



Ingestion and historization in the Data Lake

Illia Todor

/iamtodor



Data Engineer

AWS Cloud Practitioner Certified

/iamtodor



Data Engineer

AWS Cloud Practitioner Certified

HRS GROUP



HRS GROUP

Data Engineer

AWS Cloud Practitioner Certified

<https://github.com/iamtodor>

<https://stackoverflow.com/users/5151861/iamtodor>

<https://iamtodor.medium.com/>

<https://www.linkedin.com/in/iamtodor/>

HRS GROUP

provides the world's premier lodging
experience to corporations globally.



50
markets worldwide
32,000
cities globally



1,000
employees

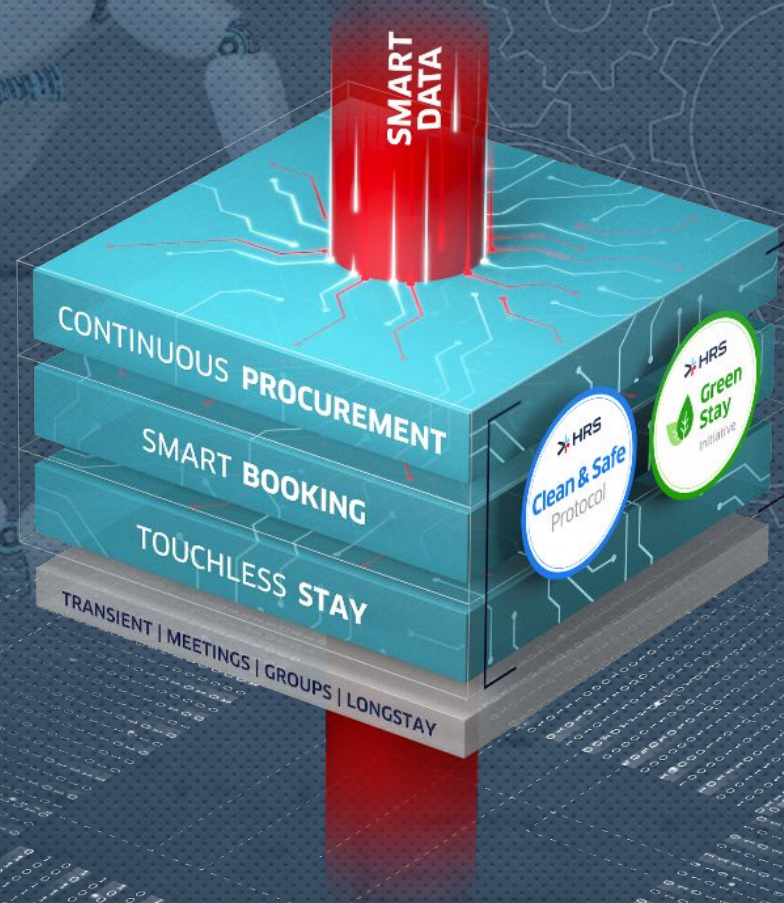
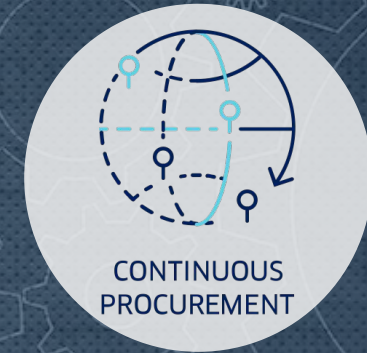


> 800.000
yearly contacts
with Hotels

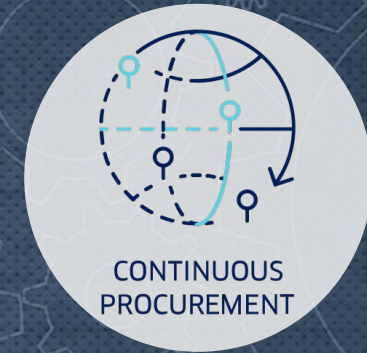


> 6,000
Global customers

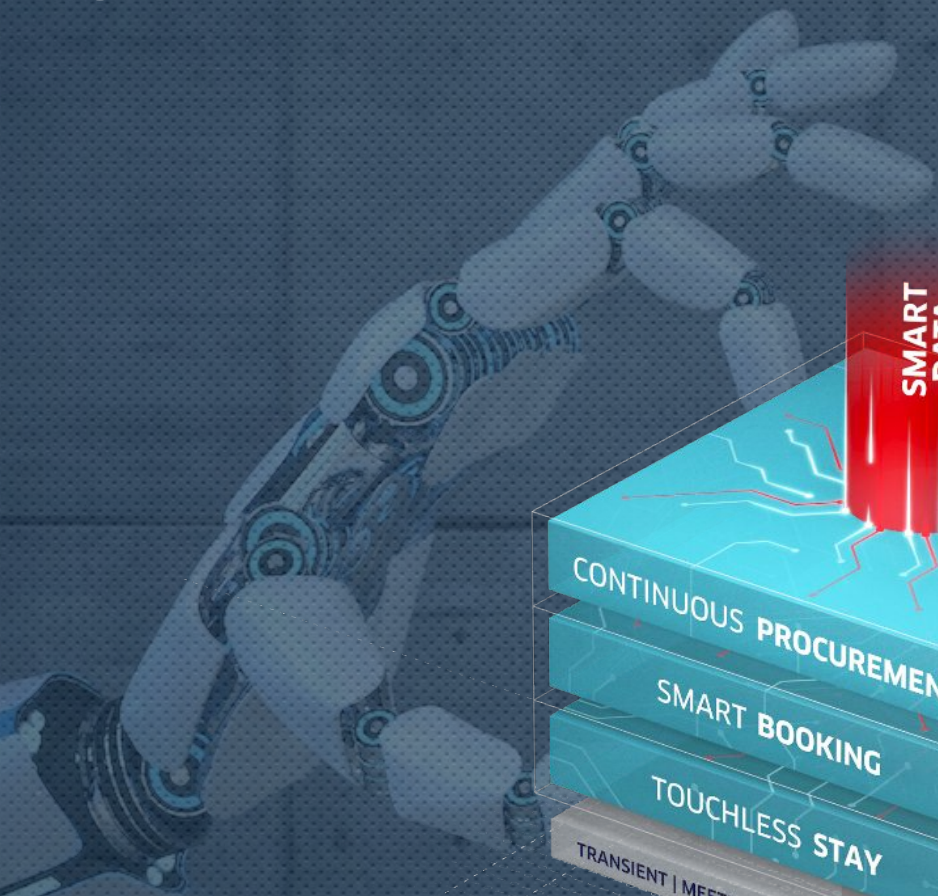
Lodging as a Service



Lodging as a Service



Lodging as a Service



Request

- DataSource
- Historical Changes
- DataLake
- DataWarehouse
- Dashboards



Historization: the process

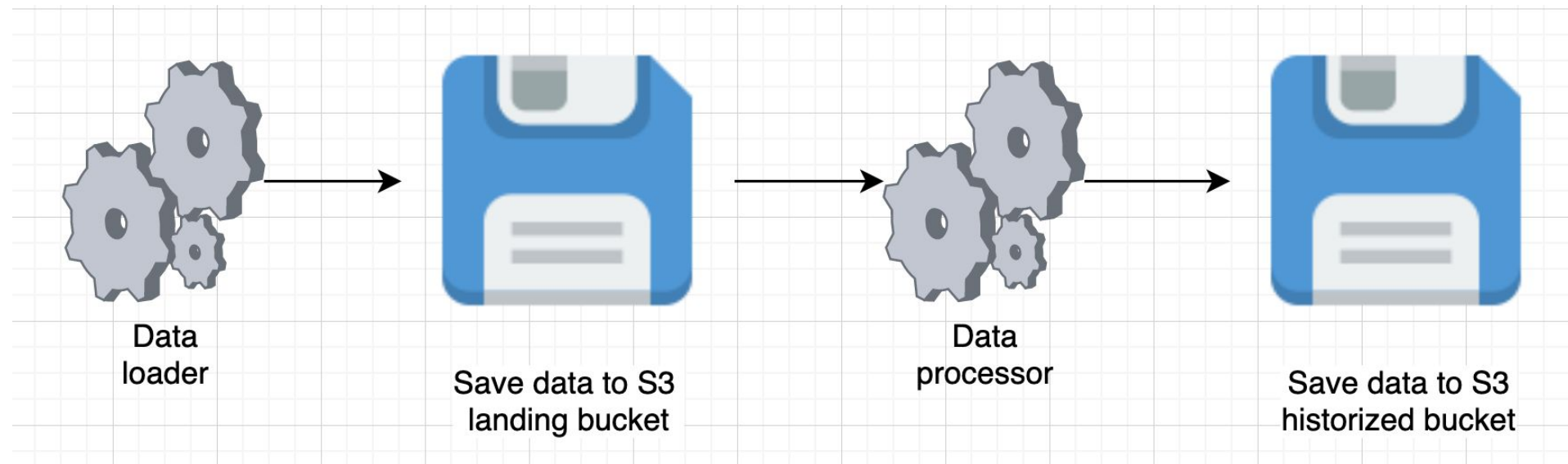


We wanted to move ingestions from DataSource into Data Lake along with applying historization approach. The data we need in the DWH for reporting.

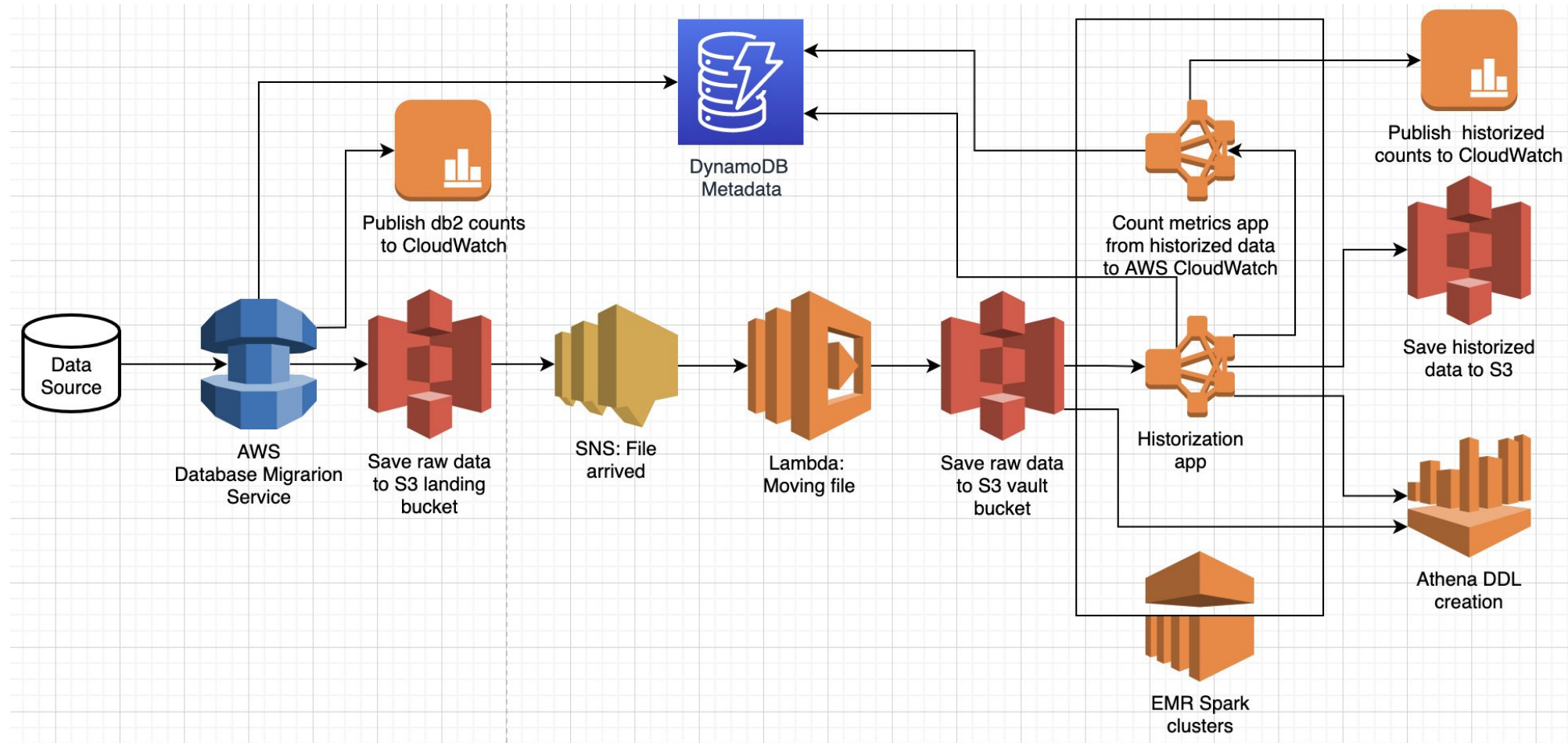
The idea of historization approach is to track historical changes.

Operation in the source DB	Field value	valid_from	valid_to
Insert	Andrew	2021-01-01	2021-01-03
Update	Mark	2021-01-03	2021-01-05
Update	Daniel	2021-01-05	2099-12-31

