

Computomic

On-premises to Databricks Data Migration Architecture



Abstract

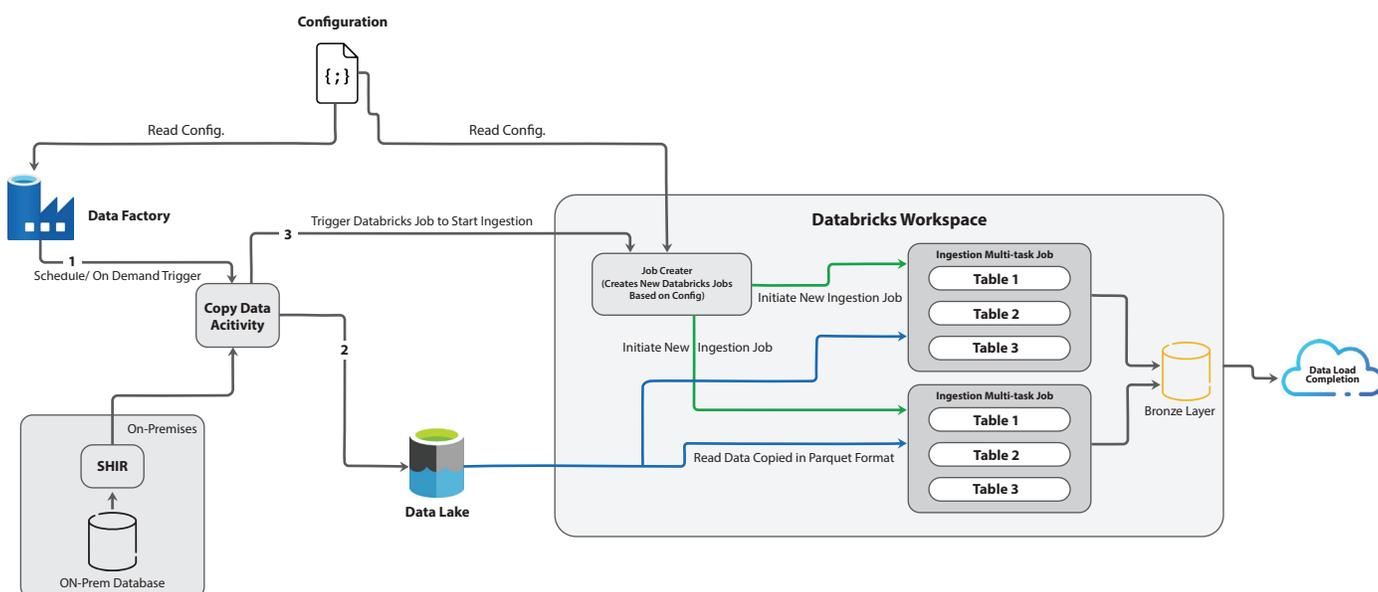
This white paper outlines a structured, phased approach for migrating complex on premises data systems to the Databricks Lakehouse, using intelligent orchestration and Azure Data Factory to bridge environments. It highlights practical architecture and data movement strategies tailored to stringent data governance and security needs. The migration method emphasizes minimizing operational risk while efficiently transferring large datasets into Databricks' scalable lakehouse architecture. Detailed configuration and pipeline orchestration techniques demonstrate how to automate parallel data extraction and loading. The paper serves as a hands on guide for teams modernizing legacy data platforms with minimal disruption.

Overview

On-premises to cloud migrations often involve significant effort in moving historical data. Our client's Oracle to Databricks project required migrating historical data from on-premises Oracle databases to the Databricks platform. This involved over 200+ tables with varying data volumes. Due to data sensitivity, the client mandated against third-party tools and preferred a solution that kept data within their network.

To address these constraints, we developed the following approach

Architecture



Leveraging the client's existing Azure infrastructure, we utilize Azure Data Factory (ADF) with a Self-Hosted Integration Runtime (SHIR) to bridge on-premises and cloud environments.

Pre-requisites

SHIR Deployment: Install Self hosted integration runtime on a Virtual Machine with connectivity to on-premises data sources.

Storage Account Setup: Configure an Azure Storage account accessible from both ADF and Databricks.

High Level Steps

Configuration Management:

- Create a configuration file with source and destination table details.
- Store configuration in a shared storage account which is accessible from ADF as well from databricks.

ADF Pipeline Orchestration:

- Trigger the ADF pipeline, passing the configuration file path as a parameter.
- The pipeline dynamically generates and executes multiple parallel copy data activities.
- These activities extract data from on-premises sources and store them in Parquet format within the shared storage account.
- Parallelism is controlled to minimize load on both source systems and SHIR.

Databricks Job Initiation:

- Upon successful data extraction, the ADF pipeline invokes the Databricks API to trigger a "job_creator" job.

Databricks Job Creation:

- The "job_creator" job reads the configuration file and dynamically creates a new Databricks multi-task job.
- This multi-task job parallelizes the loading of exported data from the shared storage account into the Databricks bronze layer.

