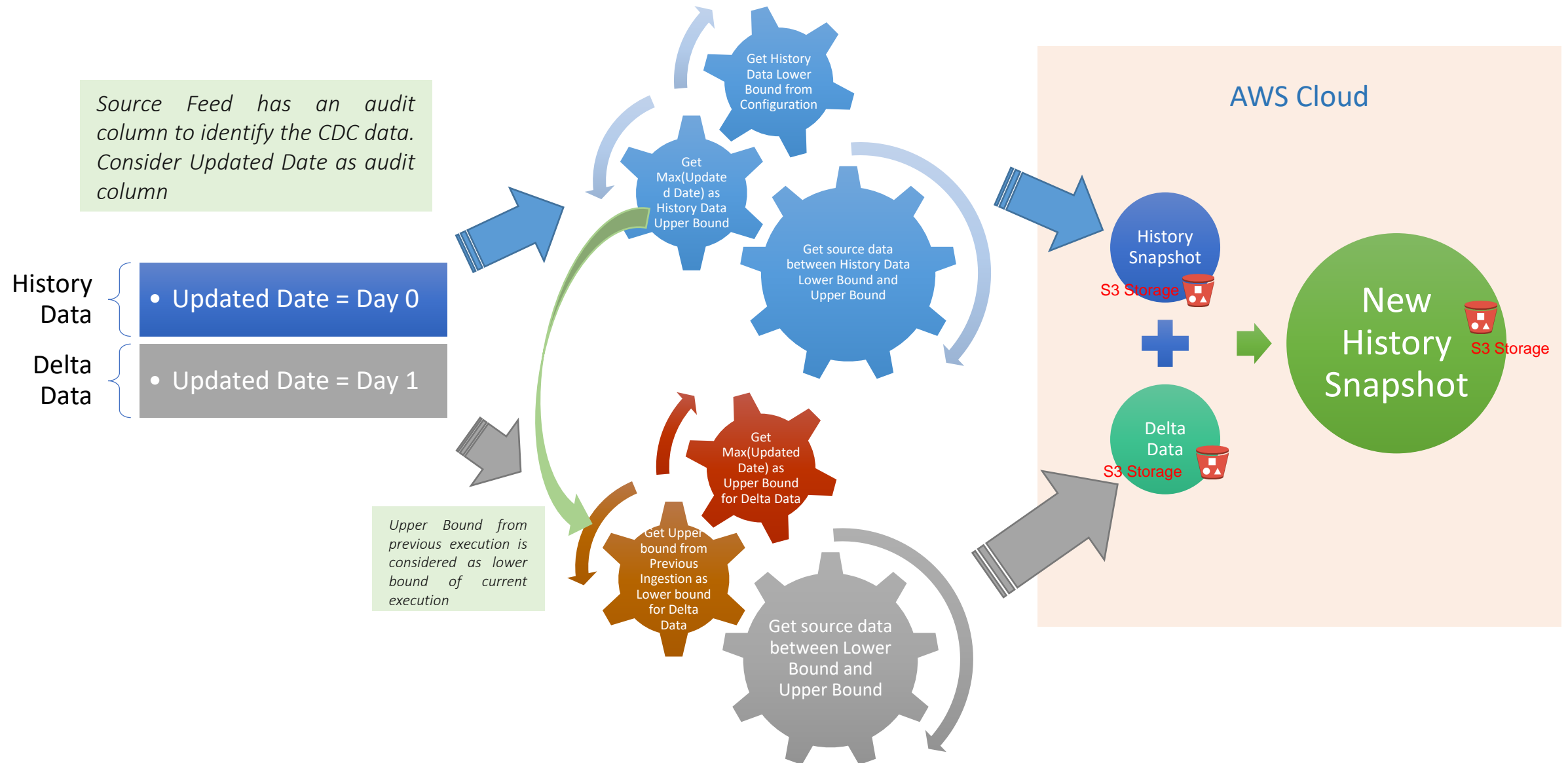


Data Flow – Feed Level

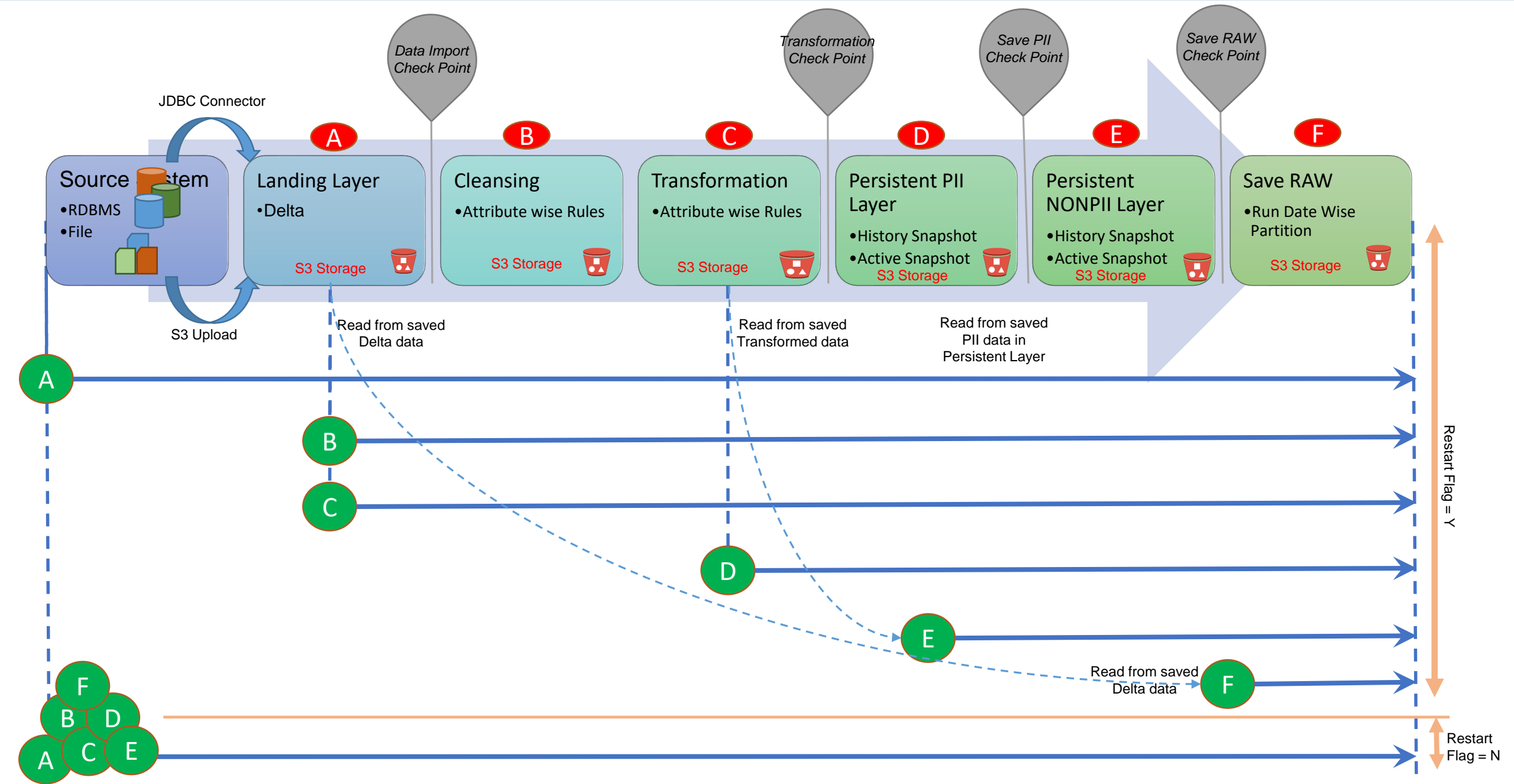


As data flows through various layers of Ingestion, checkpoints are defined and data is saved intermittently for Restart ability and recovery from failure point.