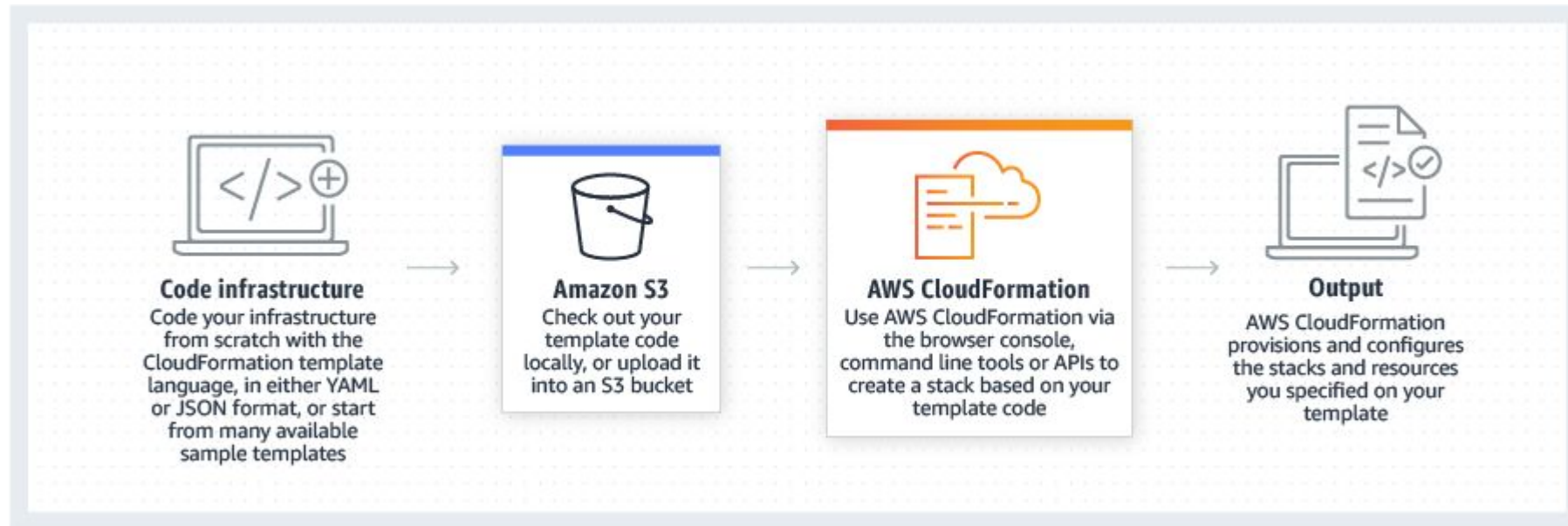
Architecture: abstract



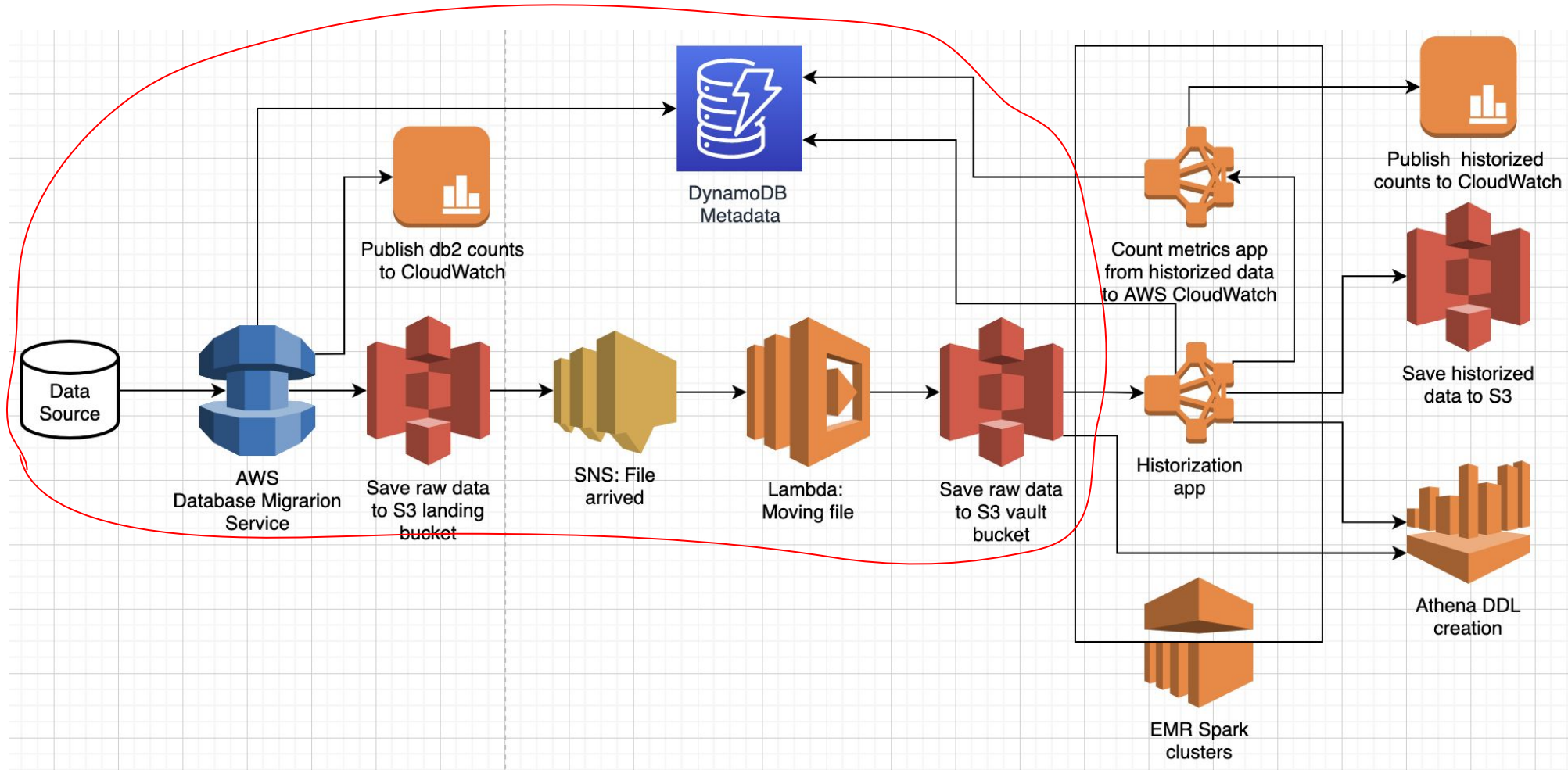
Architecture: in details



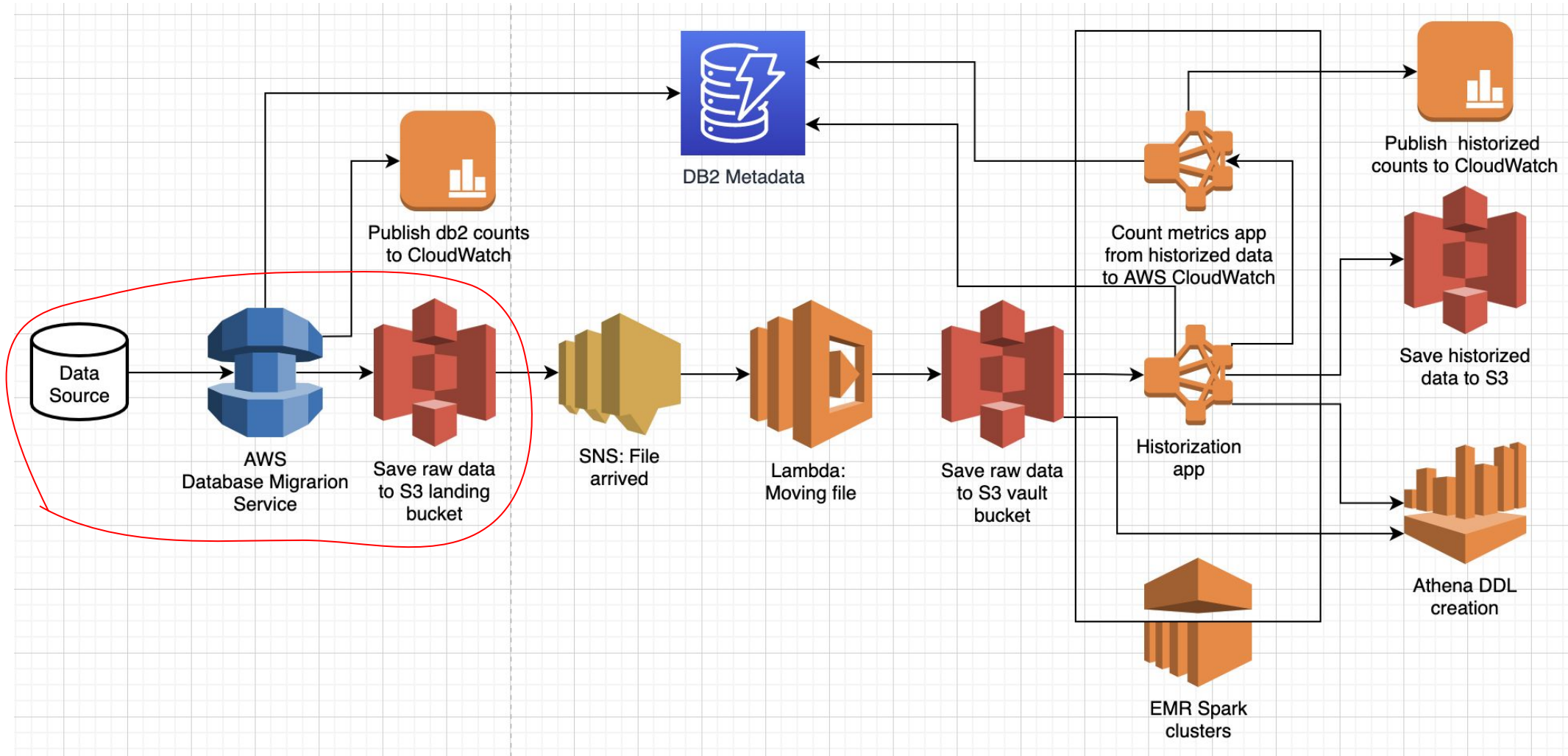
Infrastructure as a Code



First part: global overview



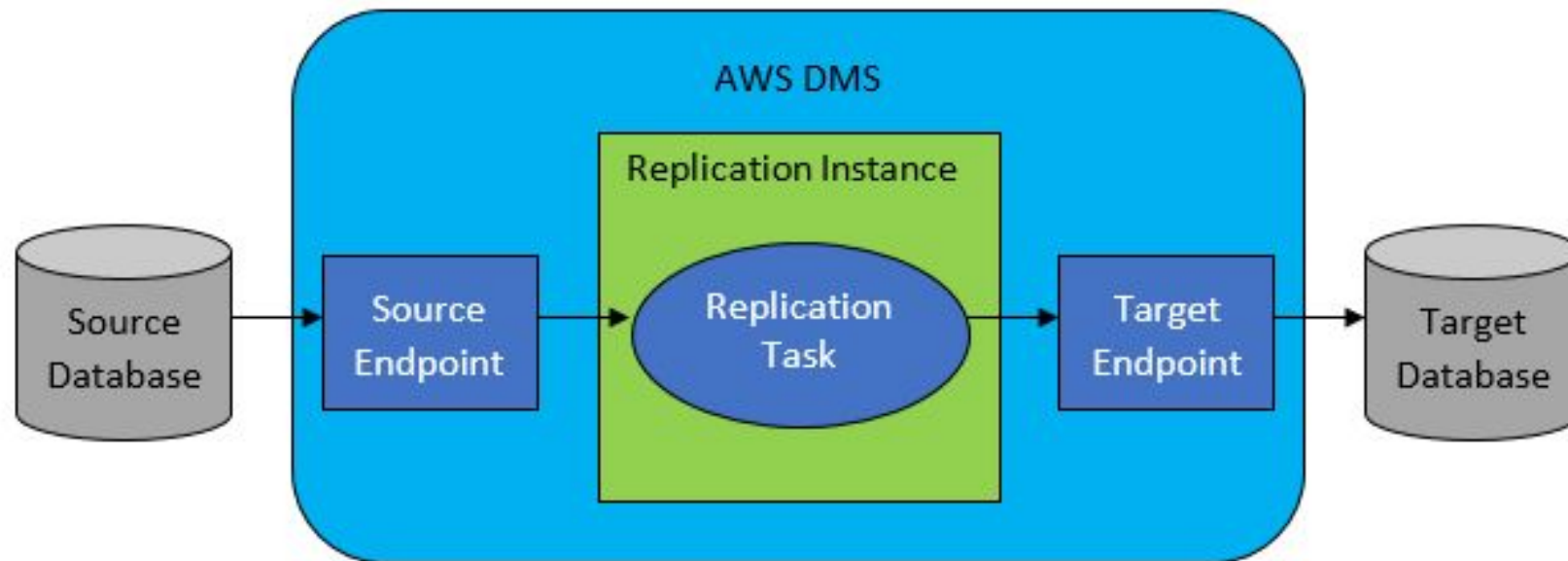
First part: DMS



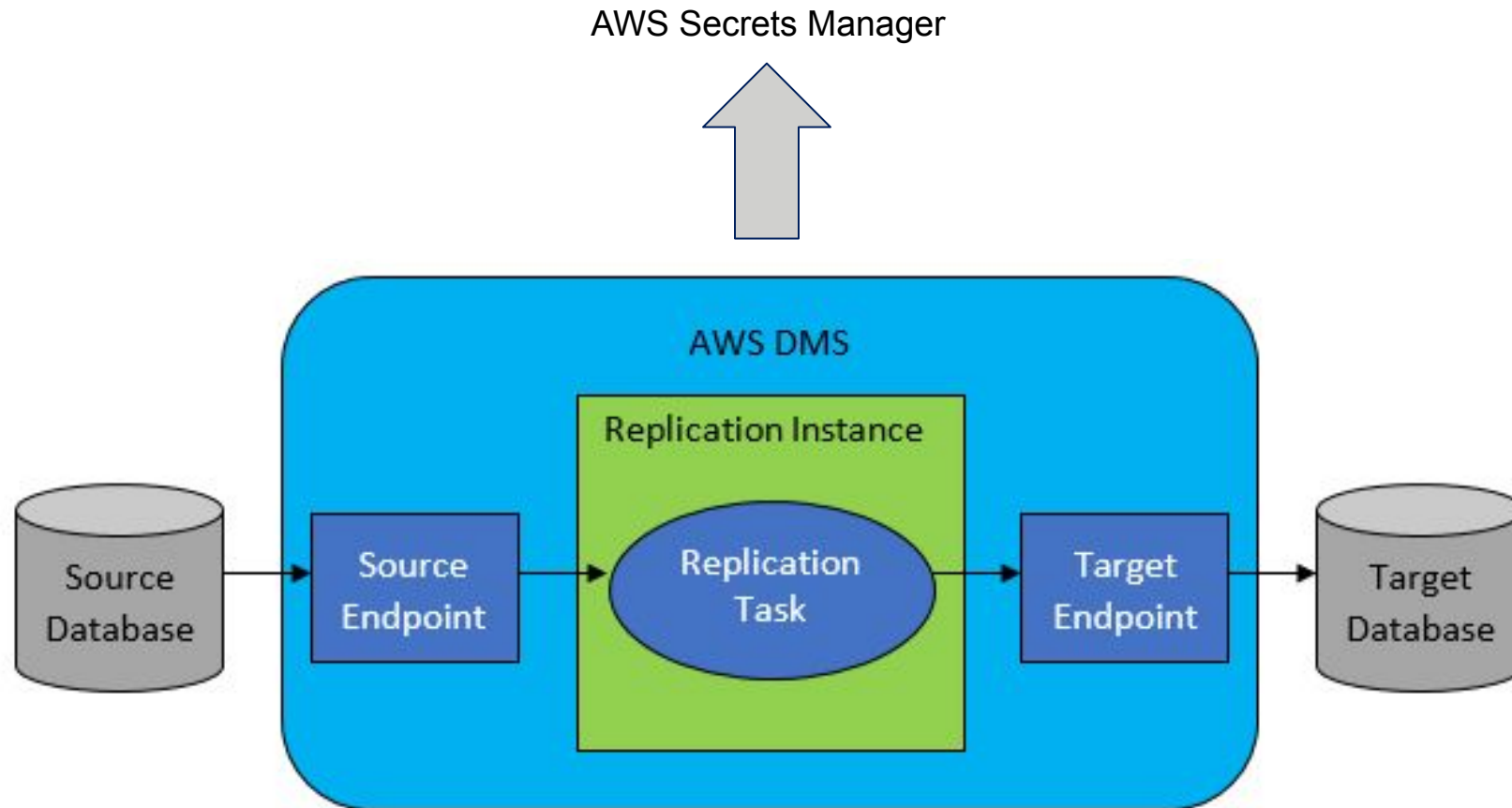
AWS DMS: Database Migration Service



AWS DMS: Database Migration Service



AWS DMS: Database Migration Service



DMS: Replication Instance



Replication instance type	vCPU	Memory (GiB)
Memory Optimized		
dms.r4.large	2	15.25
dms.r4.xlarge	4	30.5
dms.r4.2xlarge	8	61
dms.r4.4xlarge	16	122
dms.r4.8xlarge	32	244
dms.r5.large	2	16
dms.r5.xlarge	4	32
dms.r5.2xlarge	8	64
dms.r5.4xlarge	16	128

DMS: Target Endpoint



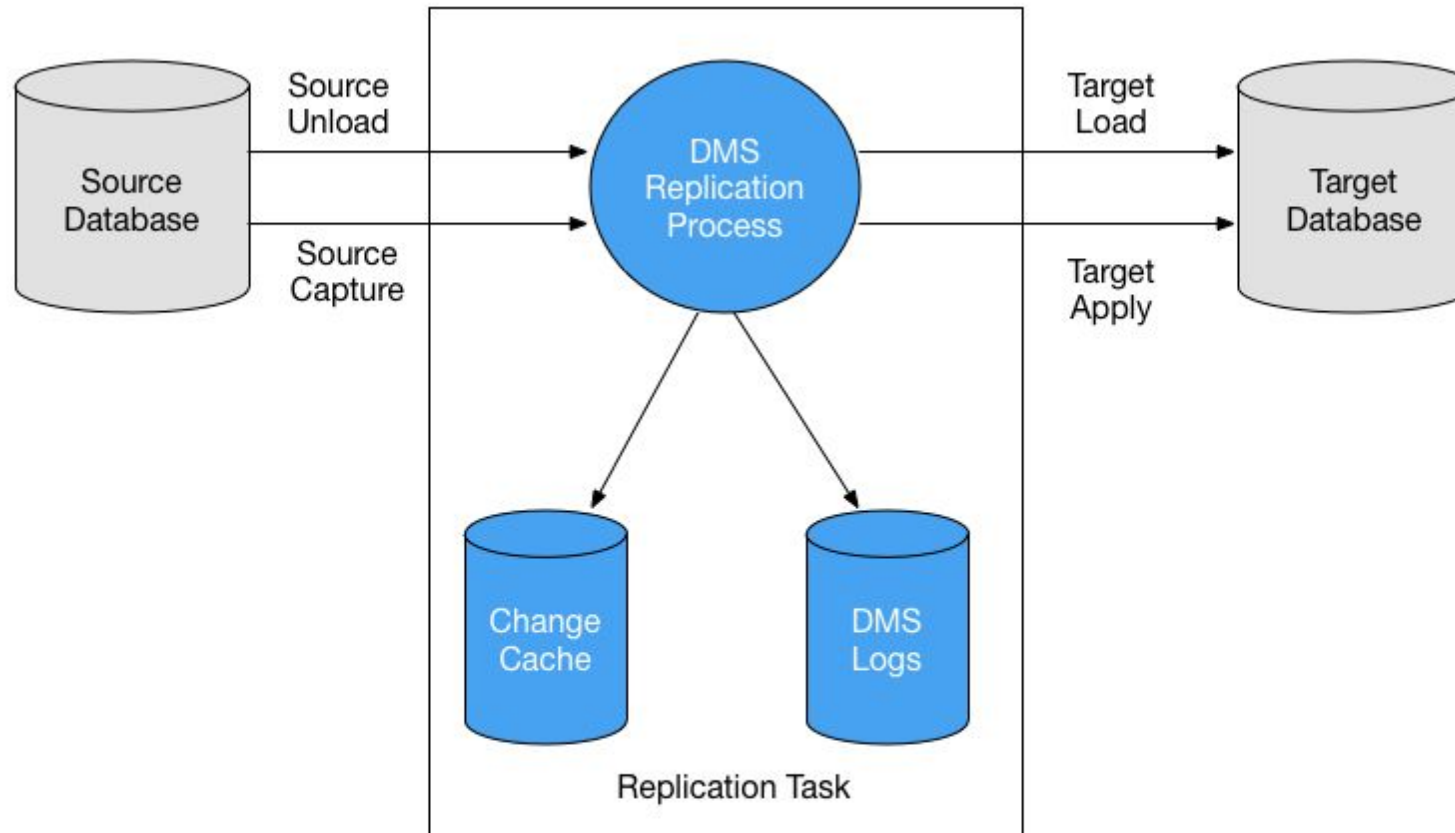
```
1 S3ProdDMS Endpoint:
2   Type: AWS::DMS::Endpoint
3   Properties:
4     EndpointIdentifier: prod-s3-landing-parquet-endpoint
5     EngineName: s3
6     EndpointType: target
7     S3Settings:
8       BucketName: 'bucket_name'
9       BucketFolder: 'prod/prefix'
10      CompressionType: gzip
11      ExtraConnectionAttributes:
12        'dataFormat=parquet;datePartitionDelimiter=DASH;datePartitionEnabled=false;datePartitionSequence=YYY
13        YMMDD;'
```


DMS: Target Endpoint



```
1 S3ProdDMSEndpoint:
2   Type: AWS::DMS::Endpoint
3   Properties:
4     EndpointIdentifier: prod-s3-landing-parquet-endpoint
5     EngineName: s3
6     EndpointType: target
7     S3Settings:
8       BucketName: 'bucket_name'
9       BucketFolder: 'prod/prefix'
10      CompressionType: gzip
11      ExtraConnectionAttributes:
12        'dataFormat=parquet;datePartitionDelimiter=DASH;datePartitionEnabled=false;datePartitionSequence=YYY
13        YMMDD;'
```


DMS: Replication Task



DMS: Replication Task



```
1 ReplicationTask:
2   Type: AWS::DMS::ReplicationTask
3   Properties:
4     ReplicationTaskIdentifier: !Sub "${Env}-{source}-{target}-{schema_name}-{table_name}"
5     MigrationType: full-load
6     ReplicationInstanceArn: !ImportValue replication-instance-arn
7     SourceEndpointArn: !ImportValue source-endpoint-arn
8     TargetEndpointArn: !ImportValue s3-prod-endpoint
9     ReplicationTaskSettings: '{"Logging": {"EnableLogging": true}}'
10    TableMappings: |
11      {
12        "rules": [
13          {
14            "rule-type": "selection",
15            "rule-id": 1,
16            "rule-name": 1,
17            "rule-action": "include",
18            "object-locator": {
19              "schema-name": "schema_name",
20              "table-name": "table_name"
21            }
22          }
23        ]
24      }
```


DMS: Replication Task



```
1 ReplicationTask:
2   Type: AWS::DMS::ReplicationTask
3   Properties:
4     ReplicationTaskIdentifier: !Sub "${Env}-{source}-{target}-{schema_name}-{table_name}"
5     MigrationType: full-load
6     ReplicationInstanceArn: !ImportValue replication-instance-arn
7     SourceEndpointArn: !ImportValue source-endpoint-arn
8     TargetEndpointArn: !ImportValue s3-prod-endpoint
9     ReplicationTaskSettings: '{"Logging": {"EnableLogging": true}}'
10    TableMappings: |
11      {
12        "rules": [
13          {
14            "rule-type": "selection",
15            "rule-id": 1,
16            "rule-name": 1,
17            "rule-action": "include",
18            "object-locator": {
19              "schema-name": "schema_name",
20              "table-name": "table_name"
21            }
22          }
23        ]
24      }
```


DMS: Replication Task



```
1 ReplicationTask:
2   Type: AWS::DMS::ReplicationTask
3   Properties:
4     ReplicationTaskIdentifier: !Sub "${Env}-{source}-{target}-{schema_name}-{table_name}"
5     MigrationType: full-load
6     ReplicationInstanceArn: !ImportValue replication-instance-arn
7     SourceEndpointArn: !ImportValue source-endpoint-arn
8     TargetEndpointArn: !ImportValue s3-prod-endpoint
9     ReplicationTaskSettings: '{"Logging": {"EnableLogging": true}}'
10    TableMappings: |
11      {
12        "rules": [
13          {
14            "rule-type": "selection",
15            "rule-id": 1,
16            "rule-name": 1,
17            "rule-action": "include",
18            "object-locator": {
19              "schema-name": "schema_name",
20              "table-name": "table_name"
21            }
22          }
23        ]
24      }
```


DMS: Replication Task



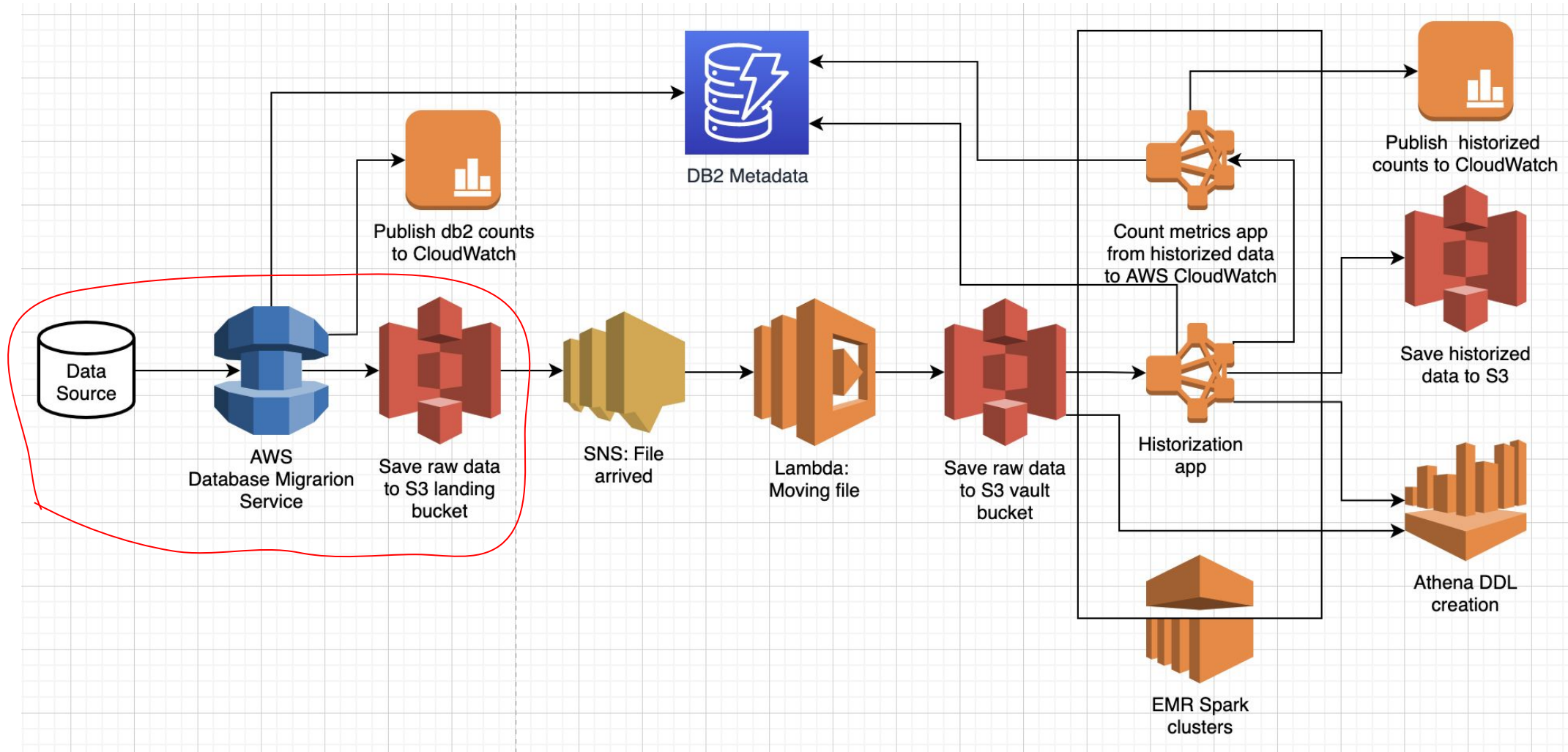
```
1 ReplicationTask:
2   Type: AWS::DMS::ReplicationTask
3   Properties:
4     ReplicationTaskIdentifier: !Sub "${Env}-{source}-{target}-{schema_name}-{table_name}"
5     MigrationType: full-load
6     ReplicationInstanceArn: !ImportValue replication-instance-arn
7     SourceEndpointArn: !ImportValue source-endpoint-arn
8     TargetEndpointArn: !ImportValue s3-prod-endpoint
9     ReplicationTaskSettings: '{"Logging": {"EnableLogging": true}}'
10    TableMappings: |
11      {
12        "rules": [
13          {
14            "rule-type": "selection",
15            "rule-id": 1,
16            "rule-name": 1,
17            "rule-action": "include",
18            "object-locator": {
19              "schema-name": "schema_name",
20              "table-name": "table_name"
21            }
22          }
23        ]
24      }
```

