# DBImport

The smart way to ingest data
into a Data Lakehouse

**Middlecon**

# Evolution of data loading for analytics



**LATE 1980'S**

## Data Warehouse

BI • Reports

Data Marts

ETL

External Data • Operational Data

**2011**

## Data Lake

Data Science • Machine Learning • Real-Time Database • Reports • BI

Data Prep and Validation

ETL

Data Marts

Data Lake

Structured, Semi-Structured and Unstructured Data

**2020**

## Lakehouse

Streaming Analytics • Data Science • Machine Learning

BI

Structured, Semi-Structured and Unstructured Data

Middlecon

# Evolution of data loading for analytics



**LATE 1980'S**

## Data Warehouse

BI   Reports

Data Marts

ETL

External Data   Operational Data

**2011**

## Data Lake

Data Science   Machine Learning   Real-Time Database   Reports   BI

Data Prep and Validation

ETL

Data Marts

Data Lake

Structured, Semi-Structured and Unstructured Data

**2020**

## Lakehouse

Streaming Analytics   Data Science   Machine Learning

BI

Structured, Semi-Structured and Unstructured Data
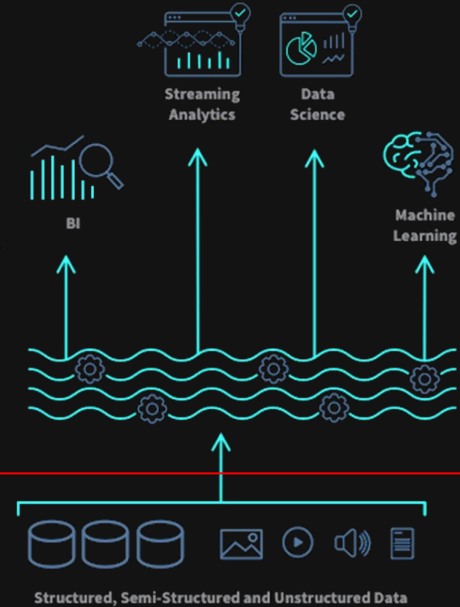
Middlecon

# Regardless of data evolution, we will always have to move and load data into a system

**THE CHALLENGE**

- Large tables take forever to ingest
- Handling incremental data synchronization
- Source system owners says no to CDC tools
- Schema Changes in the source system

Middlecon

- Support for all common database types out of the box

- Auto discovery of tables and views from source systems

- Full and incremental import and export functionality

- Automatic handling of schema changes

- Automatic handling of table and column descriptions

- Automatic handling of primary and foreign keys

**Middlecon**

# Data Ingestion Wish-List

- Change Data Tracking functions

- Audit of changed metadata

- Parallel execution and central scheduling

- Full logging and statistics of data ingestion processes

- Modern fileformat usable by many tools
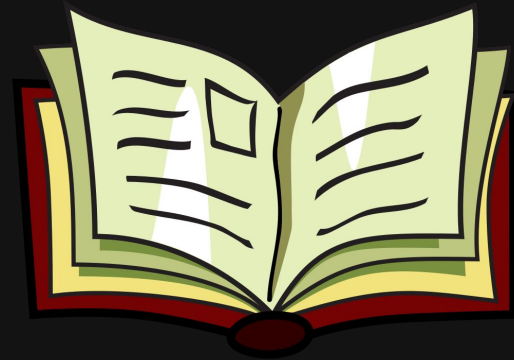
- Ease of use

# DBImport

Key functionalities

# DBImport Tool - History

Many years of experience with different sqoop, spark and jdbc data ingestions problems resulted in an OpenSource project called DBImport

Largest installation are ingesting data from over 19.000 tables in a single Hadoop environment per day, from 250 different sources

Licensed by Apache 2.0

**Middlecon**

# DBImport Tool - Goal

The goal of DBImport is to have a fast, simple but powerful tool to create a source aligned copy of the data and at the same time keeping track of what data have changed between two different points in time.