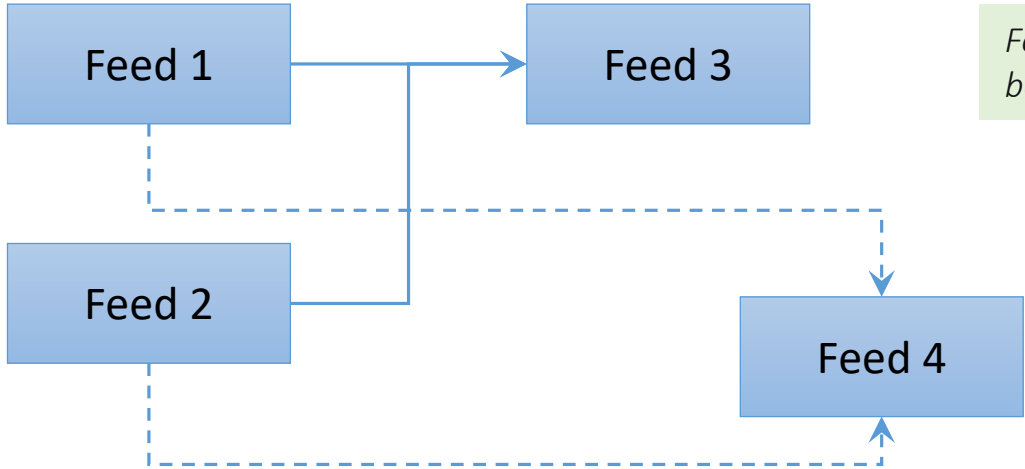
CDC Mechanism – Feed Level



Restart ability – Feed Level



Restart ability – Source Level



Feed3 has predecessor on Feed1 and Feed2. Once both of them are successful, then Feed3 will start

Feed4 has predecessor on either Feed1 and Feed2. Once either of them is successful, then Feed4 will start