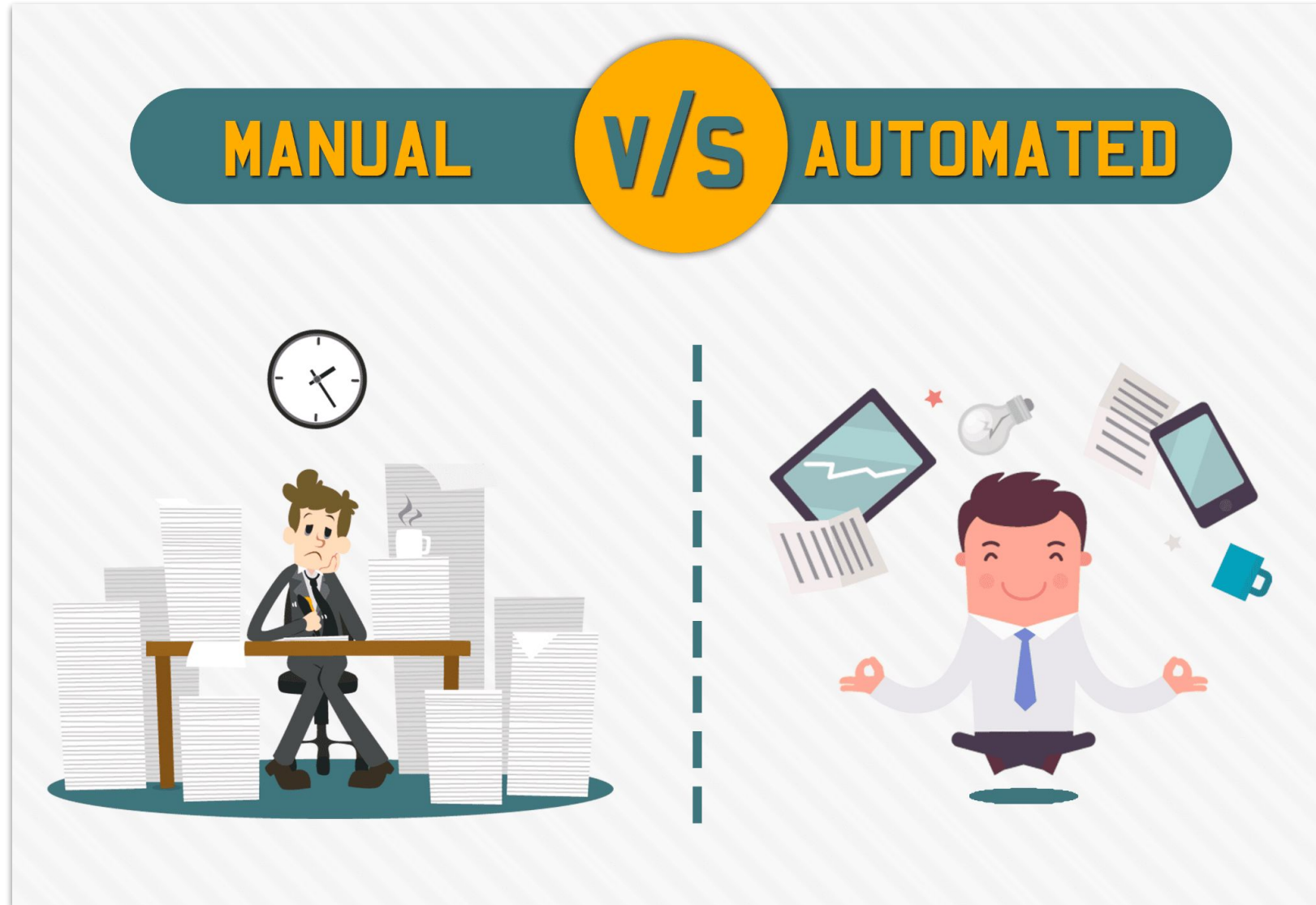

DMS: Replication Task



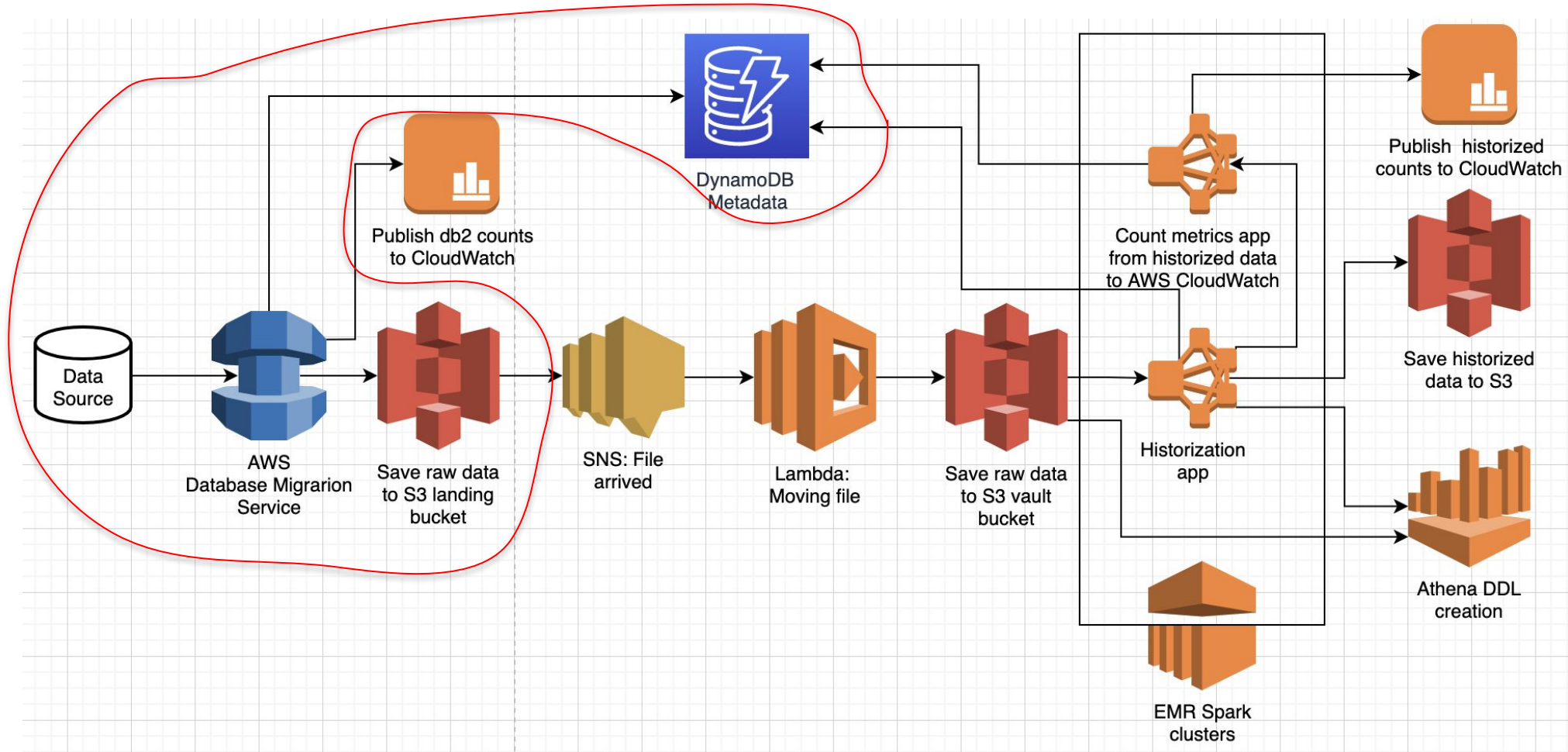
```
1 ReplicationTask:
2   Type: AWS::DMS::ReplicationTask
3   Properties:
4     ReplicationTaskIdentifier: !Sub "${Env}-{source}-{target}-{schema_name}-{table_name}"
5     MigrationType: full-load
6     ReplicationInstanceArn: !ImportValue replication-instance-arn
7     SourceEndpointArn: !ImportValue source-endpoint-arn
8     TargetEndpointArn: !ImportValue s3-prod-endpoint
9     ReplicationTaskSettings: '{"Logging": {"EnableLogging": true}}'
10    TableMappings: |
11      {
12        "rules": [
13          {
14            "rule-type": "selection",
15            "rule-id": 1,
16            "rule-name": 1,
17            "rule-action": "include",
18            "object-locator": {
19              "schema-name": "schema_name",
20              "table-name": "table_name"
21            }
22          }
23        ]
24      }
```


First part: DMS

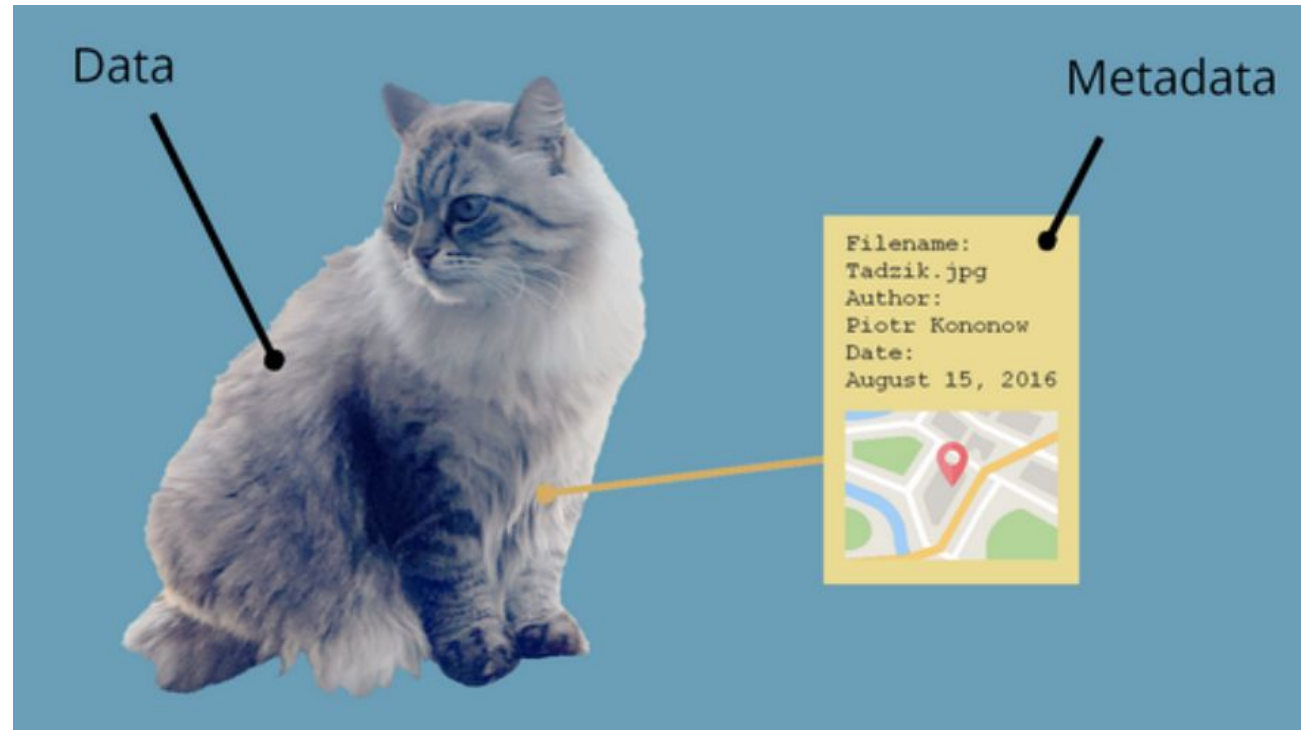




First part: Metadata management



What is metadata itself



Suitable tool for metadata



CSV

Suitable tool for metadata



CSV



DynamoDB

Suitable tool for metadata



CSV



DynamoDB



DynamoDB




```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```



```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```



```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```



```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```



```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```

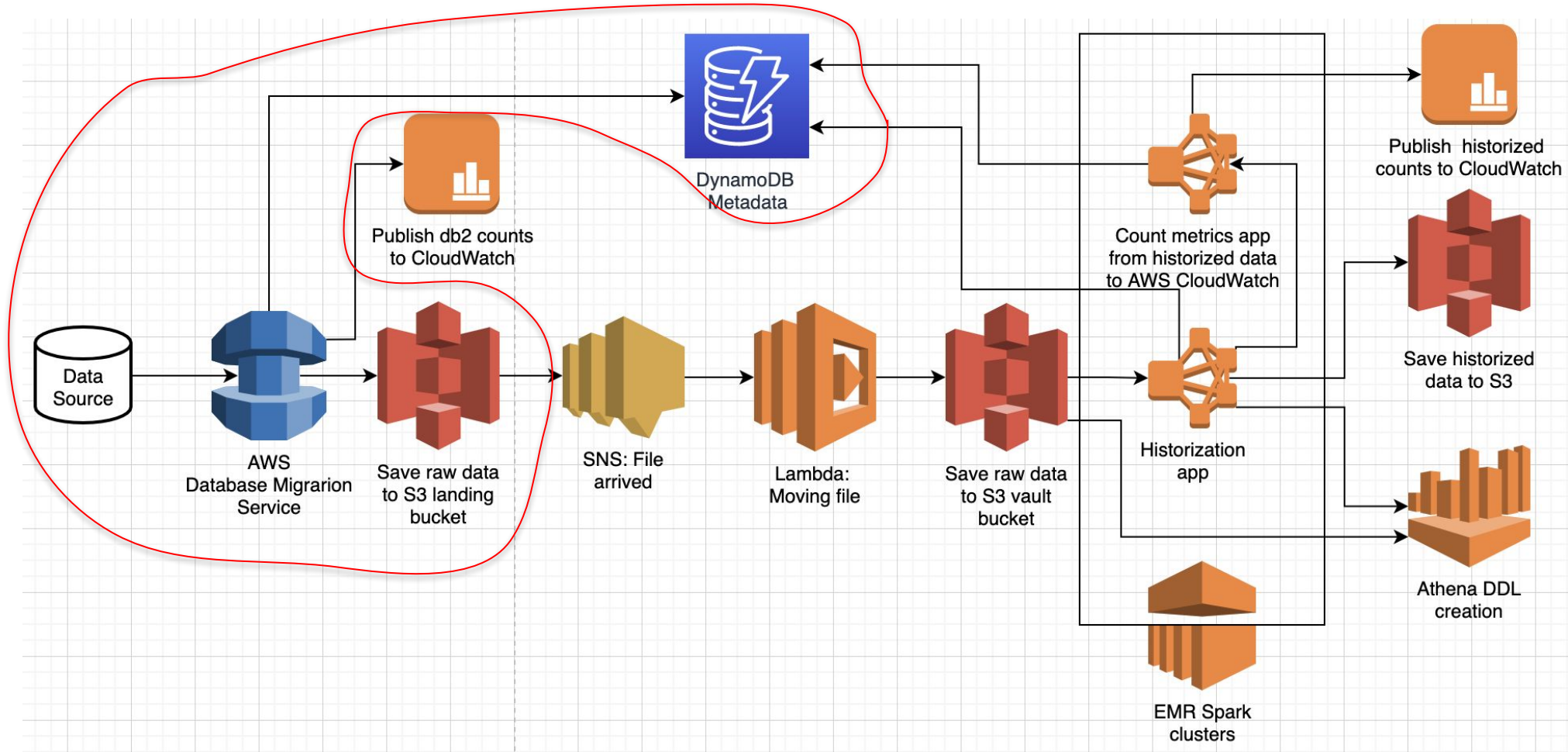


```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```

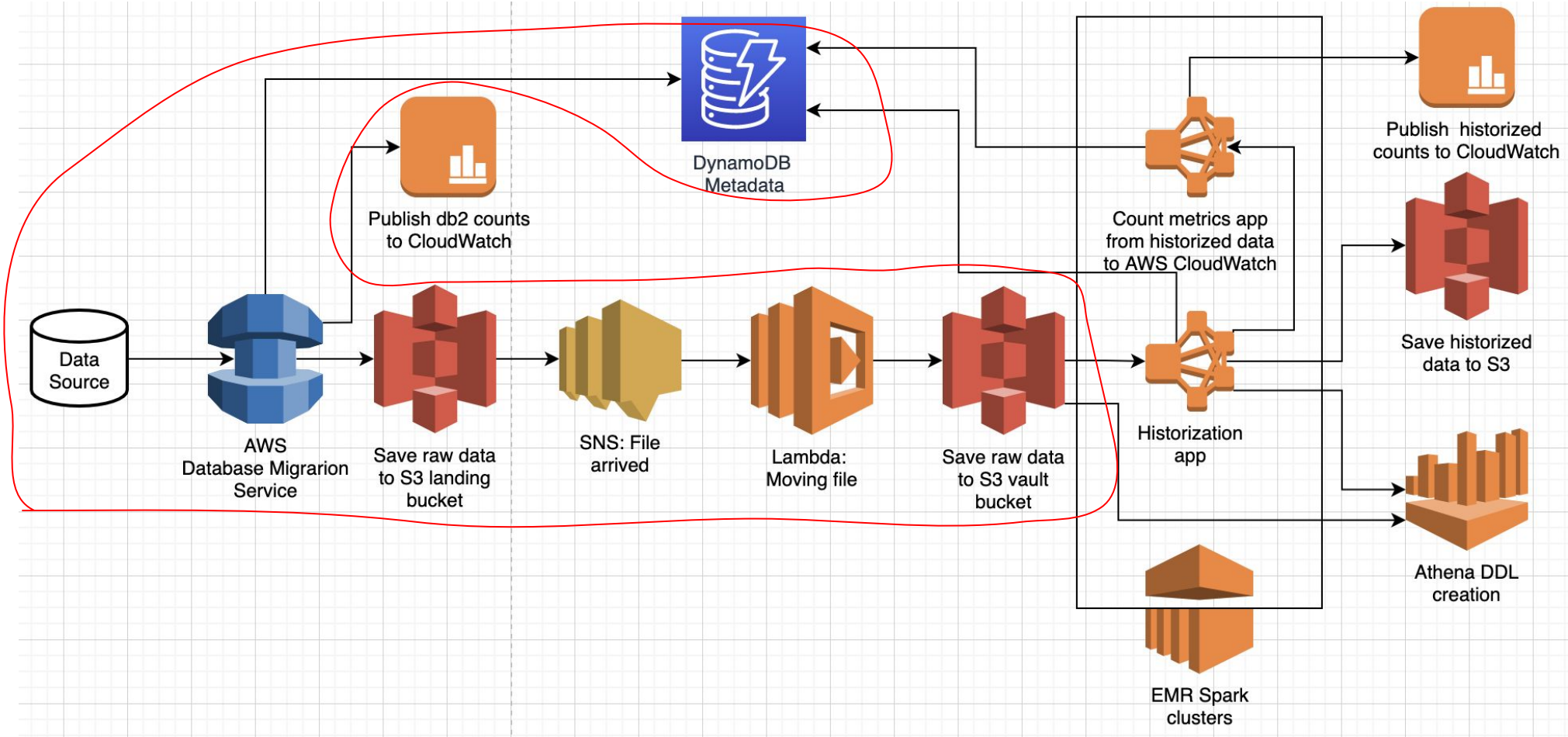


```
1 {  
2   "SchemaName": "schema_name",  
3   "TableName": "table_name",  
4   "Source": "db2",  
5   "Target": "s3",  
6   "StartValue": "1970-01-01 00:00:00.000000",  
7   "EndValue": "2021-05-11 00:00:00.000000",  
8   "PrimaryKeyFields": "pk_field",  
9   "TimestampColumnName": "",  
10  "OutputFormat": "parquet",  
11  "BrokenDateTimeColumns": "",  
12  "DecimalWithNullColumns": "",  
13  "IsAthenaTableRequired": true,  
14  "IsTargetPartitioned": false  
15 }
```