The target is always a Data Lake or Data Lakehouse and the source aligned data storage is usable for all AI, Machine learning, BI and exploration of data within the organization.

**Middle**con

Support for Oracle, MsSQL, MySQL, DB2 UDB,
DB2 AS400, MongoDB, Snowflake, Progress,
SQL Anywhere and PostgreSQL databases

Both Full extraction and Incremental extraction is supported through standard SQL.

Additionally, both MSSQL Change Data Tracking and Oracle Flashback Query can be used to extract data from source system

Uses sqoop or spark in the background for transferring data from source system

# DBImport - Fileformats and Storage

DBImport will save data in Orc, Parquer or
Iceberg format on either HDFS or Ozone.

# DBImport - Load and transform tool

Once data is loaded from source system, transformation is handled by Spark or Hive

Data accessible as files on HDFS or Ozone or through Hive and Impala with standard SQL

# DBImport Tool - Supported export systems

Exports data stored in Hive.

Support for Oracle, MsSQL, MySQL, DB2 UDB, DB2 AS400, MongoDB, Snowflake, Progress, SQL Anywhere and PostgreSQL databases. Also support for creating files on AWS S3

# DBImport - History Data

If selected, DBImport will keep track of what data has been changed in a source system table and log this in a separate history table.

This history table can then be used to incrementally process data even if the source system does not support incremental data loads or keep track of changed data.

It can also be used to see how data is changed over time.

| Id | Type | Stock | datalake_iud | datalake_timestamp |
|----|--------|-------|--------------|---------------------|
| 1 | apple | 205 | I | 2019-09-08 02:18:45 |
| 2 | orange | 155 | I | 2019-09-08 02:18:45 |
| 3 | banana | 40 | I | 2019-09-08 02:18:45 |
| 4 | pear | 70 | I | 2019-09-08 02:18:45 |
| 5 | kiwi | 65 | I | 2019-09-08 02:18:45 |
| 1 | apple | 200 | U | 2019-09-09 02:20:40 |
| 2 | orange | 125 | U | 2019-09-09 02:20:40 |
| 5 | kiwi | 65 | D | 2019-09-09 02:20:40 |

Middlecon

# DBImport - Auto Discovery and source system changes

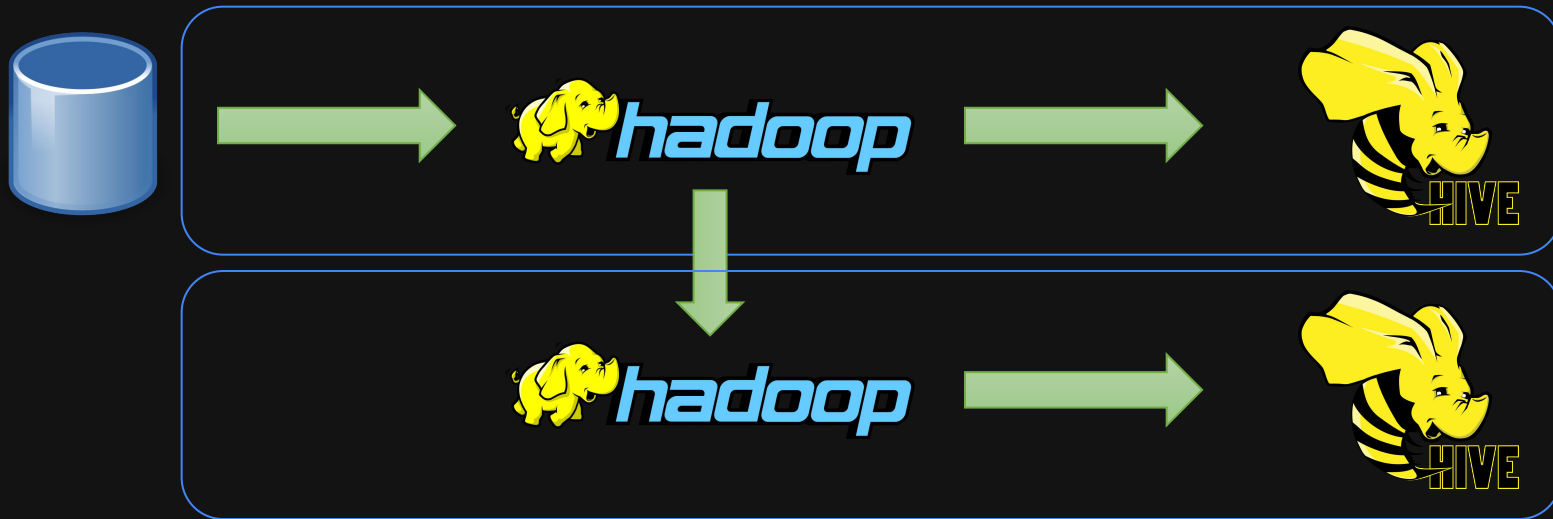Auto discovery of tables and views from all supported source systems.