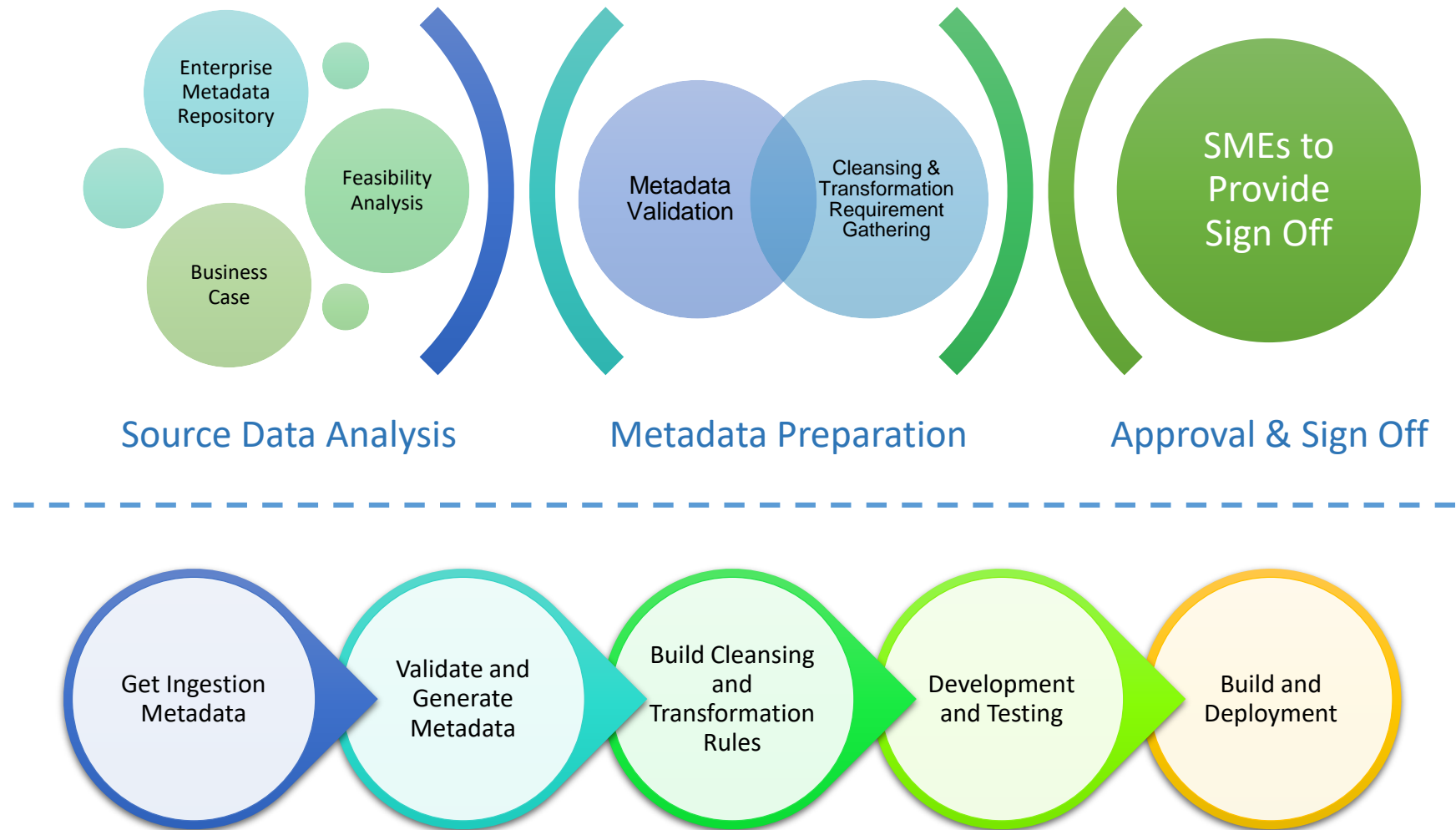
Failure Scenario	Status	On Restart (Restart Flag = Y)
Feed1	Success	Process will only build PySpark job for Feed2 from failure point and Feed3 from beginning
Feed2	Failed	
Feed3	Skipped	
Feed4	Success	

Failure Scenario	Status	On Restart (Restart Flag = Y)
Feed1	Success	Process will only build PySpark job for Feed3 from failure point and execute Feed3 ingestion
Feed2	Success	
Feed3	Failed	
Feed4	Success	

Failure Scenario	Status	On Restart (Restart Flag = Y)
Feed1	Success	Process will only build PySpark job for Feed4 from failure point and execute Feed4 ingestion
Feed2	Success	
Feed3	Success	
Feed4	Failed	