Sample Configuration File

(can be extended to include multiple other settings)

```
[
  {
    "source_database": "edw",
    "source_table": "customer",
    "source_filter_condition": "",
    "target_catalog": "uson",
    "target_schema": "bz_edw",
    "target_table": "customer",
    "isActive": "true",
    "group_id" : "edw_group_1"
  },
  {
    "source_database": "rpt",
    "source_table": "orders",
    "source_filter_condition": " WHERE TRUNC(created_date) = TRUNC(SYSDATE- 1);",
    "target_catalog": "uson",
    "target_schema": "bz_rpt",
    "target_table": "orders",
    "isActive": "true",
    "group_id" : "edw_group_2"
  }
]
```

ADF Pipeline

The screenshot shows the Azure Data Factory pipeline editor for a pipeline named 'pl_export_on_pre...'. The left sidebar shows 'Factory Resources' with a tree view containing Pipelines, Datasets, Data flows, Power Query, and Templates. An arrow points to the 'ds_sql_database' dataset with the label 'Reuse datasets'. The main editor shows a pipeline flow: Lookup (read_config) -> Filter (Select Active) -> ForEach (Load data) -> Execute Pipeline (Start Databricks... pl_initiate_databricks...). The 'Output' tab is active, showing a table of pipeline run results.

Activity name	Activity status	Run start
read_config	Succeeded	12/27/2024, 10:58:02 AM
Select Active	Succeeded	12/27/2024, 10:58:20 AM
Load data	Succeeded	12/27/2024, 10:58:21 AM
Copy data to storage	Succeeded	12/27/2024, 10:58:22 AM
Copy data to storage	Succeeded	12/27/2024, 10:58:22 AM
Start Databricks Job	Succeeded	12/27/2024, 10:59:05 AM

An arrow points to the two 'Copy data to storage' rows with the label 'Run Parallel'.

Output In Storage Account

The screenshot shows the 'landingzone' container in Azure Storage. The 'Overview' tab is selected, showing the location path: 'landingzone / raw_data / uson / bz_edw / customers'. An arrow points to this path with the label 'folder structure based on config file'. Below the path, there is a table of blobs.

Name	Modified	Access tier	Archive status	Blob type
[.]				
dbo_00000.customers.parquet	12/27/2024, 10:58:52 AM	Hot (Inferred)		Block blob

The data reading process can be optimized based on the specific source system. For instance, techniques like physical table partitions or dynamic range partition properties within the copy data activity can be leveraged to enhance performance.

For detailed guidance on executing Databricks jobs from Azure Data Factory, refer to this comprehensive blog post: <https://techcommunity.microsoft.com/blog/analticsonazure/leverage-azure-databricks-jobs-orchestration-from-azure-data-factory/3123862>

This blog post offers a step-by-step approach to constructing a modular ADF pipeline that can execute any Databricks job using built-in ADF activities and managed identity authentication. It also covers integrating the ADF Managed Identity as a contributor within your Databricks workspace.

Upload to Databricks

Once the data is available in the storage account as Parquet files, Databricks is utilized to load it. As depicted in the architecture diagram, two workflows are involved: `job_creator` and `ingestion_multitask_job`.

job_creator: This job, triggered by the ADF pipeline via the JOBS API with a JSON configuration file path as input, dynamically creates a new `ingestion_multitask_job`. This process leverages the Databricks SDK to generate and execute notebook tasks, enabling parallel data loading for multiple tables.

Job Creator Code:

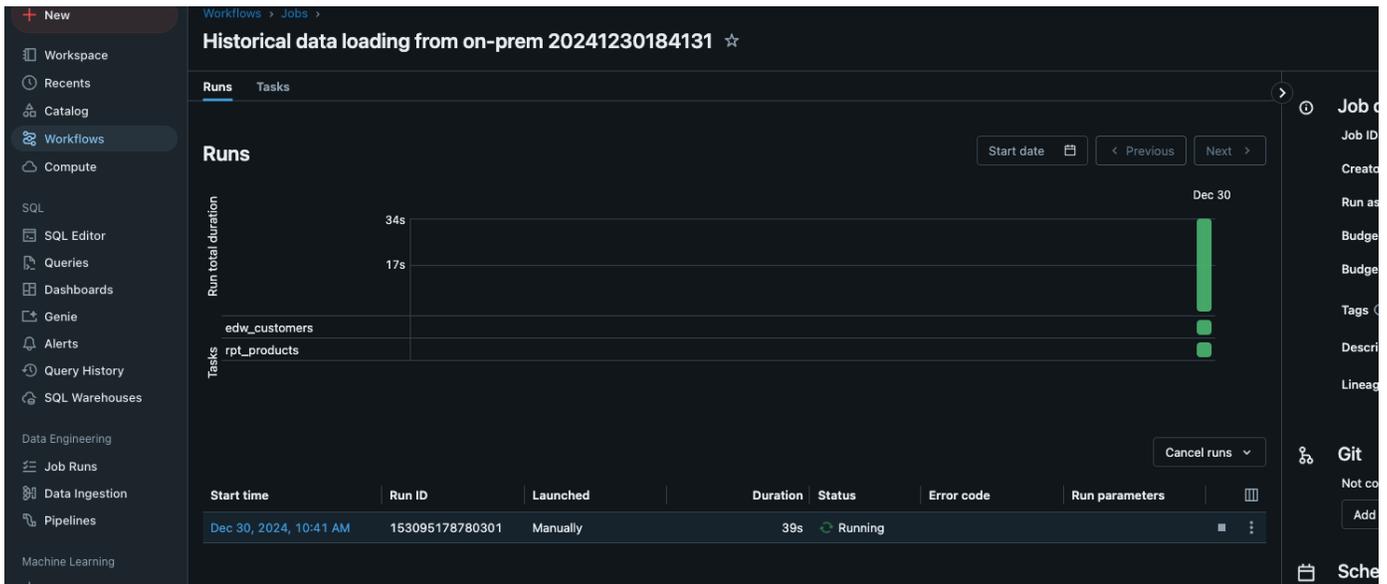
```

10:41 AM (2s) 4
1 from databricks.sdk import WorkspaceClient
2 from databricks.sdk.service.jobs import Task, NotebookTask
3 from datetime import datetime, timezone
4
5 import json
6 from datetime import datetime, timezone
7
8 ### Read json configuration file and save it to json_string
9 ### Load settings in settings
10 settings = json.loads(json_string)
11
12 # Loop through settings and create notebook tasks
13 tasks = []
14 for row in settings:
15     task = Task(
16         task_key=f'{row["source_database"]}_{row["source_table"]}',
17         notebook_task=NotebookTask(
18             notebook_path="/Workspace/testsandesh_data_load_on_prem/historical_data_loader",
19             base_parameters={
20                 "source_database": row["source_database"],
21                 "source_table": row["source_table"],
22                 "target_catalog": row["target_catalog"],
23                 "target_schema": row["target_schema"],
24                 "target_table": row["target_table"],
25             }
26         )
27     )
28     tasks.append(task)
29
30 ### Adding a UTC timestamp to make sure job name is unique.
31 current_time = datetime.now(timezone.utc)
32
33 job_settings = {
34     "name": f"Historical data loading from on-prem {current_time.strftime('%Y%m%d%H%M%S')}",
35     "tasks": tasks
36 }
37 workspace = WorkspaceClient()
38 job = workspace.jobs.create(**job_settings)
39 logger.info(f"Job successfully created: {job}")
40 # Start the job
41 workspace.jobs.run_now(job_id=job.job_id)
42 logger.info(f"Job successfully started: {job}")

INFO:hist.job_creator:Job successfully created: CreateResponse(job_id=1102262524521763)
INFO:hist.job_creator:Job successfully started: CreateResponse(job_id=1102262524521763)

```

Once ADF executes this job it will create a multi task job as below.



The provided JSON configuration specifies two tables: customers and products. Consequently, the **DataLoader** notebook job includes two tasks, each responsible for loading data into the corresponding table.

The following code snippet within the **DataLoader** notebook demonstrates the basic logic for reading Parquet files and loading them into the respective destination tables.

```

Historical data loader

▶ 10:20 AM (<1s) 2
1 import logging
2 logger = logging.getLogger("hist.dataloader")
3 logger.setLevel(logging.INFO)
4 logger.info("Starting the dataloader")
5
INFO:hist.dataloader:Starting the dataloader

▶ 10:14 AM (<1s) 3
1 ##### Read parameters
2 source_database = dbutils.widgets.get("source_database")
3 source_table = dbutils.widgets.get("source_table")
4 target_catalog = dbutils.widgets.get("target_catalog")
5 target_schema = dbutils.widgets.get("target_schema")
6 target_table = dbutils.widgets.get("target_table")
7
8 ##### Construct input source path and destination table
9 destination_table_name = f"{target_catalog}.{target_schema}.{target_table}"
10 source_path = f"/Volumes/landingzone/{source_database}/{source_table}"
11 logger.info(f"Loading data to {source_path} =====> {destination_table_name}")
12
13 ##### Load the data in the target table
14 try:
15     df = spark.read.parquet(source_path)
16     df.write.mode("overwrite").saveAsTable(destination_table_name)
17     logger.info(f"Data loading is successful for {destination_table_name}")
18 except Exception as e:
19     logger.exception(e)
20     raise e
    
```

Computomic

Results

This project successfully demonstrated a robust and secure method for migrating historical data from on-premises Oracle databases to the Databricks platform on Azure. By leveraging Azure Data Factory with a Self-Hosted Integration Runtime, we achieved a controlled and efficient data extraction process while adhering to the client's requirement for on-premises data access and security.

The implementation of a dynamic configuration-driven approach, combined with the parallel execution capabilities of both ADF and Databricks, significantly accelerated the migration process. This solution provides a scalable and flexible framework for future data migrations and ongoing data integration between on-premises and cloud environments

Author

- **Sandesh Daddi** - Resident Solutions Architect

www.computomic.ai