All changes on column and table changes including keys and comments is identified and changed in target tables. All changes is logged and accessible by third-party tools.
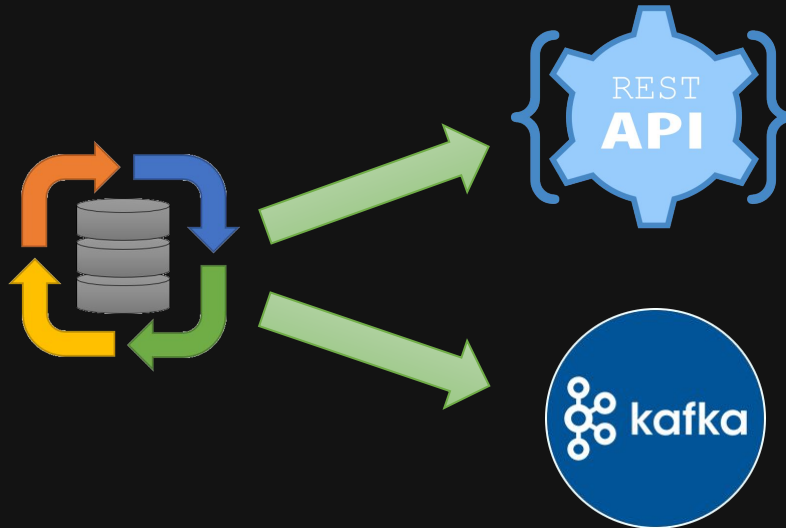
Middlecon

# DBImport Tool - Multi Cluster Ingestions

Ability to keep two system updated with only one extraction from source system

# DBImport Tool - Notification

After each task is completed, a notification will be posted through a Rest API or on a Kafka topic.
JSON Data includes statistics and status of the job

# DBImport Tool - Statistics

All completed tasks will store statistics about how the task was executed. Both time and the amount of data that was moved.

During the import, there is verifications three times on the imported data to make sure that the data that was fetched is the same as on the source system.

If the data is not correct, or due to other technical reason the job failed the task will be restarted from last successful state

Middlecon

# DBImport Tool - Data Anonymization

During transport of data from source system, data can be anonymized before it is stored on any permanent storage device.

Supports Hash, 'Replace with star' and 'Show first 4 chars only' methods. Selectable on column level

Middlecon

# DBImport Tool - Airflow Integration

Airflow is an OpenSource scheduling and Workflow management system

# DBImport Tool - Airflow Integration

DBImport generates Airflow DAG's and writes them to the DAG directory of Airflow.