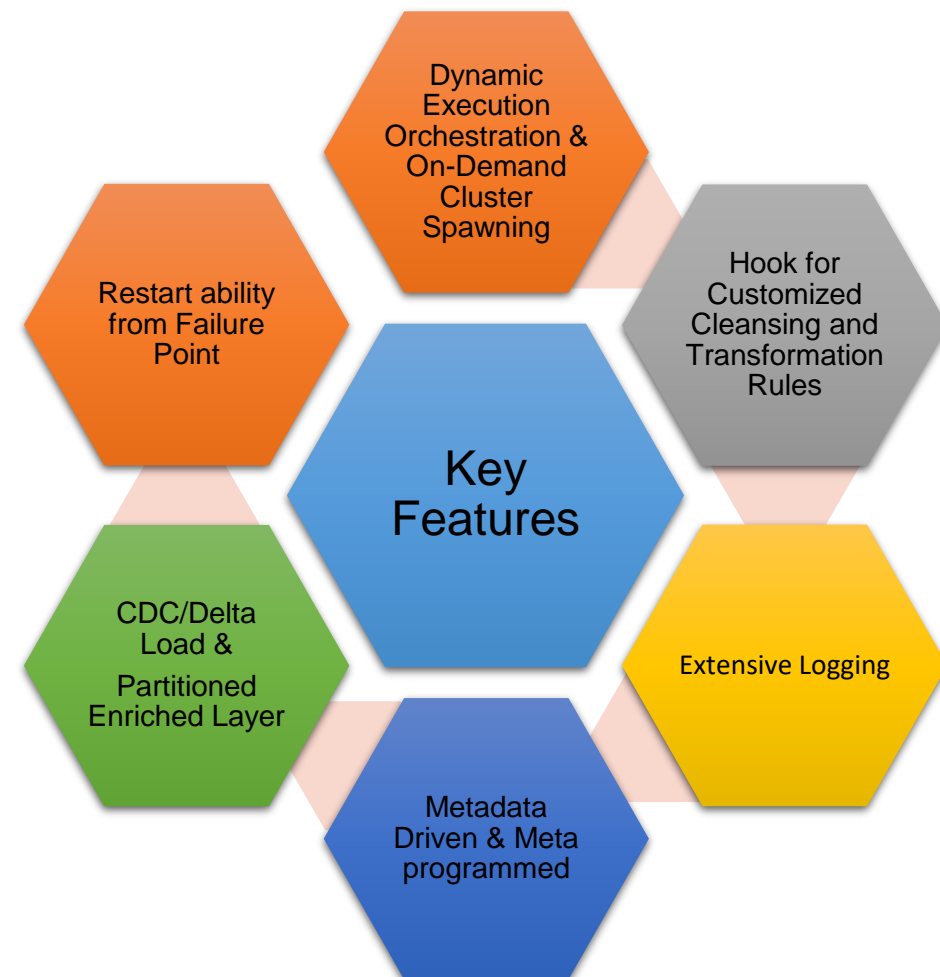
Failure Scenario	Status	On Restart (Restart Flag = Y)
Feed1	Success	Process wont find any failure, so it will exit.
Feed2	Success	
Feed3	Success	
Feed4	Success	

Development Life Cycle



Current Features

S#	Available Features
1.	<i>Supports RDBMS – Oracle, SQL Server, TEXT file – Fixed width & delimited</i>
2.	<i>Supports data ingestion to Raw & Enriched layer</i>
3.	<i>Supports PII and NON PII Ingestion</i>
4.	<i>Feed wise Logging mechanism</i>
5.	<i>Support data archiving for files on HDFS after processing into Raw & Enriched</i>
6.	<i>Support for execution at source or feed level</i>
7.	<i>Support for hive table partitioning in the Raw & Enriched layer</i>
8.	<i>Dynamically build's the PySpark Code and Spawn EMR Cluster, Executes all the ingestion Source/Feed Level. Once all the ingestion is completed in S3, terminates the cluster.</i>
9.	<i>Supports Cleansing and Transformations as configured in Metadata</i>
10.	<i>Builds the Job Orchestration configuration dynamically from Metadata configuration and executes the ingestion for different feeds in that fashion Within a Source, inter-dependency between feeds (Parallel/Serial) can be configured</i>
11.	<i>Supports RAW and Enriched Layer storage in TEXT and Parquet format</i>
12.	<i>Supports CDC/Delta Load, Full Load</i>
13.	<i>Restart ability within a Source</i>
14.	<i>Restart ability within a Feed Ingestion from Failure Point</i>
15.	<i>RDBMS Ingestion - Build History and Current Snapshot in Enriched Layer</i>
16.	<i>Metadata Generator – Generate JSON metadata configuration files from pre-defined excel spread sheet</i>



Upcoming Features - Planned

S#	Upcoming Features
1.	<i>Ingestion of JSON, AVRO File Ingestion</i>
2.	<i>Auditing mechanism</i>
3.	<i>Support ORC format storage for RAW and Enriched Layer</i>
4.	<i>Azure enablement</i>
5.	<i>UI Development</i>