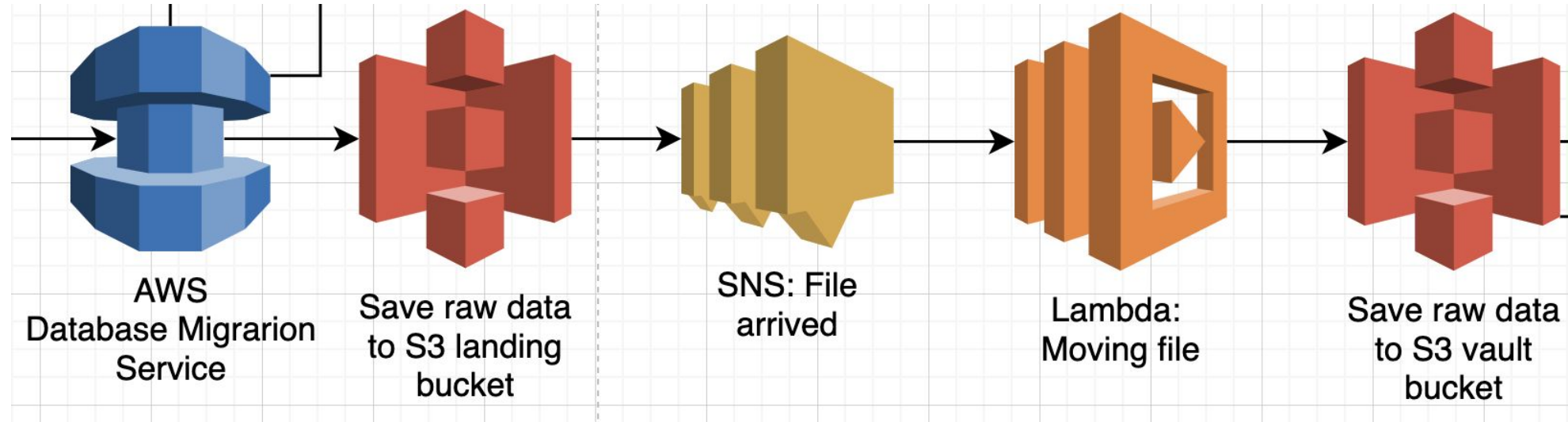

First part: Metadata management



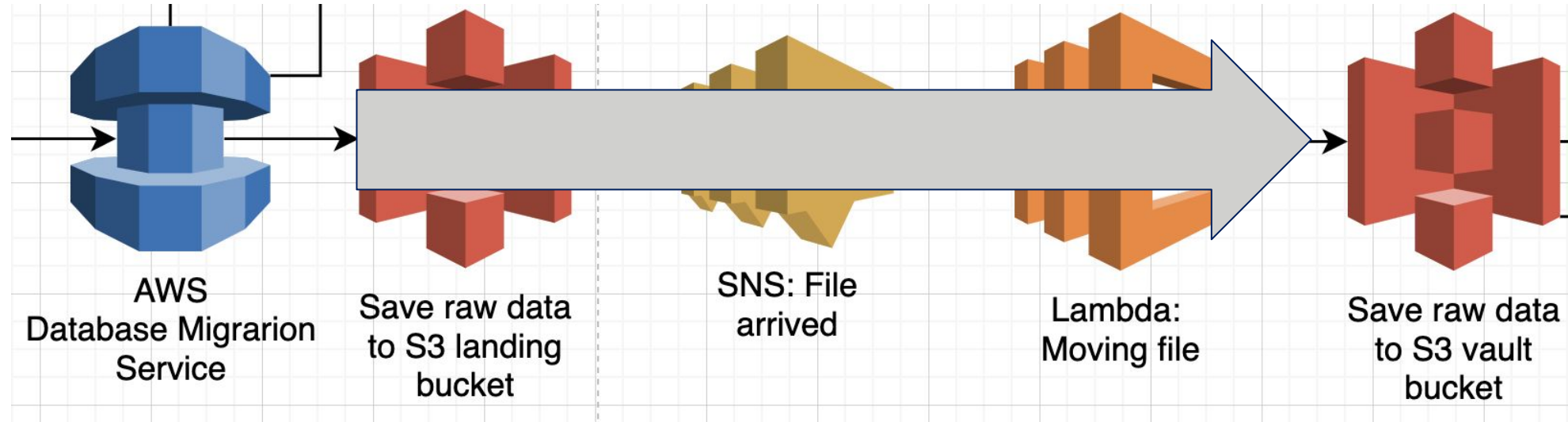
First part: DMS saving data



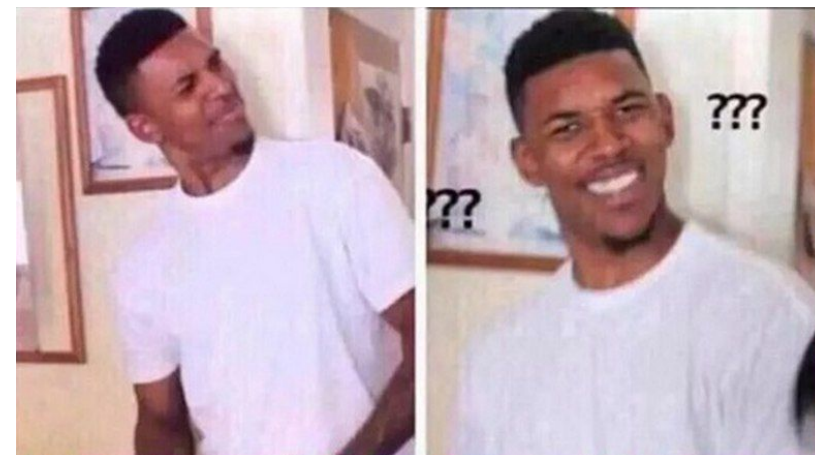
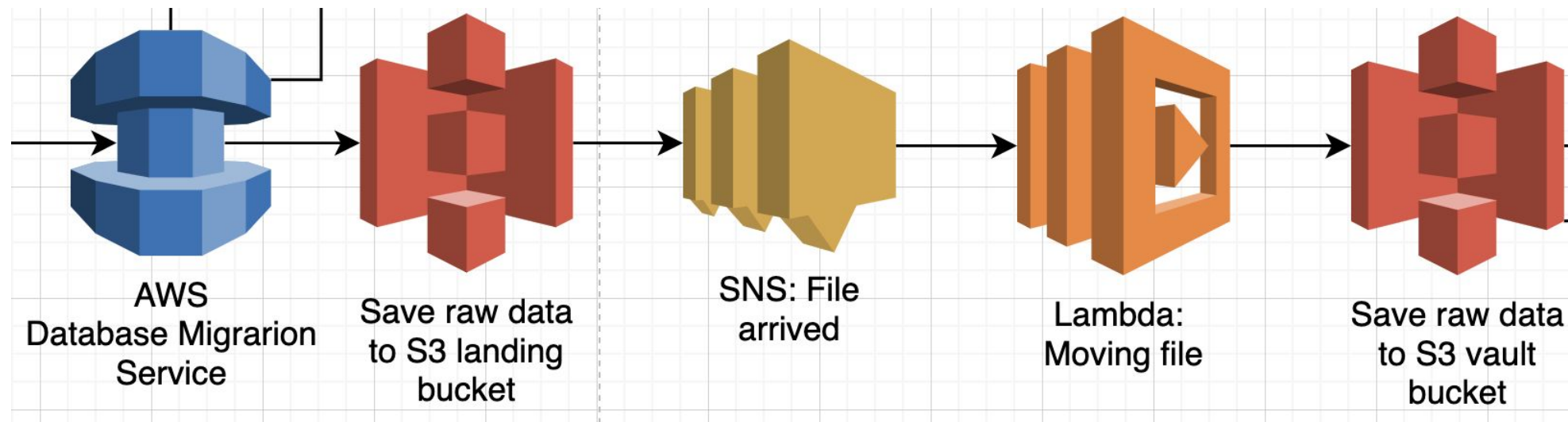
DMS: saving data



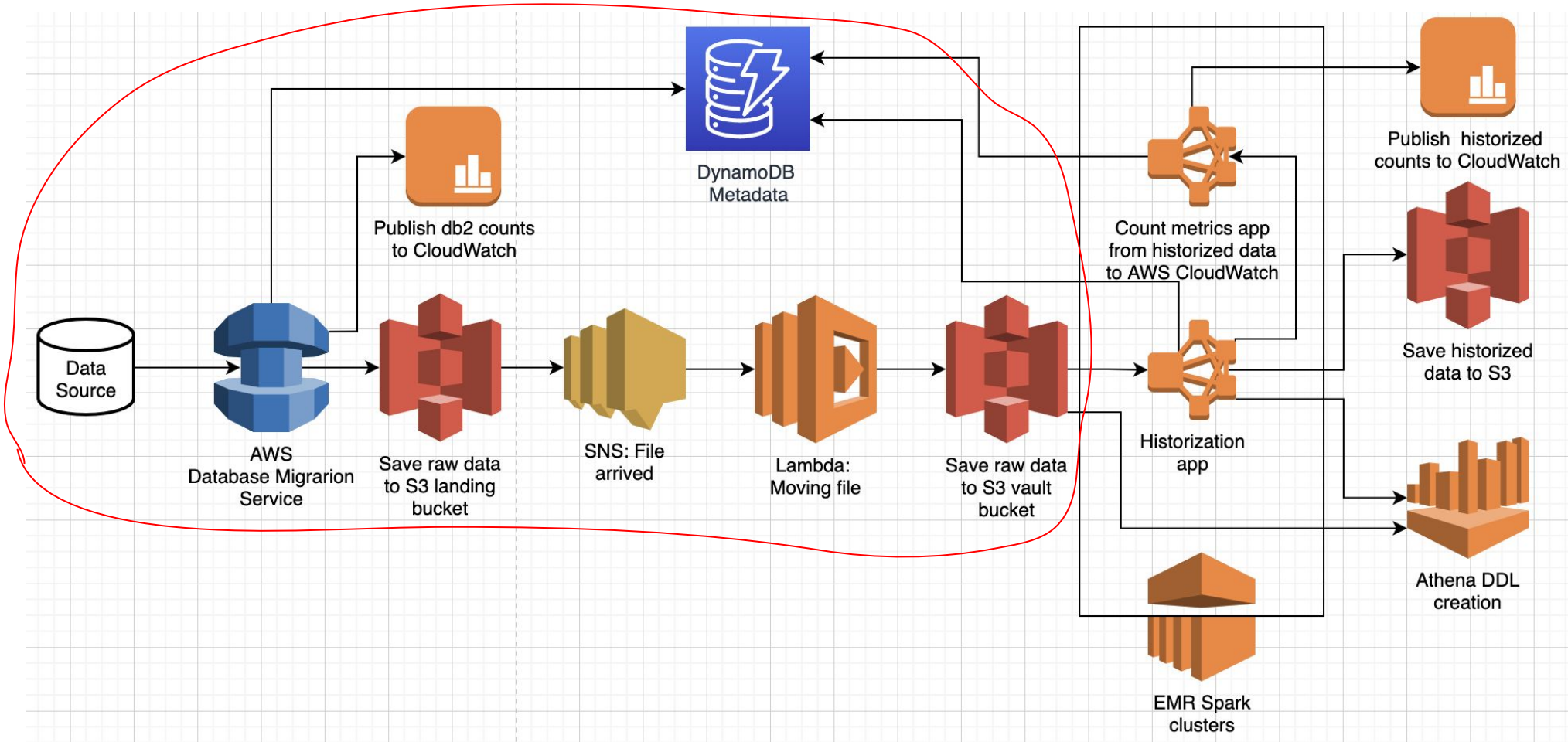
DMS: saving data

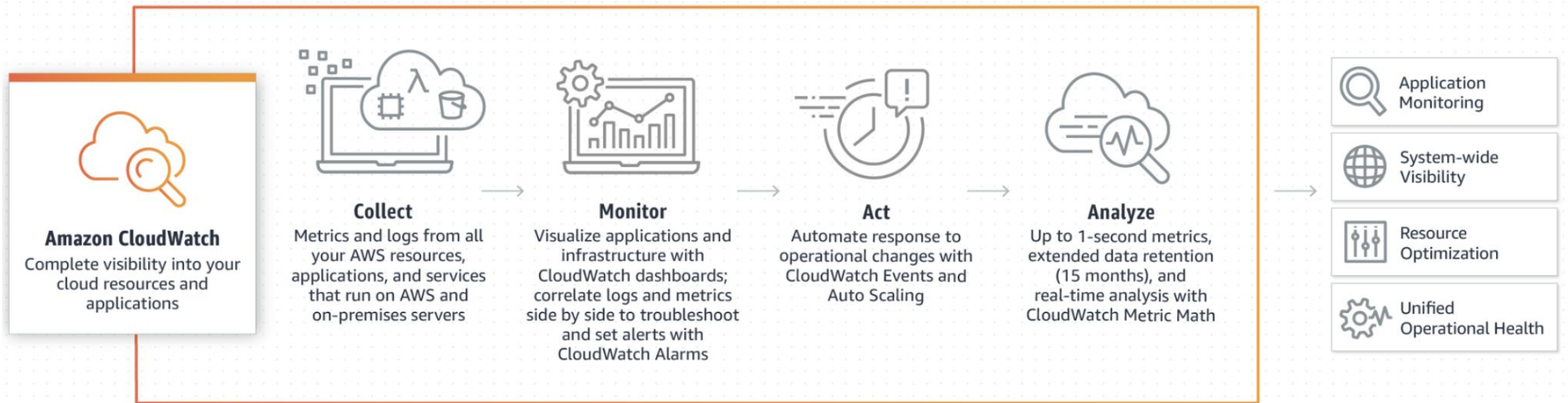


DMS: saving data



First part: DMS







Data Platform product
owner



Data Platform product
owner

He is happy when he
observes metrics

AWS CloudWatch



All metrics Graphed metrics (6/84) Graph options Source

Math expression ? Dynamic labels ?

Statistic: Average Period: (multiple) Remove all

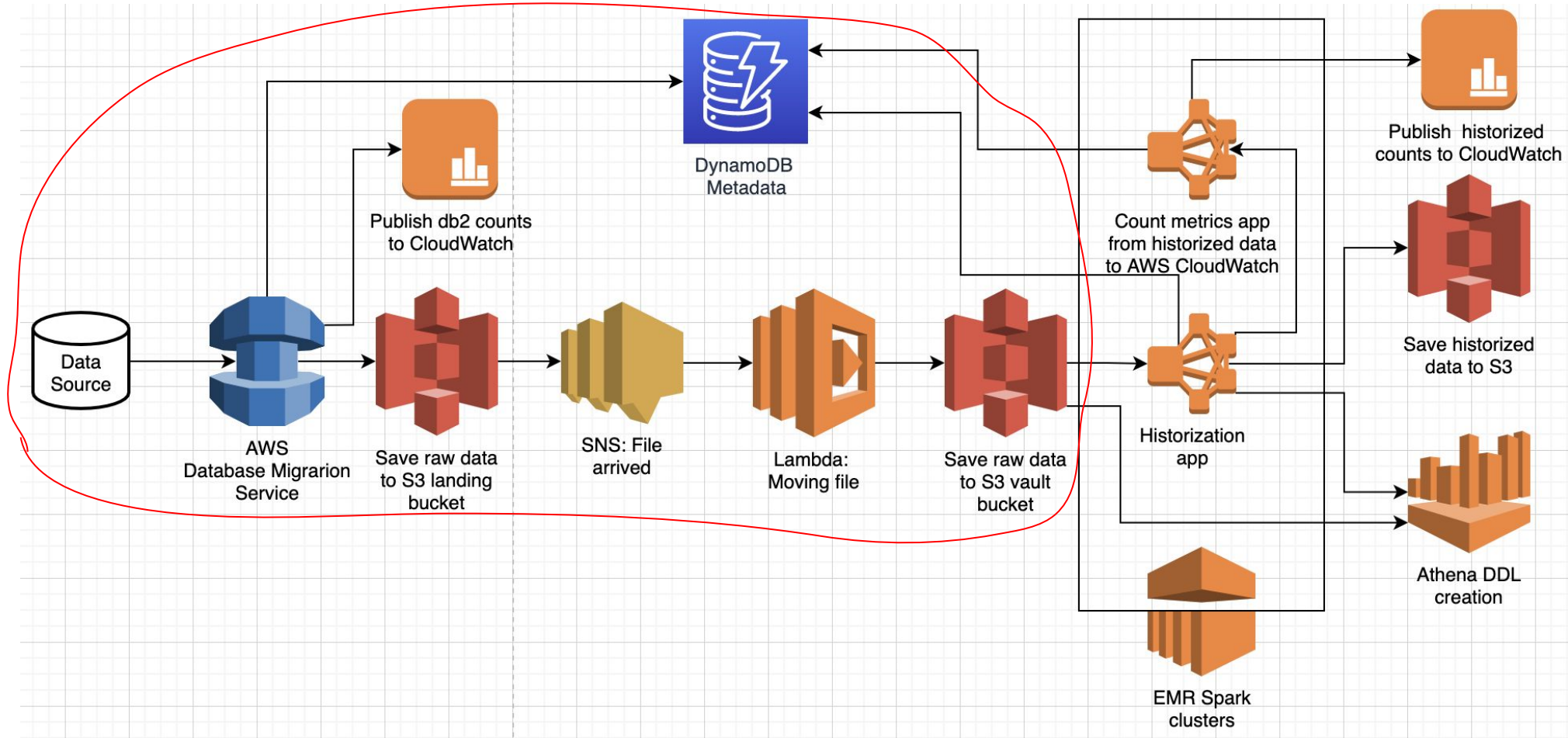
	Label	Details	Statistic	Period	Y Axis	Actions
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	
<input checked="" type="checkbox"/>		Historization • SourceRecordCount • TableName	Average	30 Days	< >	



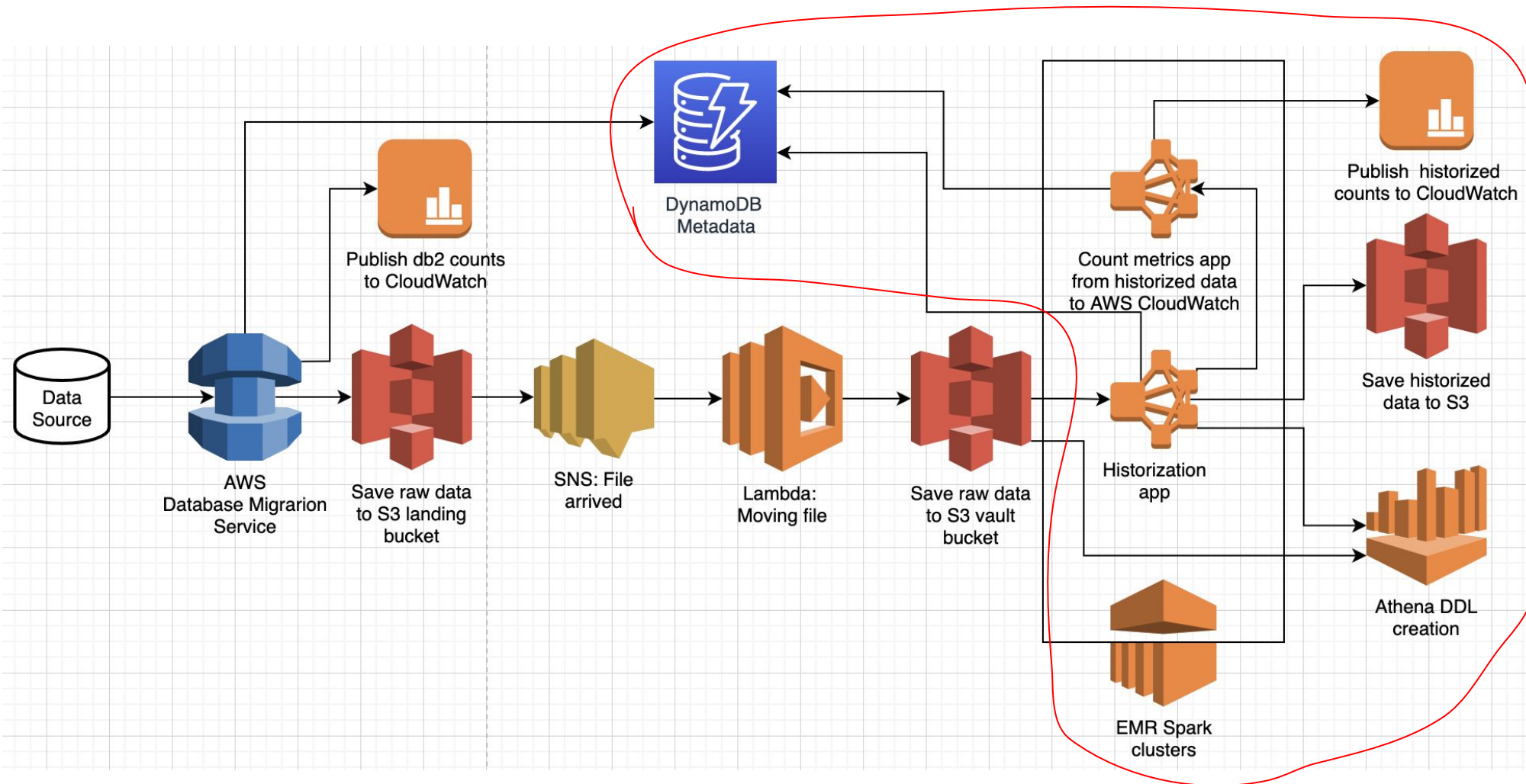
Data Platform product owner

He is happy when he observes metrics

First part: DMS



Second part: Historization



Historization: components



amazon
EMR



amazon
EMR



Amazon Athena

Historization: Athena



The screenshot displays the AWS Athena Query Editor interface. On the left sidebar, under 'Data source', 'AwsDataCatalog' is selected. Under 'Database', 'sampledb' is selected. In the 'Tables (2)' section, 'superstore' is highlighted with a red arrow. The main query editor shows a SQL statement to create an external table named 'superstore' in the 'sampledb' database. The query is as follows:

```
1 CREATE EXTERNAL TABLE IF NOT EXISTS sampledb.superstore (  
2   `order` int,  
3   `product` string,  
4   `qty` int,  
5   `amount` int  
6 )  
7 ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe'  
8 WITH SERDEPROPERTIES (  
9   'serialization.format' = ',',  
10  'field.delim' = ',',  
11 ) LOCATION 's3://ad-athena-datasets/data/'  
12 TBLPROPERTIES ('has_encrypted_data'='false');
```

Below the query editor, the 'Run query' button is highlighted with a red arrow. The 'Results' section at the bottom shows the message 'Query successful.' with a red arrow pointing to it. The interface also includes a search bar at the top, a navigation menu with 'Athena', 'Query editor', 'Saved queries', 'History', 'Data sources', and 'Workgroup : primary', and a 'Settings' menu with 'Tutorial', 'Help', and 'What's new'.