It also generates pools in Airflows configuration database, one per source database hostname and one per DAG
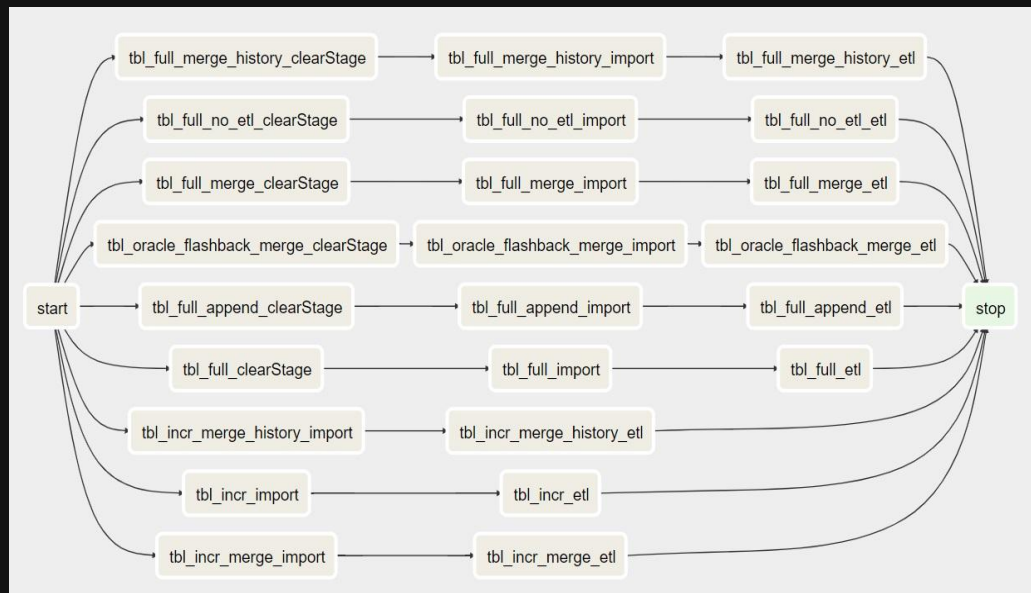


MiddleCon

# DBImport Tool - Airflow Integration

Each DAG contains the required steps for all tables being imported in the DAG. This usually is all the tables from the source database that is being imported.
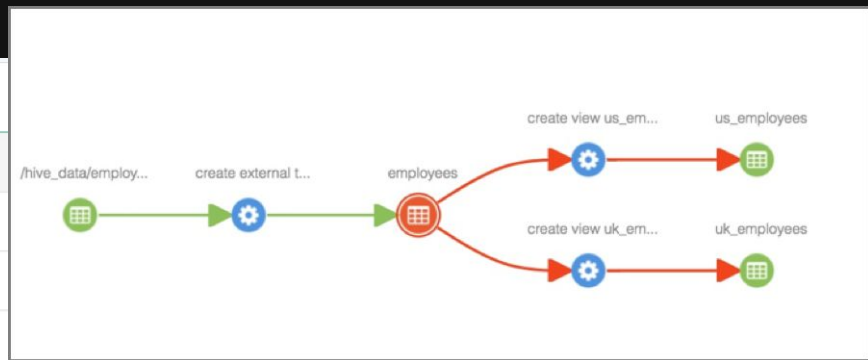


Middlecon

# DBImport Tool - Atlas

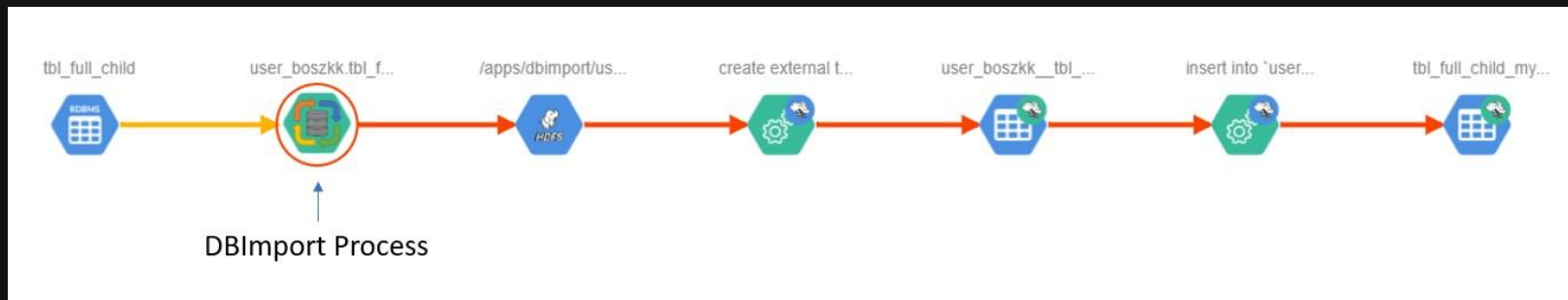Apache Atlas is an open-source metadata
and big data governance framework