Historization: components



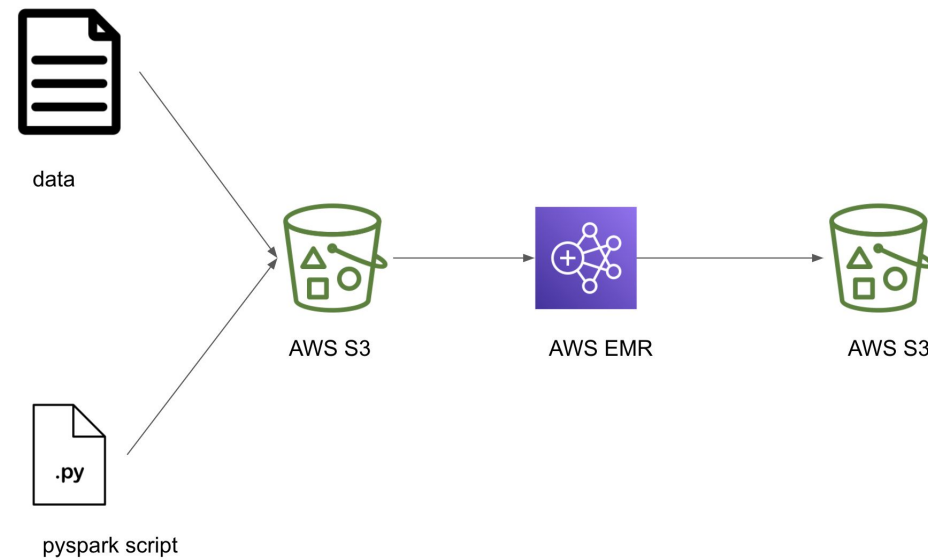
amazon
EMR



Amazon Athena

Spark jobs on EMR cluster

How it looks like in general



Historization: script input params

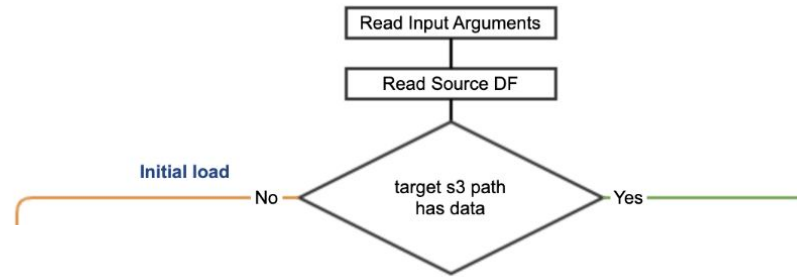


Input parameters

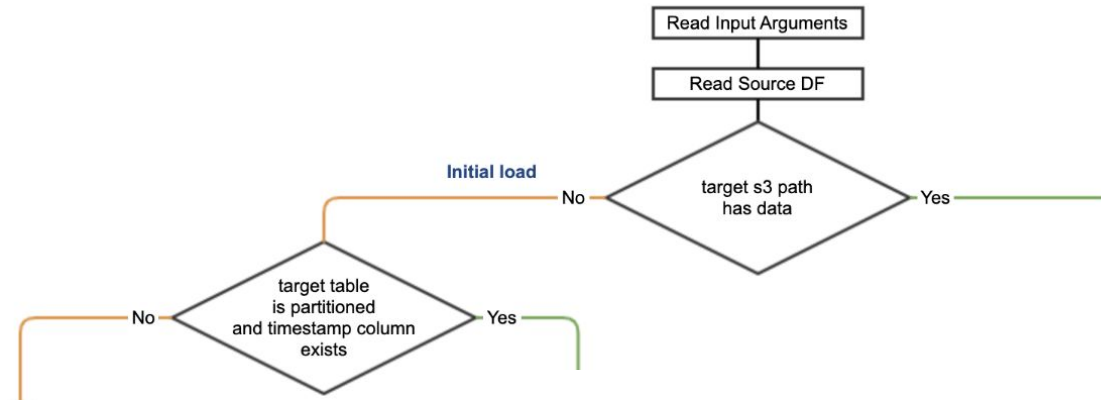
Name	Required	Type	Description
--source_s3_path	True	String	Source path for DMS outputs
--target_s3_path	True	String	Target path to the Delta Table
--load_date	True	String	Load date for source data in format YYYY-MM-DD
--primary_key_fields	True	String	List of fields (comma-separated string) used as primary key for source table
--timestamp_column_name	True	String	Timestamp column which will be used for historization to define valid_from and valid_to fields, default 'CTS'
--is_target_partitioned	False	Bool	True if target table is partitioned, default is 'True'
--is_athena_table_required	False	Bool	True if you want to create/update Athena table and you're sure that table doesn't contain PII data, default is 'False'

```
1 {
2   "SchemaName": "schema_name",
3   "TableName": "table_name",
4   "Source": "db2",
5   "Target": "s3",
6   "StartValue": "1970-01-01 00:00:00.000000",
7   "EndValue": "2021-05-11 00:00:00.000000",
8   "PrimaryKeyFields": "ID_VALUE",
9   "TimestampColumnName": "",
10  "OutputFormat": "parquet",
11  "BrokenDateTimeColumns": "",
12  "DecimalWithNullColumns": "",
13  "IsAthenaTableRequired": true,
14  "IsTargetPartitioned": false
15 }
```

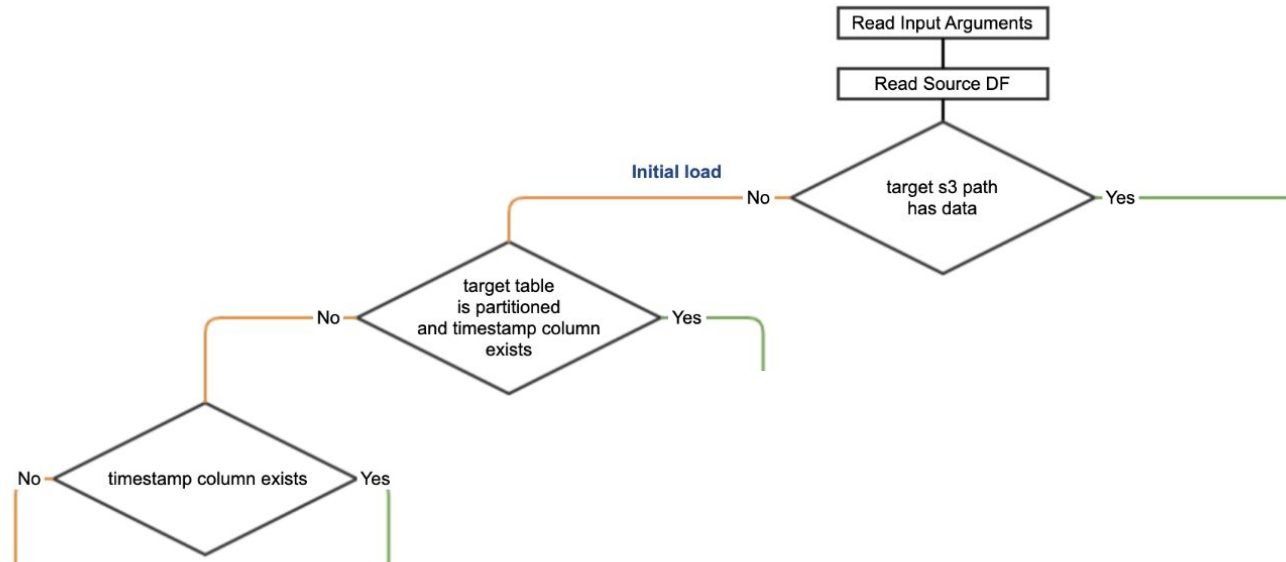

Historization: algorithm



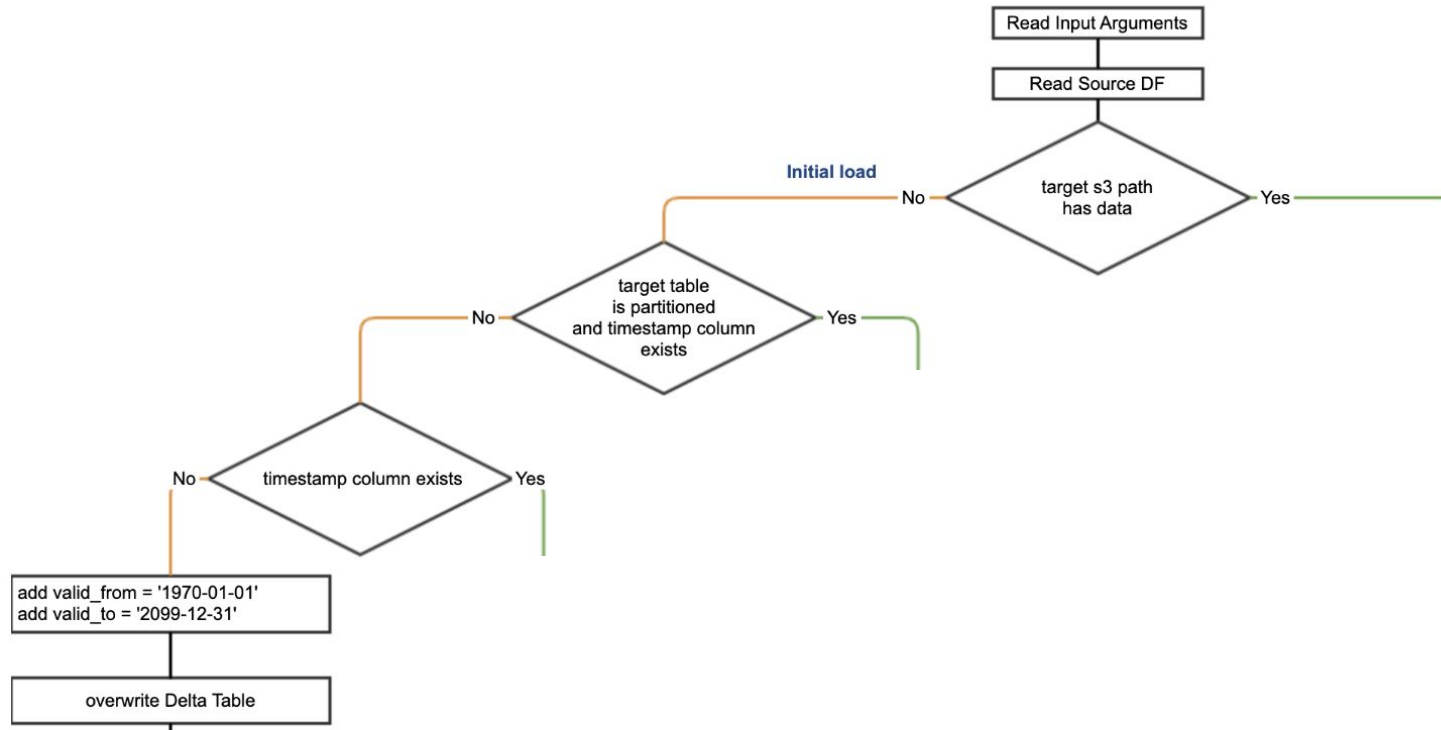
Historization: algorithm



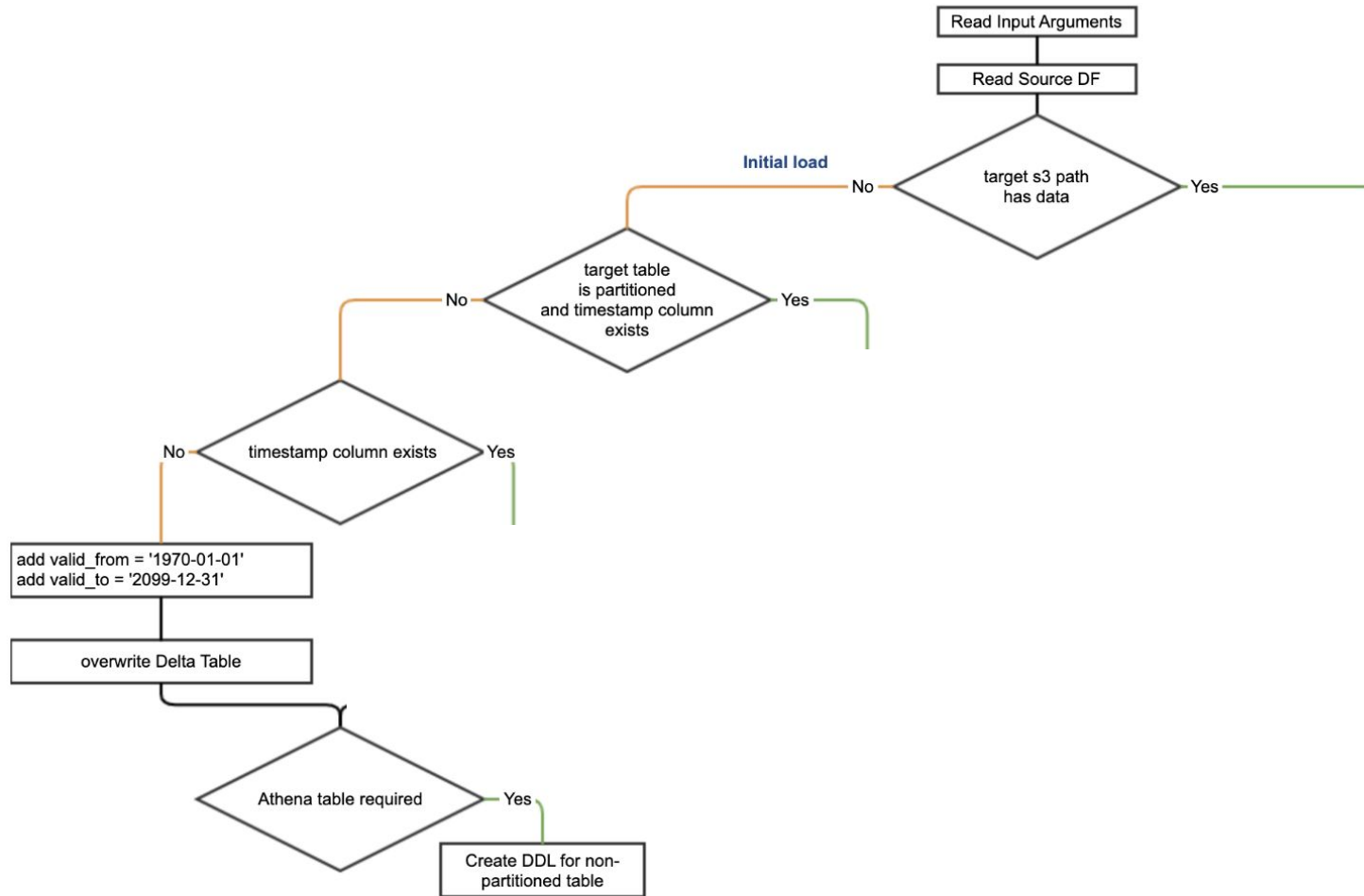
Historization: algorithm



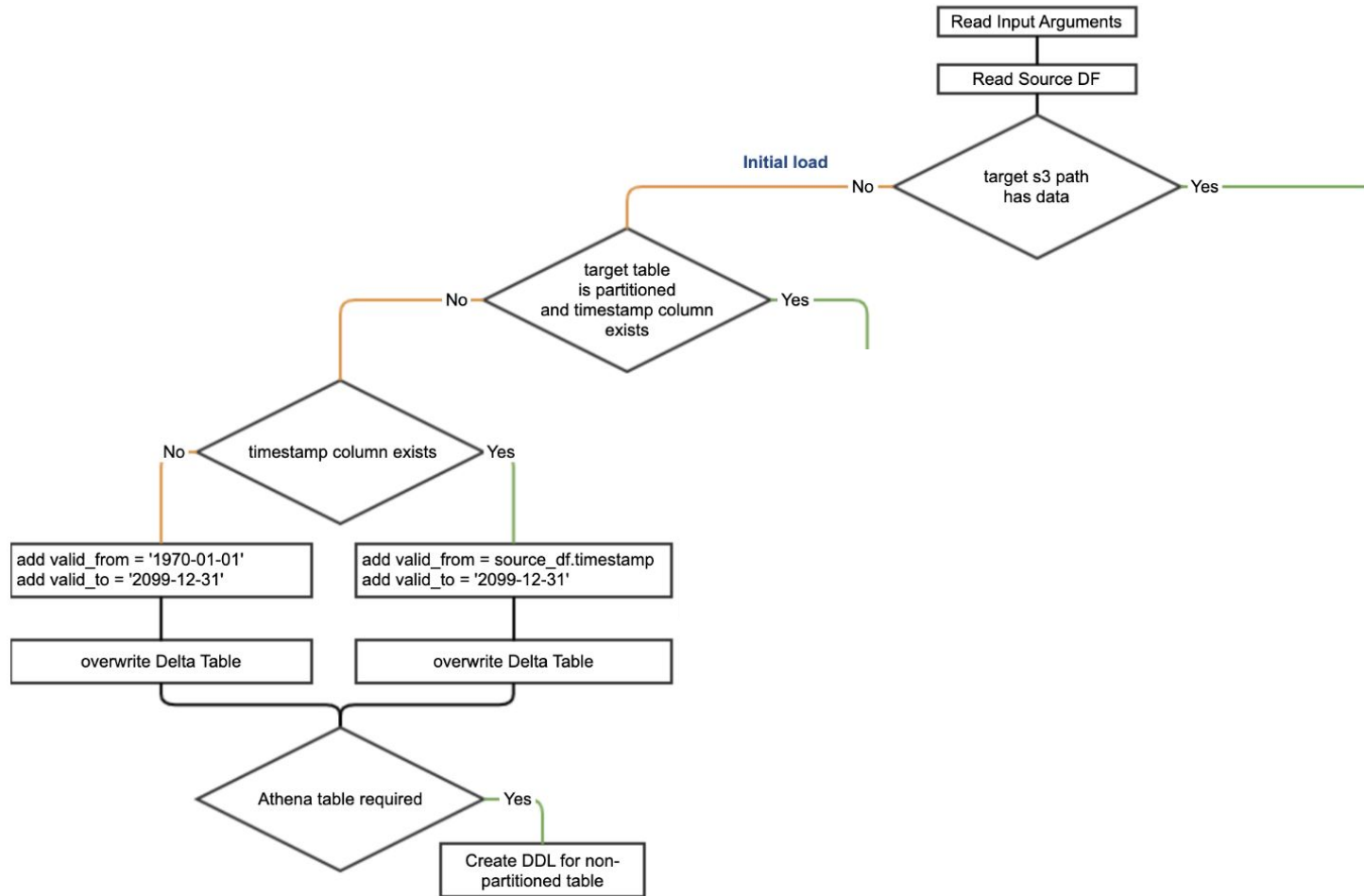
Historization: algorithm



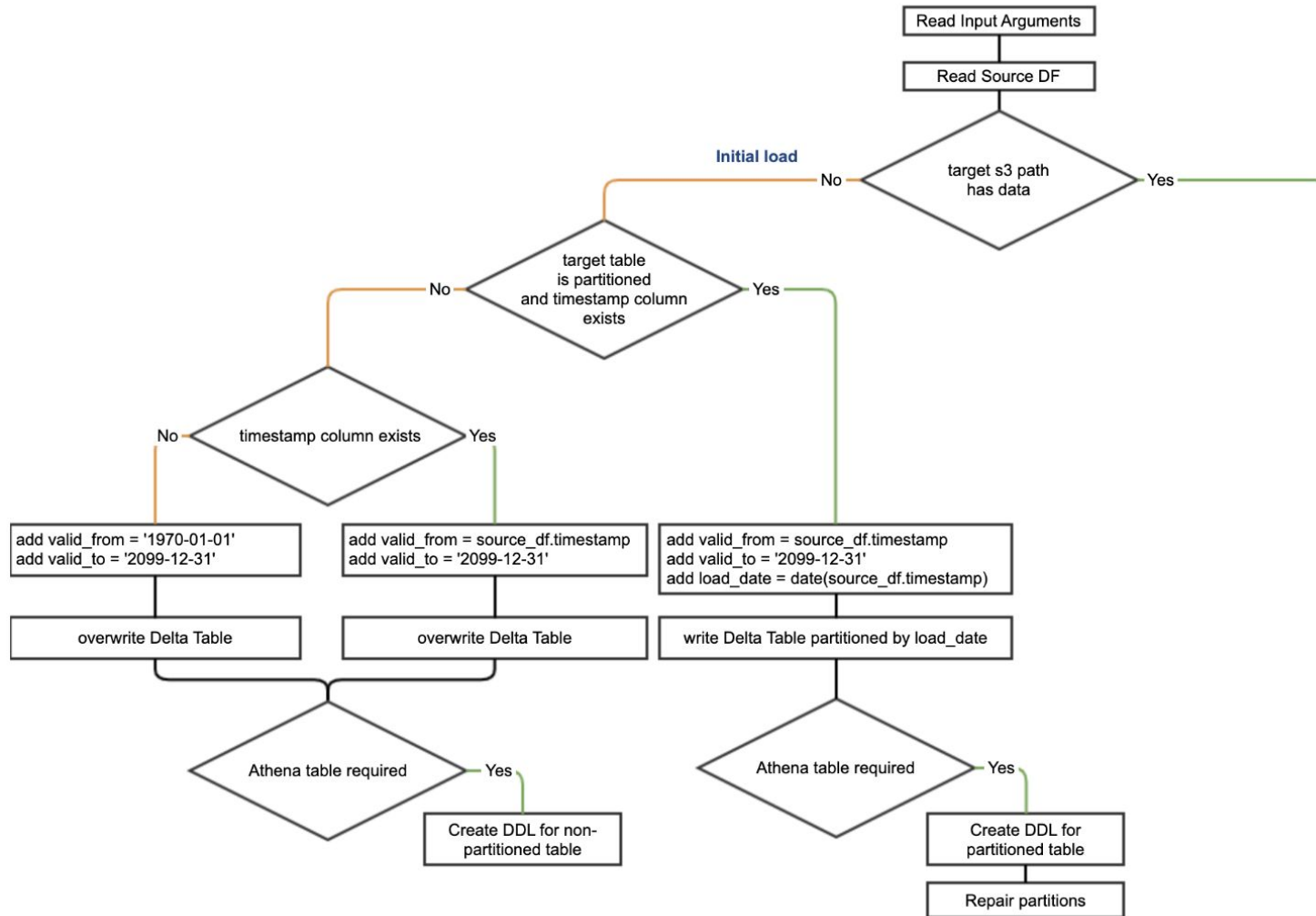
Historization: algorithm



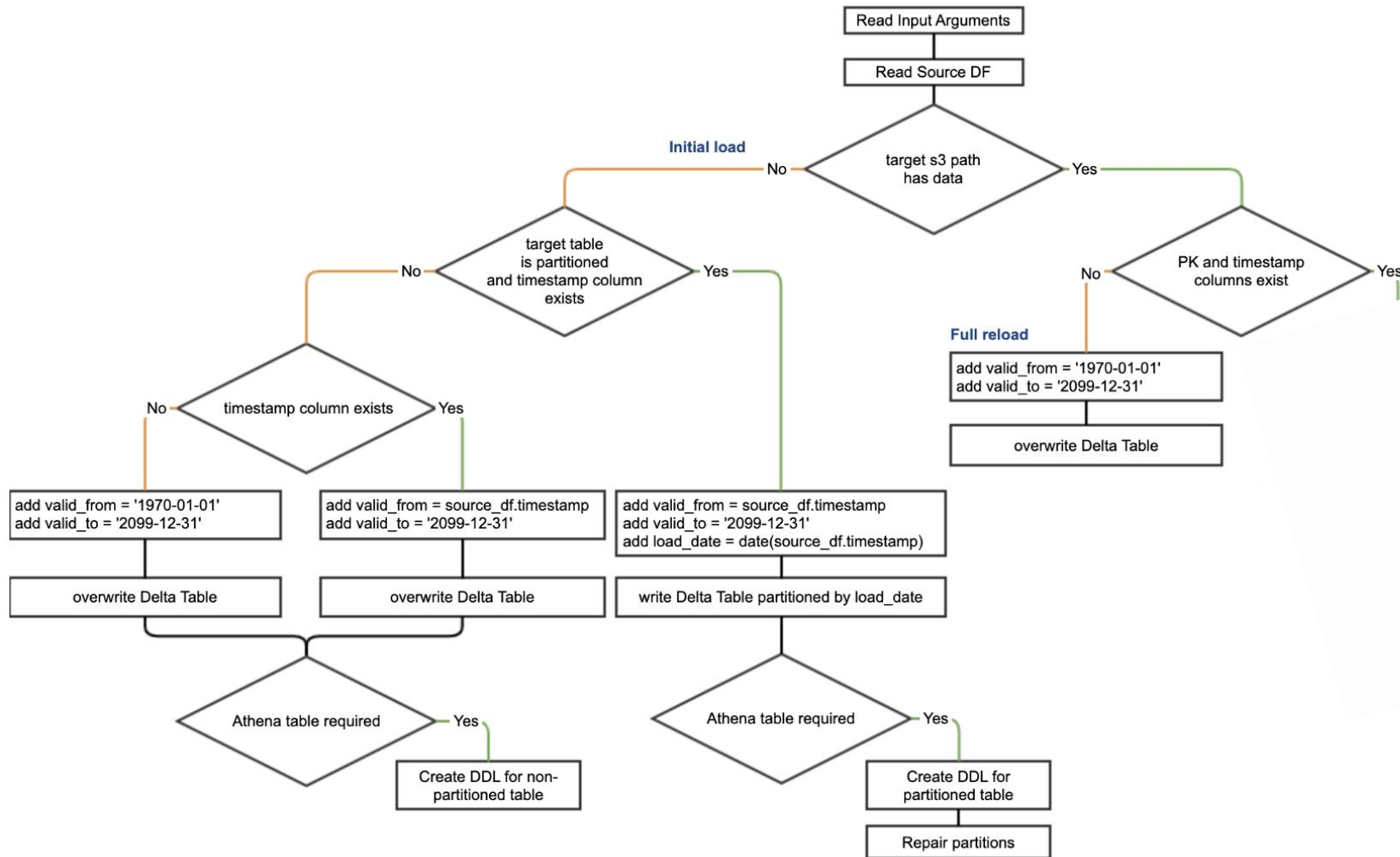
Historization: algorithm



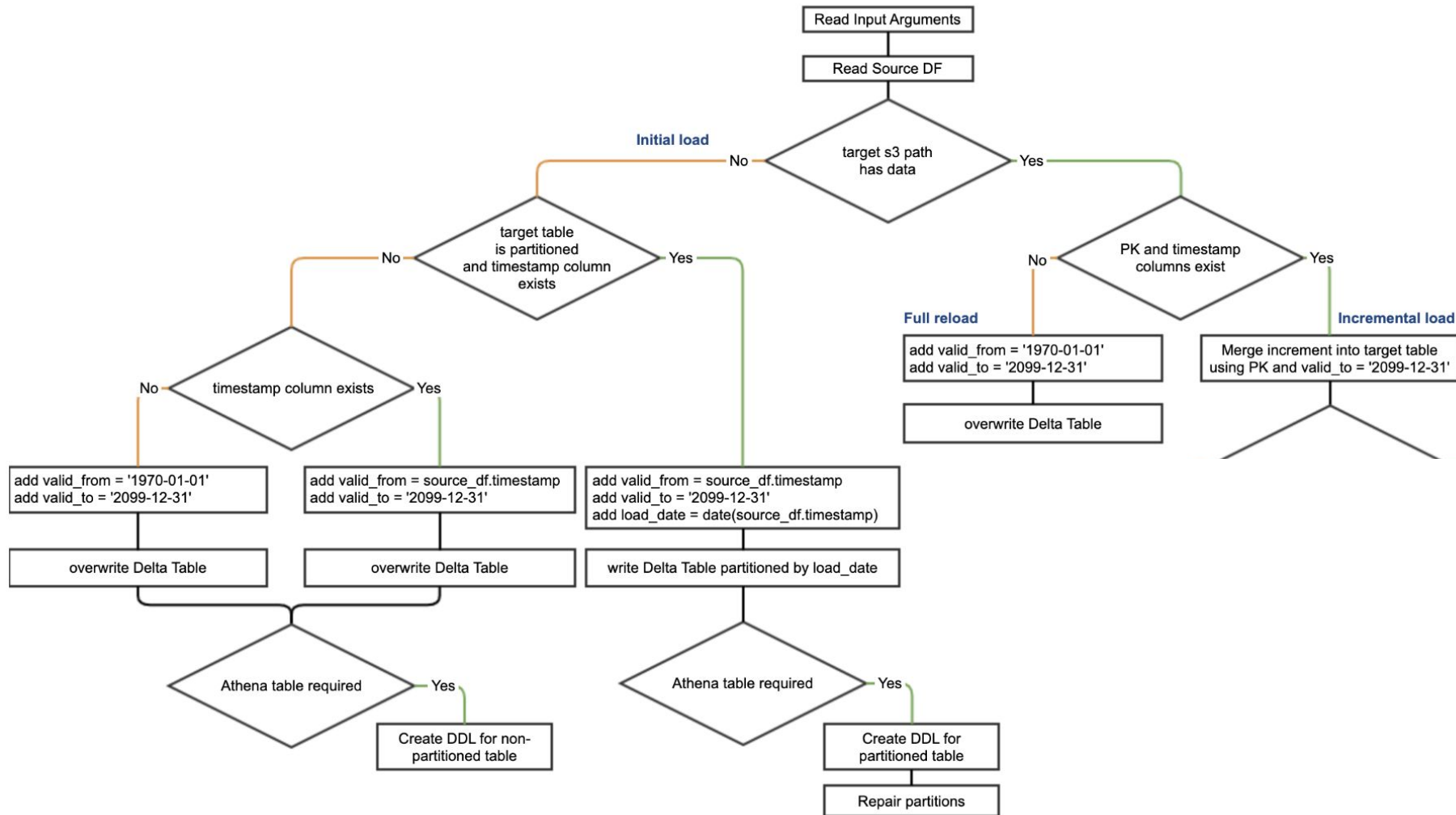
Historization: algorithm



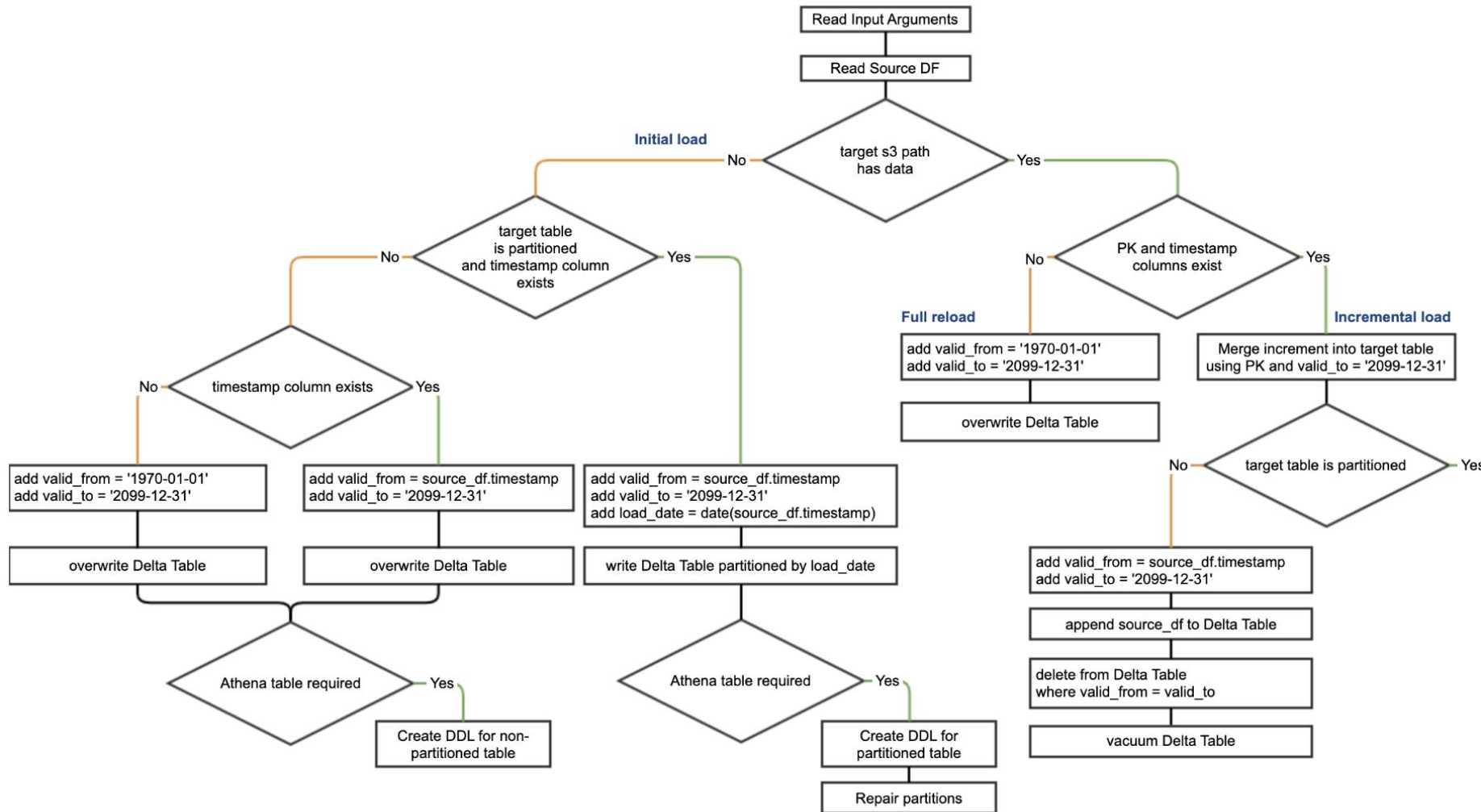
Historization: algorithm



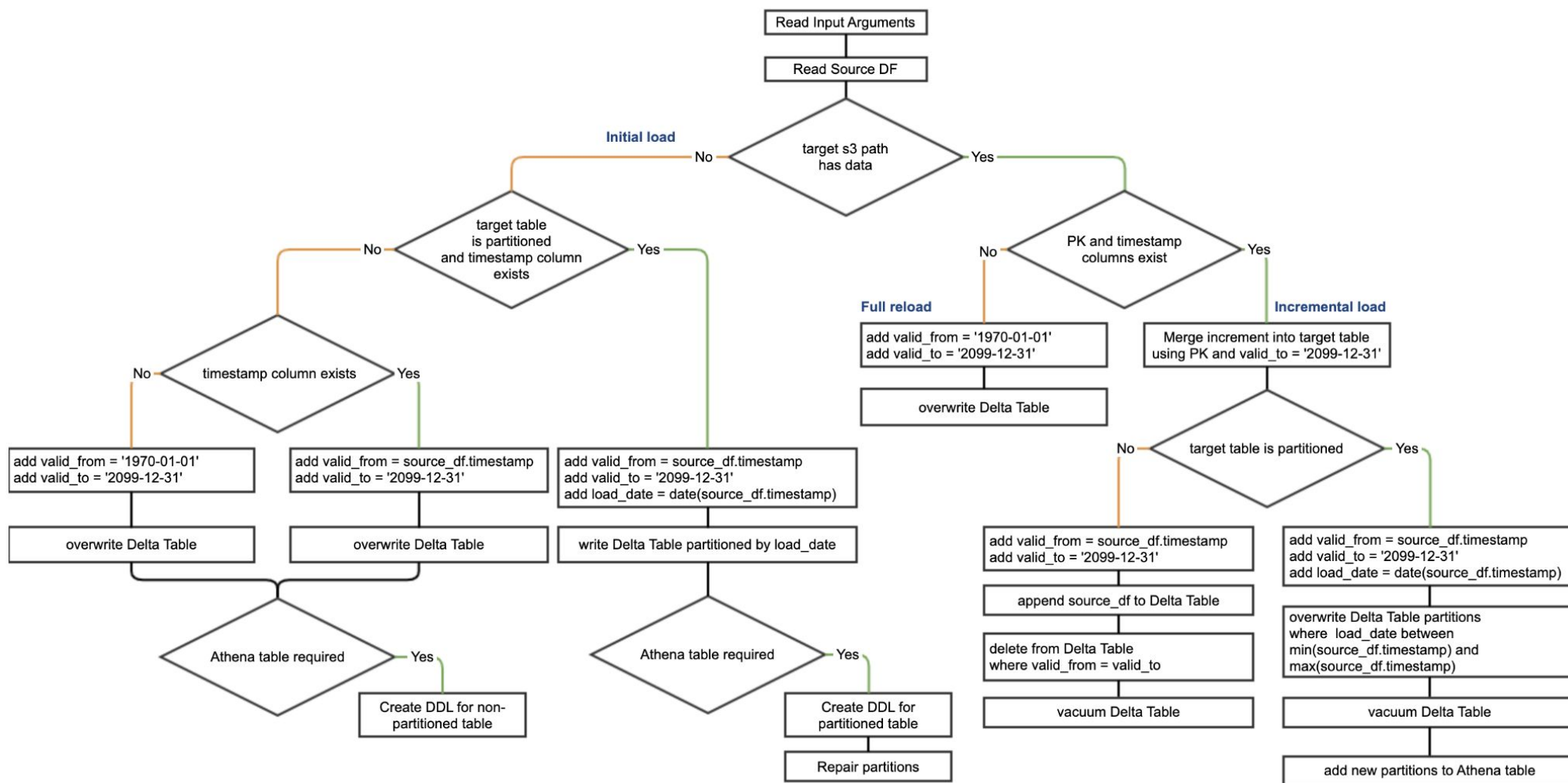
Historization: algorithm



Historization: algorithm



Historization: algorithm



Workflow management platform

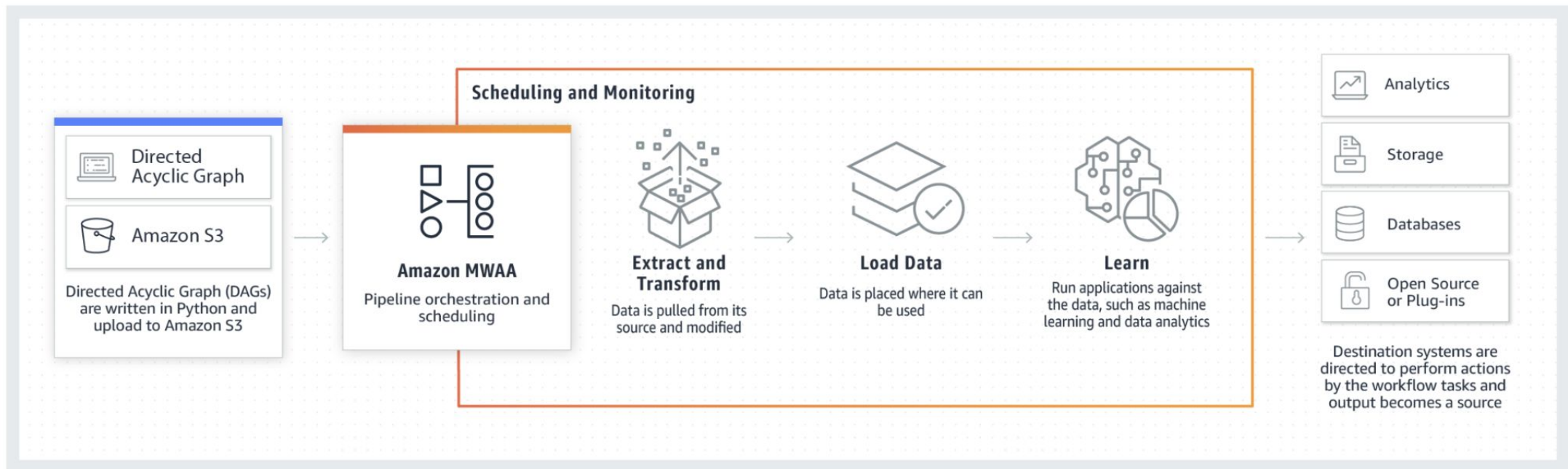


The workflow is orchestrated by MWAA

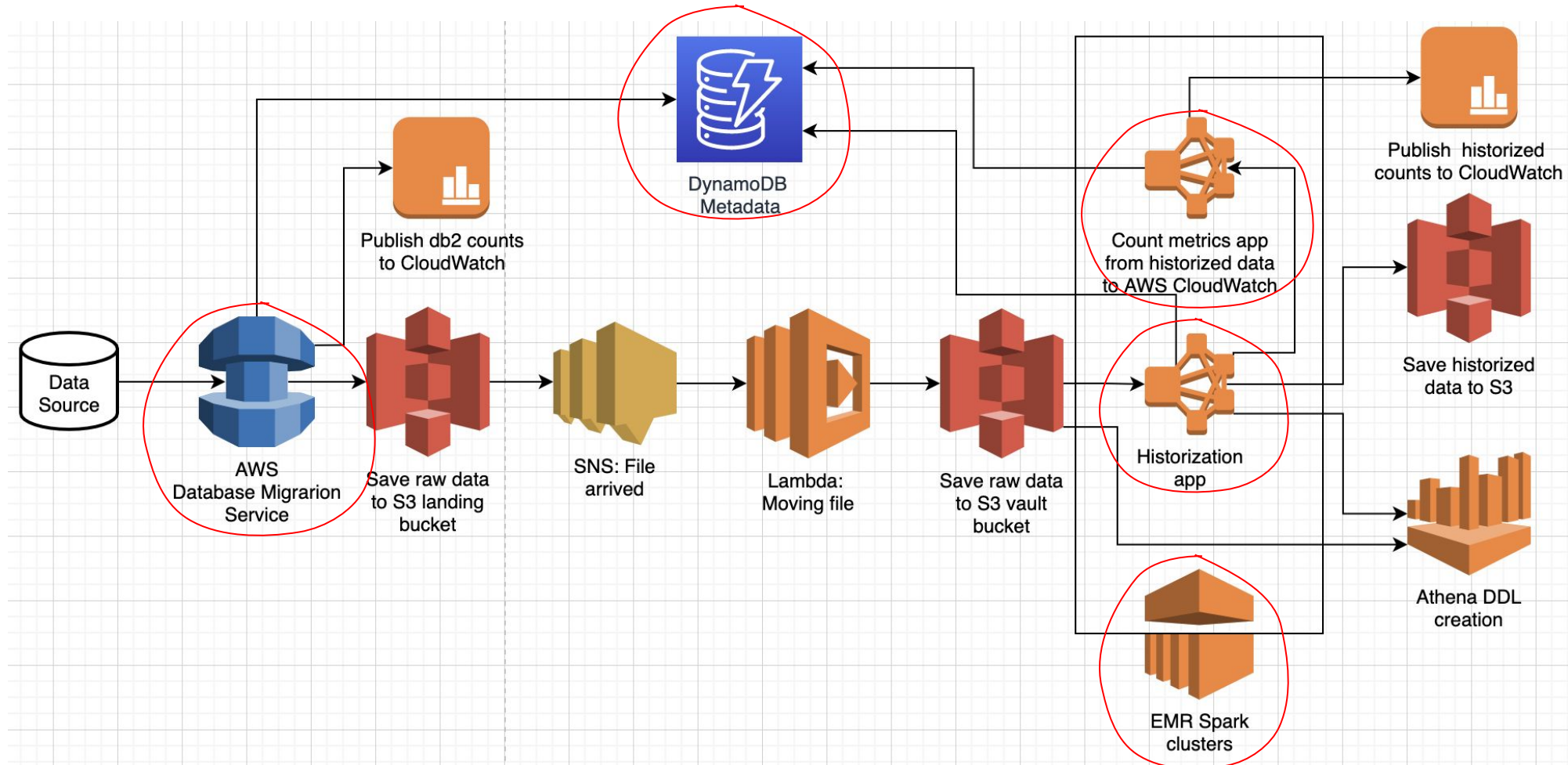


Workflow management platform

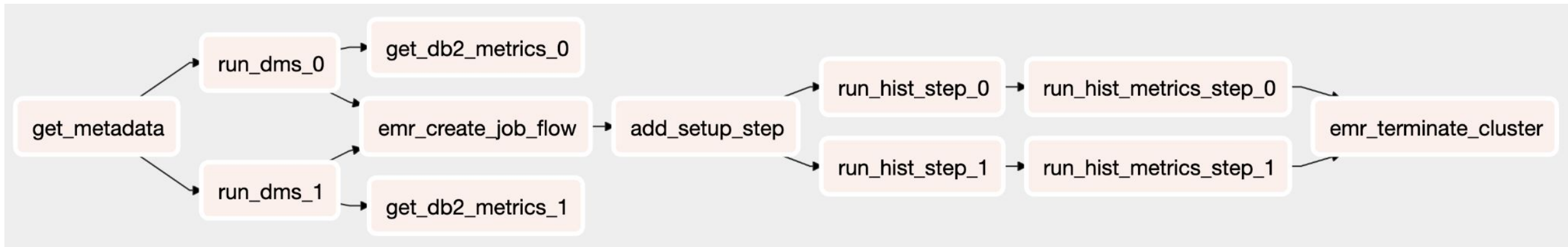
The workflow is orchestrated by MWAA



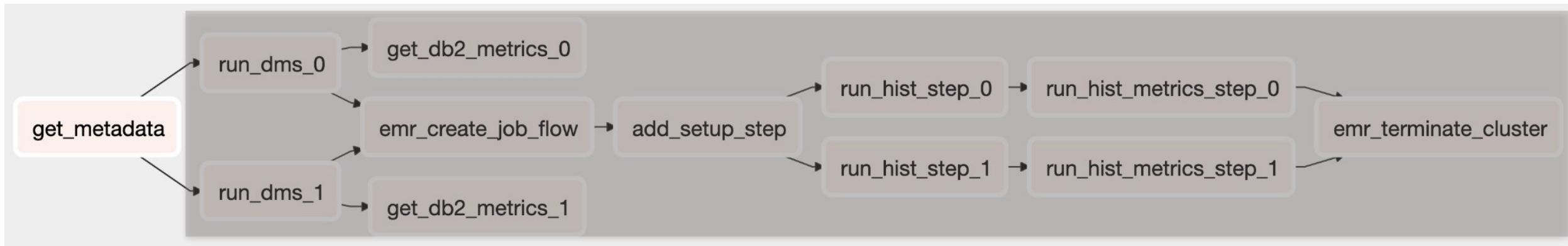
Airflow: Dag first iteration



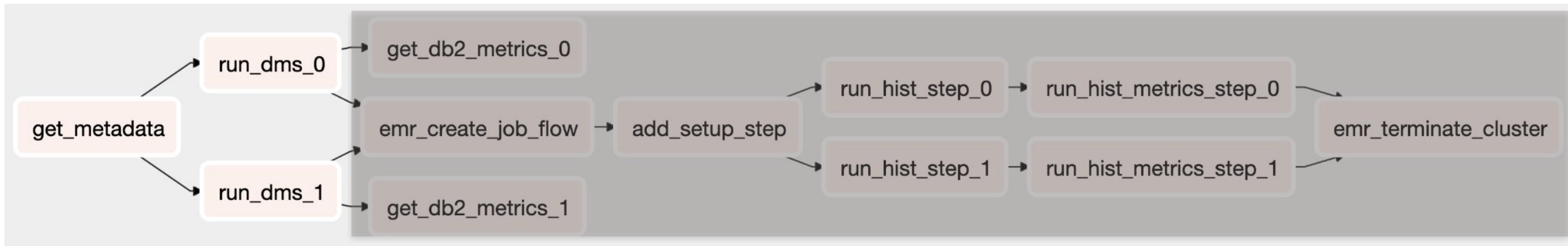
Airflow: Dag first iteration



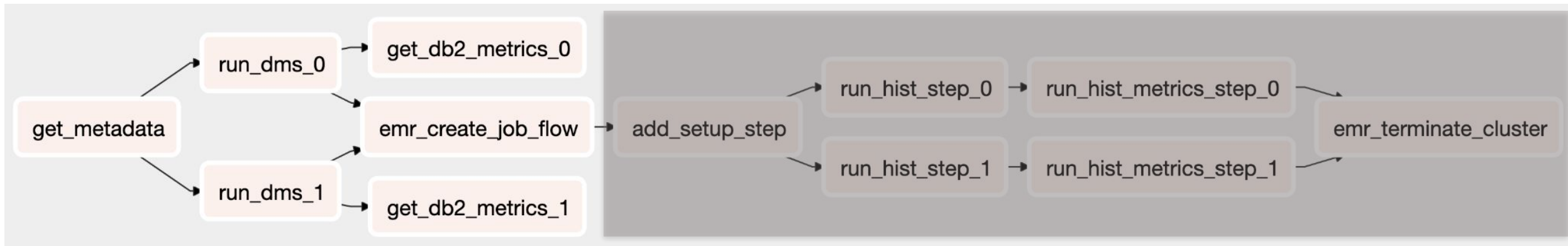
Airflow: Dag first iteration



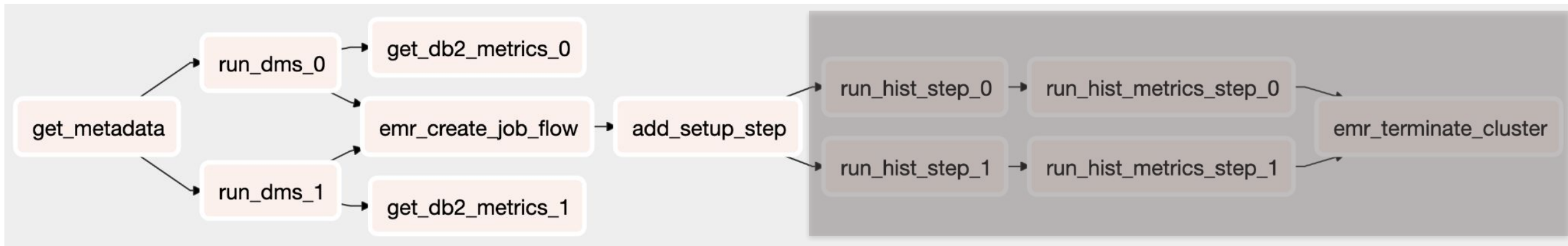
Airflow: Dag first iteration



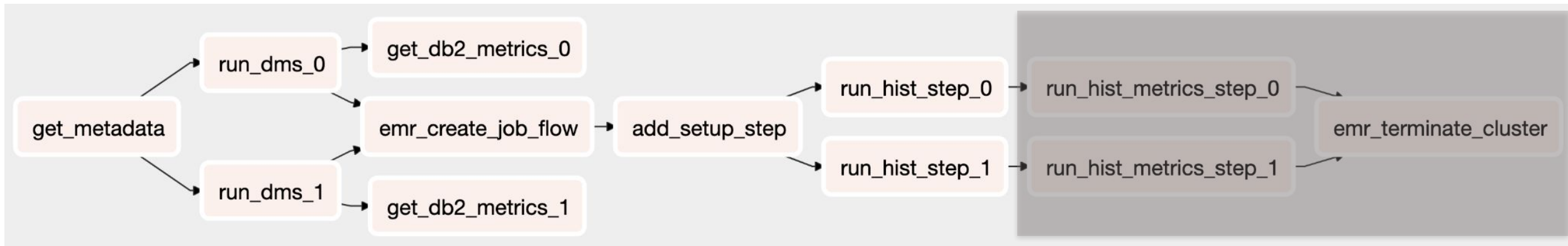
Airflow: Dag first iteration



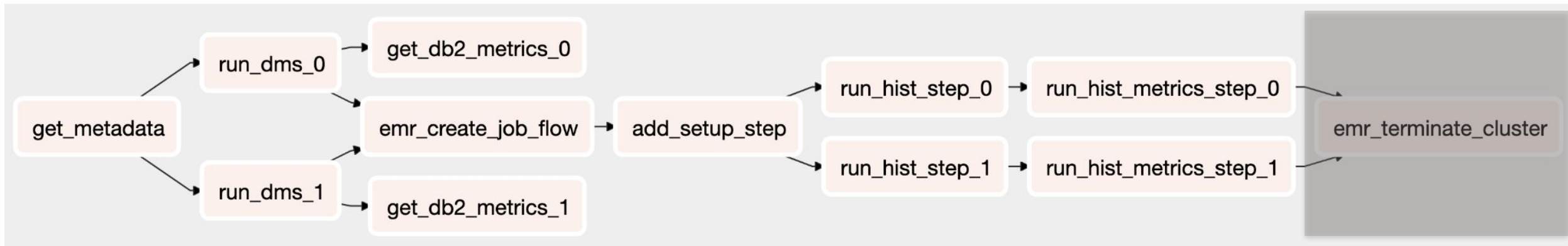
Airflow: Dag first iteration



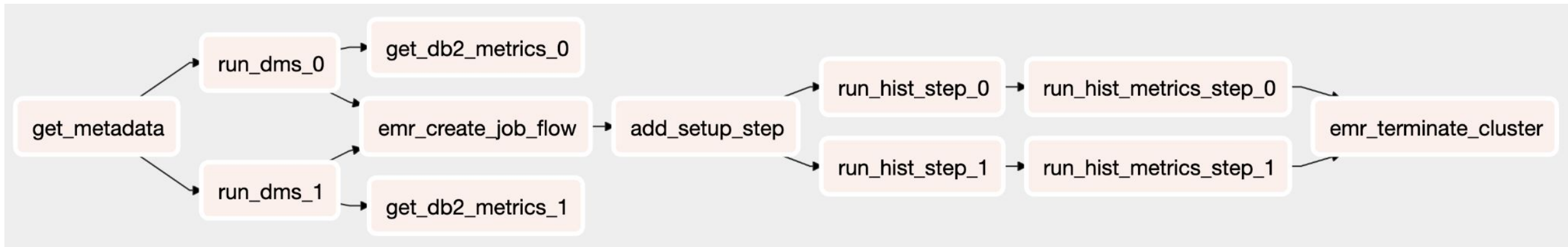
Airflow: Dag first iteration



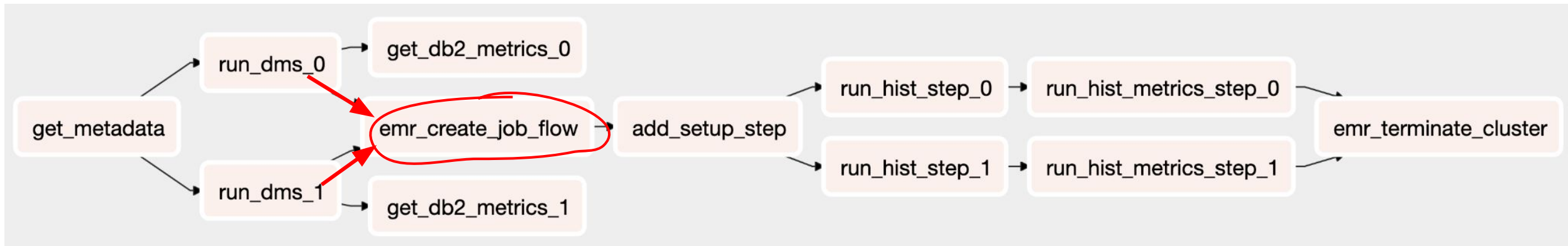
Airflow: Dag first iteration



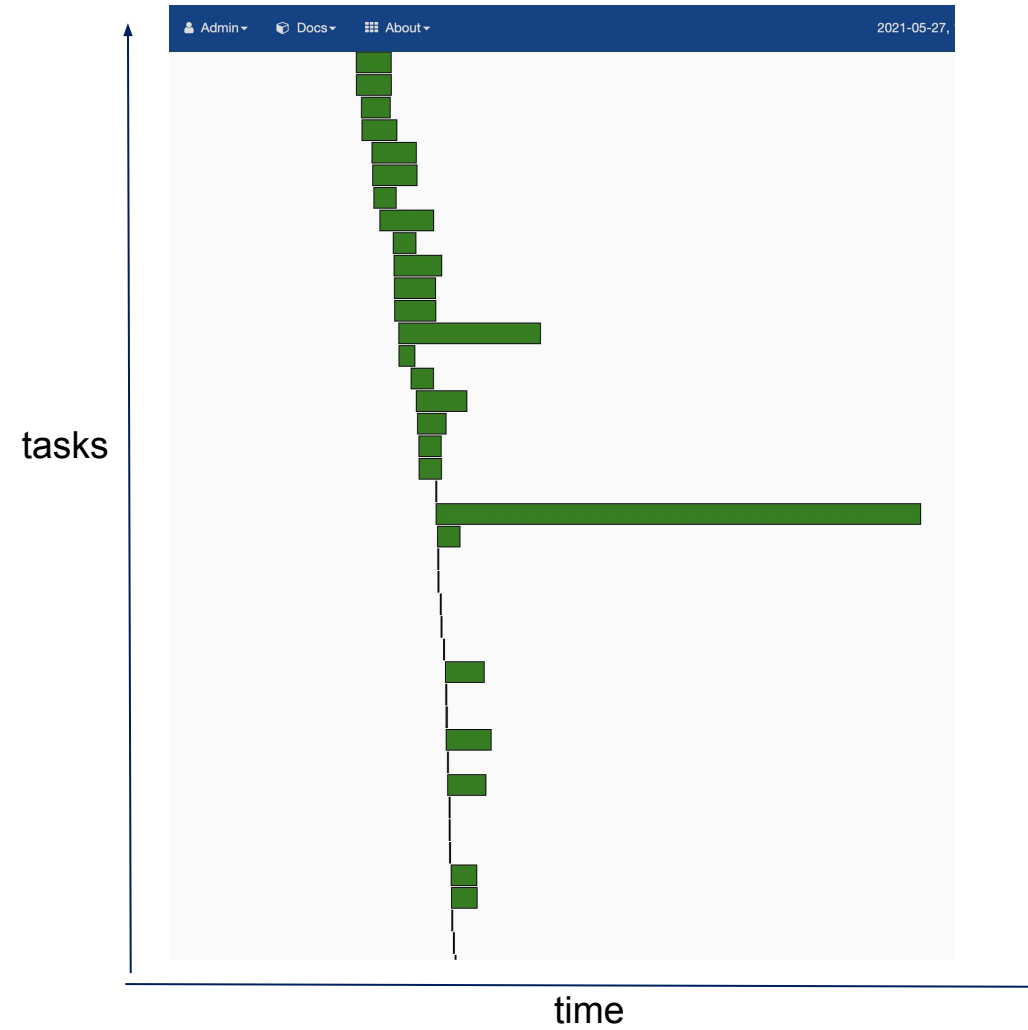