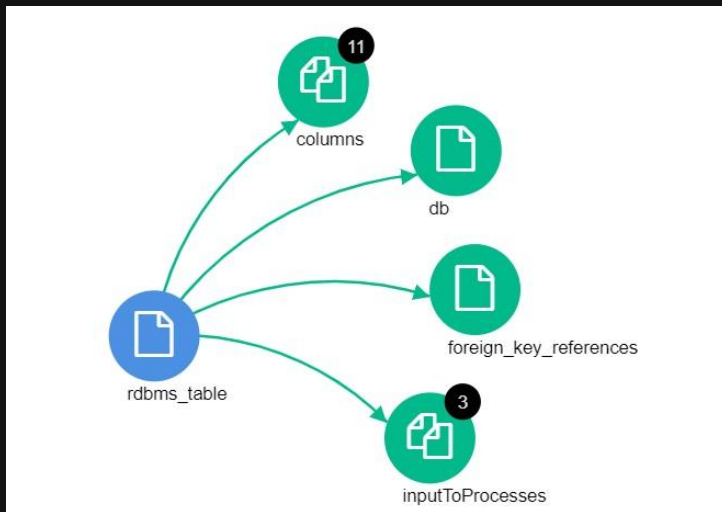| | Name | Owner | Type | Type | Classifications | |
|---|---|---|---|---|---|---|
| ☐ | providername | hive | hive_column | string | VENDOR_PII | + |
| ☐ | emailaddress | hive | hive_column | string | PII ✖ | + |
| ☐ | ccnumber | hive | hive_column | string | PII ✖ | + |
| ☐ | nationalid | hive | hive_column | string | PII ✖ | + |
| ☐ | nationalid | hive | hive_column | string | PII ✖ | + us_customers |
| ☐ | ssn | hive | hive_column | string | FINANCE... ✖ | + tax_2015 |
| ☐ | providername | hive | hive_column | string | VENDOR... ✖ | + provider_summary |
| ☐ | ssn | hive | hive_column | string | FINANCE... ✖ | + tax_2010 |

# DBImport Tool - Atlas Integration

Full lineage is available with the help of the custom DBImport Process Atlas type.

# DBImport Tool - Atlas Integration

DBImport utilizes the rdbms_* types already existing in Atlas to create a complete source system object. Source system object is complete with all information and relationship for the imported / exported table

# DBImport Tool - Atlas Integration

Source system information can be added in two steps.

- When an import or export is executed, the information for that specific table is created / updated in Atlas.

- When the DBImport server is discovering ALL tables on a specific connection, regardless if it's imported or not to Hadoop

**Middlecon**

# DBImport

Customer Use-Case

Large Swedish company in manufacturing industry uses DBImport for close to all batch ingestion into their Datalake / Data lakehouse.

- 19.000 tables per day
- 800 GB per day
- 41 billion rows per day
- 450 Airflow jobs

Total imported size over time is 1.1 PB with DBImport

Middlecon

# DBImport

Current and Future development
and improvements

- Python based Windows Client
  *In Development*

- Spark as ETL engine
  *In Development - Beta version available*

- Container based cloud application
  *On Road-Map*

**Middle**con

# Questions?

MiddleCON

# DBImport

# Like to contribute?
Please contact Berry Österlund

**Middlecon**

# DBImport

Git

https://github.com/Middlecon/DBImport

Documentation

https://dbimport.readthedocs.io/en/latest/

Lead developer

berry.osterlund@middlecon.se