Airflow: Dag first iteration



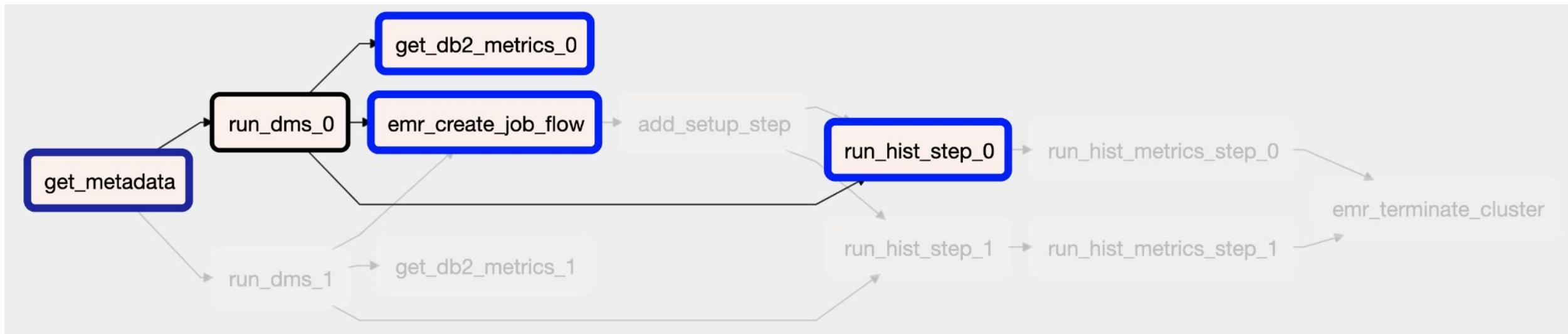
Airflow: Dag first iteration



Dag first iteration: Gantt diagramm



Airflow: Dag second iteration



Airflow: Monitoring and alerting




```
1 # Airflow DAG Config
2 DEFAULT_ARGS = {
3     'email': ['airflow@example.com'],
4     'email_on_failure': True,
5     'email_on_retry': False,
6     'on_failure_callback': send_zoom_failure_message,
7     'retries': 3,
8     'retry_delay': timedelta(minutes=3)
9 }
```



```
1 # Airflow DAG Config
2 DEFAULT_ARGS = {
3     'email': ['airflow@example.com'],
4     'email_on_failure': True,
5     'email_on_retry': False,
6     'on_failure_callback': send_zoom_failure_message,
7     'retries': 3,
8     'retry_delay': timedelta(minutes=3)
9 }
```



```
1 # Airflow DAG Config
2 DEFAULT_ARGS = {
3     'email': ['airflow@example.com'],
4     'email_on_failure': True,
5     'email_on_retry': False,
6     'on_failure_callback': send_zoom_failure_message,
7     'retries': 3,
8     'retry_delay': timedelta(minutes=3)
9 }
```


Airflow: points to improve



Cluster: [REDACTED] **Terminated with errors** Master node was terminated due to not enough capacity in the Spot Instance pool

EMR EC2 Spot Instances

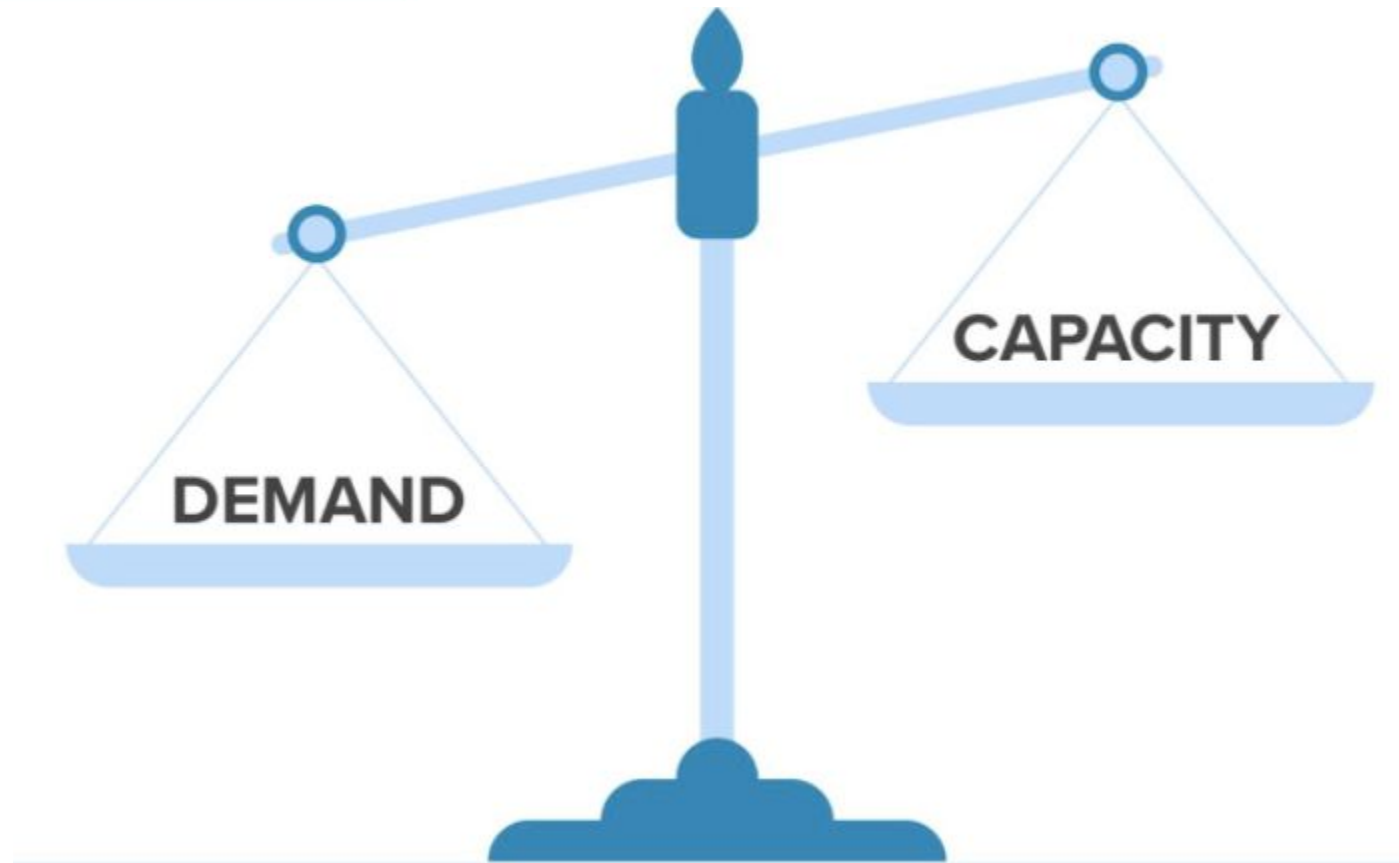
Name your own price for EC2 Compute

- A market where price of compute changes based upon Supply and Demand
- When Bid Price exceeds Spot Market Price, instance is launched
- Instance is terminated (with 2 minute warning) if market price exceeds bid price
- Unused On-Demand Instances

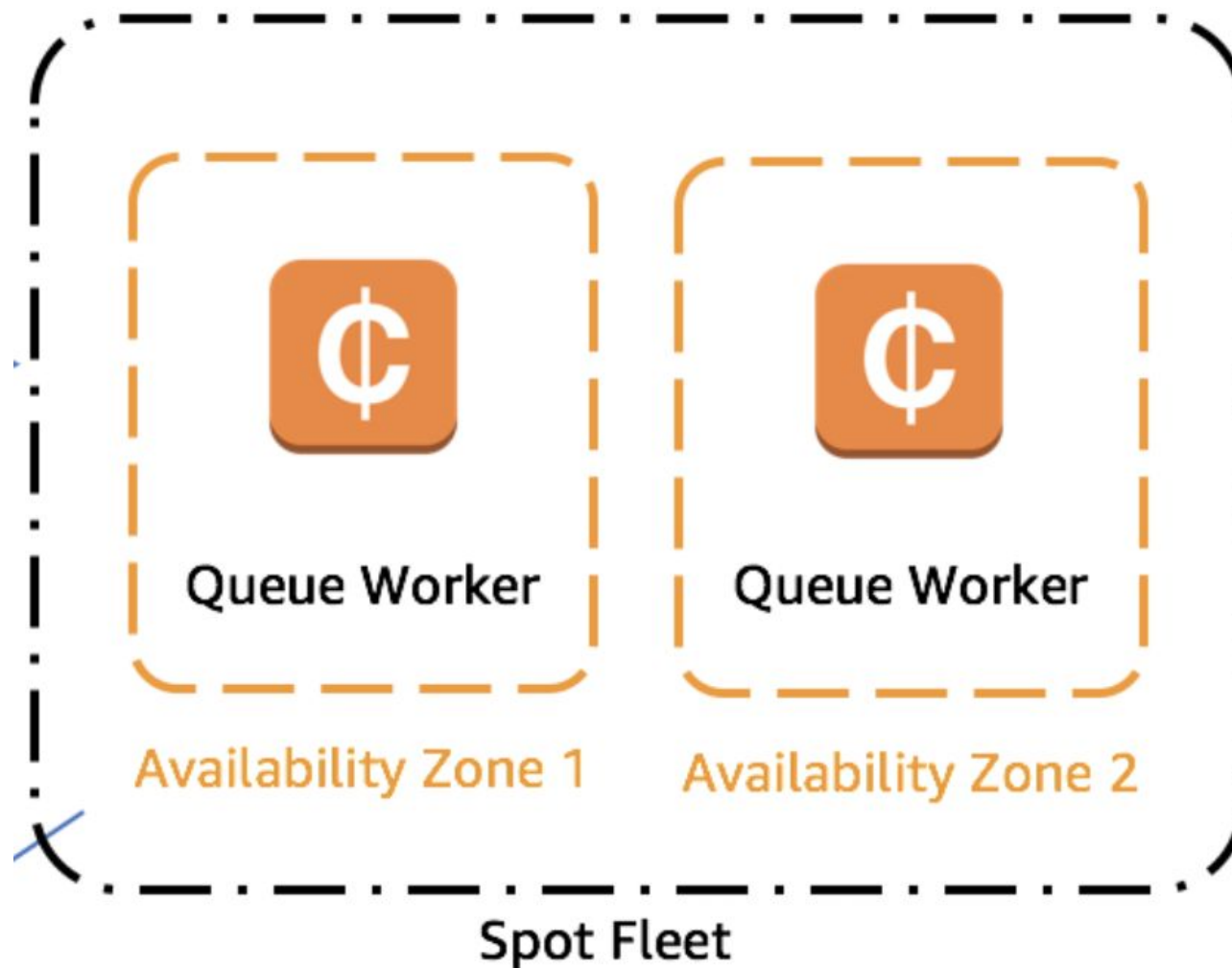


Airflow: points to improve

Cluster: [REDACTED] **Terminated with errors** Master node was terminated due to not enough capacity in the Spot Instance pool



- The instance fleets configuration for a cluster offers the widest variety of provisioning options for EC2 instances
- With instance fleets, you specify target capacities for On-Demand Instances and Spot Instances within each fleet.



That's it?



That's it?



Data type issues



- `time_millis`

- time_millis

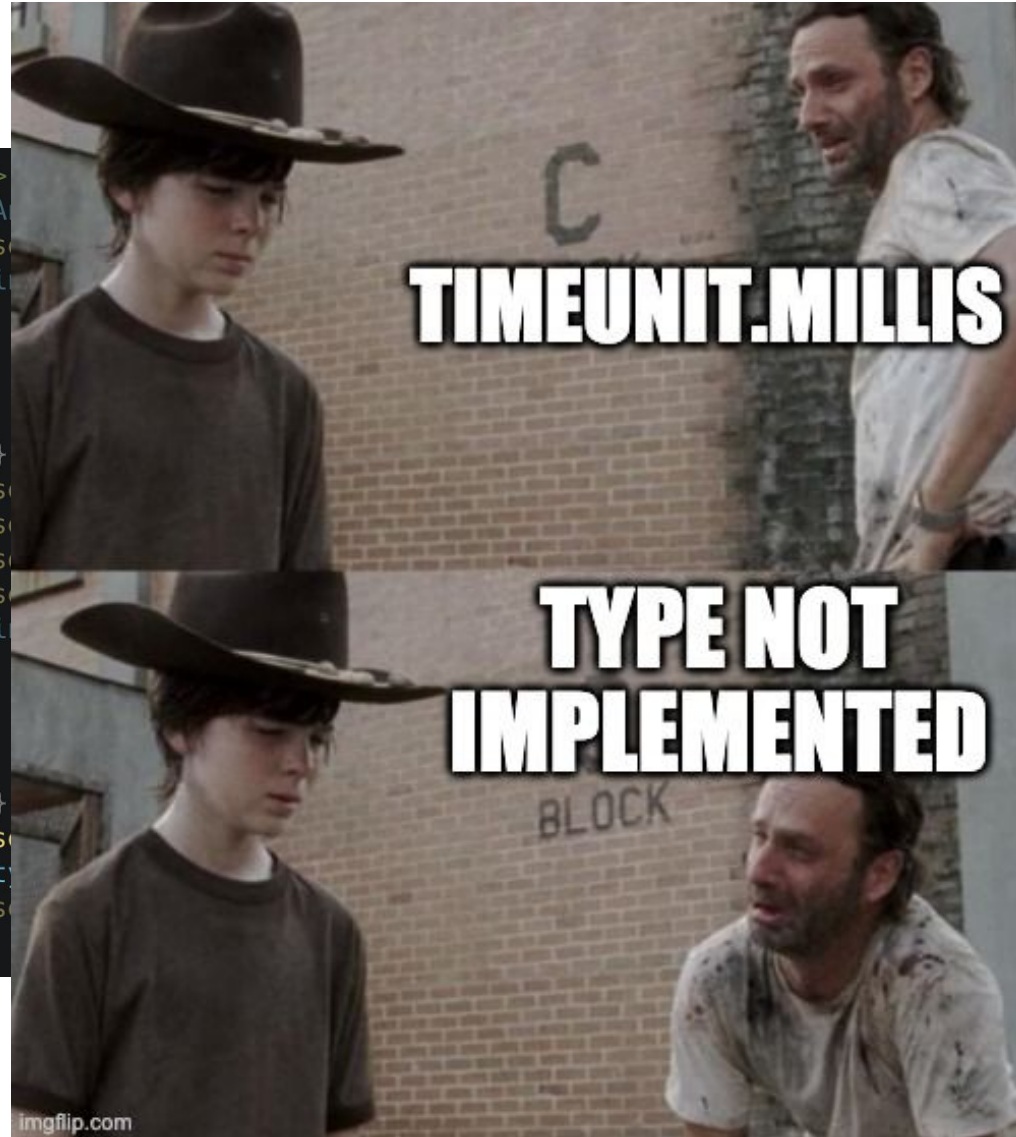
[Spark ParquetSchemaConverter.scala](#)

```
1 case INT32 =>
2     typeAnnotation match {
3         case intTypeAnnotation: IntLogicalTypeAnnotation if intTypeAnnotation.isSigned =>
4             intTypeAnnotation.getBitWidth match {
5                 case 8 => ByteType
6                 case 16 => ShortType
7                 case 32 => IntegerType
8                 case _ => illegalType()
9             }
10        case null => IntegerType
11        case _: DateLogicalTypeAnnotation => DateType
12        case _: DecimalLogicalTypeAnnotation => makeDecimalType(Decimal.MAX_INT_DIGITS)
13        case intTypeAnnotation: IntLogicalTypeAnnotation if !intTypeAnnotation.isSigned =>
14            intTypeAnnotation.getBitWidth match {
15                case 8 => ShortType
16                case 16 => IntegerType
17                case 32 => LongType
18                case _ => illegalType()
19            }
20        case t: TimestampLogicalTypeAnnotation if t.getUnit == TimeUnit.MILLIS =>
21            typeNotImplemented()
22        case _ => illegalType()
23    }
```


Data type issues

- time_millis

```
1 case INT32 =>  
2   typeA  
3   cas  
4   i  
5  
6  
7  
8  
9   }  
10  cas  
11  cas  
12  cas  
13  cas  
14  i  
15  
16  
17  
18  
19   }  
20  cas  
21  t  
22  cas  
23 }
```



```
tion.isSigned =>
```

```
(X_INT_DIGITS)  
tion.isSigned =>
```

```
LLIS =>
```


Data type issues

```
63
64 SELECT
65     min( ) AS min_date1,
66     max( ) AS max_date1
67 FROM
68
```

- time_millis
- inadequate dates

Data type issues

```
63  
64 SELECT  
65     min( [REDACTED] ) AS min_date1,  
66     max( [REDACTED] ) AS max_date1  
67 FROM [REDACTED]  
68
```

- time_millis
- inadequate dates



MIN_DATE1
0006-05-12

Data type issues

```
63  
64 SELECT  
65     min( [REDACTED] ) AS min_date1,  
66     max( [REDACTED] ) AS max_date1  
67 FROM [REDACTED]  
68
```

- time_millis
- inadequate dates



 MIN_DATE1 
0006-05-12



Data type issues

```
63  
64 SELECT  
65     min( ) AS min_date1,  
66     max( ) AS max_date1  
67 FROM  
68
```

- time_millis
- inadequate dates



MIN_DATE1	MAX_DATE1
0006-05-12	2580-05-15



Data type issues

```
63  
64 SELECT  
65     min( ) AS min_date1,  
66     max( ) AS max_date1  
67 FROM  
68
```

- time_millis
- inadequate dates



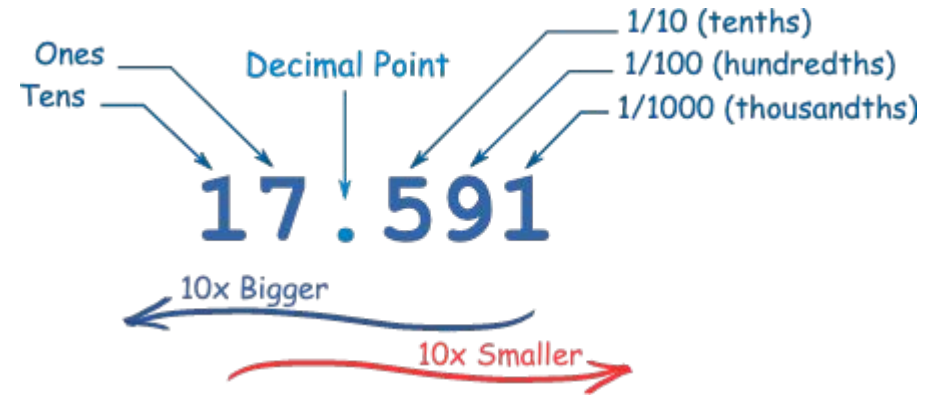
MIN_DATE1	MAX_DATE1
0006-05-12	2580-05-15



Data type issues



- `time_millis`
- inadequate dates
- decimal



Data type issues: solution



Data type issues: solution

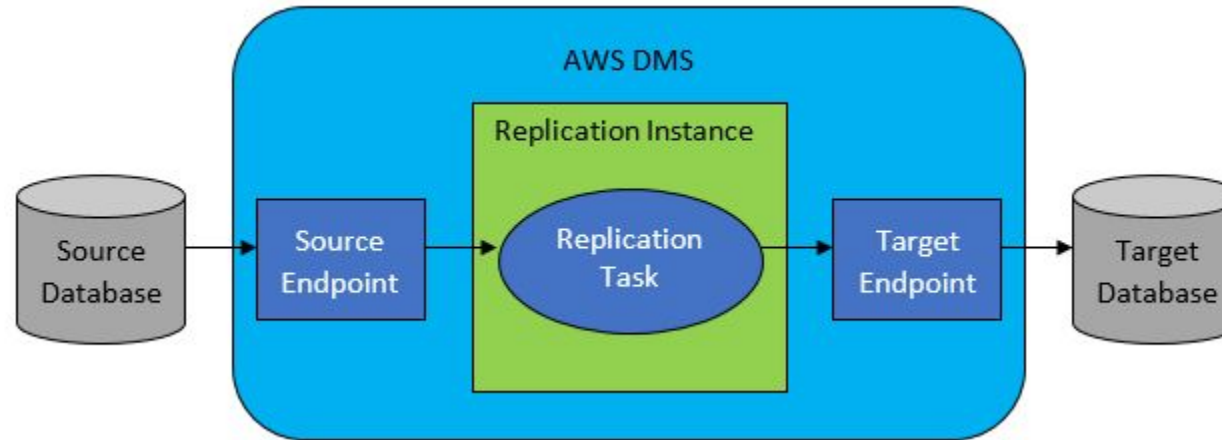


DMS: Source Endpoint



```
1 DataSourceEndpoint:
2   Type: AWS::DMS::Endpoint
3   Properties:
4     EndpointIdentifier: prod-data-source-endpoint
5     EngineName: data_source
6     EndpointType: source
7     Username: '{{resolve:secretsmanager:service:SecretString:username}}'
8     Password: '{{resolve:secretsmanager:service:SecretString:password}}'
9     ServerName: '{{resolve:secretsmanager:service:SecretString:host}}'
10    Port: '{{resolve:secretsmanager:service:SecretString:port}}'
11    DatabaseName: '{{resolve:secretsmanager:service:SecretString:database}}'
12    ExtraConnectionAttributes: 'executeTimeout=3600;'
```


DMS: ReplicationTask fun part



DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

DataSource schema: my_schm

DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

```
In [11]: len('my_schm')  
Out[11]: 7
```

DataSource schema: my_schm

DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

```
In [11]: len('my_schm')  
Out[11]: 7
```

DataSource schema: my_schm

known issue when using DB2 as source with schema name having (underscore) and less than 8 characters selection rule fails

DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

```
In [11]: len('my_schm')  
Out[11]: 7
```

DataSource schema: my_schm

known issue when using DB2 as source with schema name having () underscore and less than 8 characters selection rule fails



DMS: ReplicationTask fun part

No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

```
In [11]: len('my_schm')  
Out[11]: 7
```

DataSource schema: my_schm

known issue when using DB2 as source with schema name having () underscore and less than 8 characters selection rule fails

We are in process of updating the documentation with this limitation.

Being said that there is a workaround to select one specific table by changing "rule-action" from "include" to "explicit"



DMS: ReplicationTask fun part



No tables were found at task initialization. Either the selected table(s) no longer exist or no match was found for the table selection pattern

```
In [11]: len('my_schm')
Out[11]: 7
```

DataSource schema: my_schm

known issue when using DB2 as source with schema name having () underscore and less than 8 characters selection rule fails

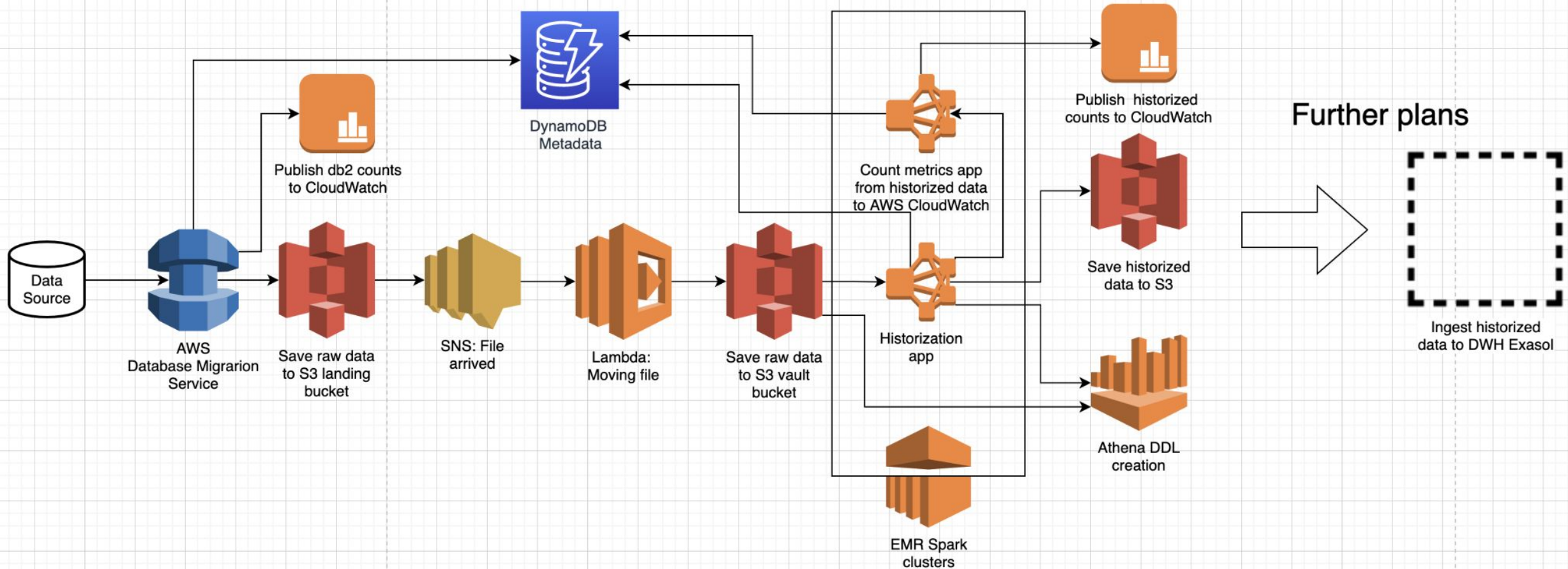


We are in process of updating the documentation with this limitation.

Being said that there is a workaround to select one specific table by changing "rule-action" from "include" to "explicit"

rule-action	include, exclude, explicit	A value that includes or excludes the object or objects selected by the rule. If explicit is specified, you can select and include only one object that corresponds to an explicitly specified table and schema.
-------------	----------------------------	--

Further work



Achievements





- We have built ETL process for full and incremental load data from Data Source to Data Lake

Achievements



- We have built ETL process for full and incremental load data from Data Source to Data Lake
- We have automated and auto scaled process

Achievements



- We have built ETL process for full and incremental load data from Data Source to Data Lake
- We have automated and auto scaled process
- About 70 tables are running day by day



- We have built ETL process for full and incremental load data from Data Source to Data Lake
- We have automated and auto scaled process
- About 70 tables are running day by day
- We have just one entry point for adding new tables to the process - all other steps are catching up and reusing this metadata

Achievements



- We have built ETL process for full and incremental load data from Data Source to Data Lake
- We have automated and auto scaled process
- About 70 tables are running day by day
- We have just one entry point for adding new tables to the process - all other steps are catching up and reusing this metadata
- We enabled process monitoring and alarming on data and process issues

Achievements



- We have built ETL process for full and incremental load data from Data Source to Data Lake
- We have automated and auto scaled process
- About 70 tables are running day by day
- We have just one entry point for adding new tables to the process - all other steps are catching up and reusing this metadata
- We enabled process monitoring and alarming on data and process issues
- Our data stored in effective way and easy to access with different tools

Learning points to you



Q&A: Your voice matters



[Link to the anonymous survey](#)

Metadata in DynamoDB



Partition key:

TableID is the following composition **<Source>-<Target>-<SchemaName>-<TableName>**

Metadata in DynamoDB



Partition key:

TableID is the following composition **<Source>-<Target>-<SchemaName>-<TableName>**

Sort key: **TableName**

