



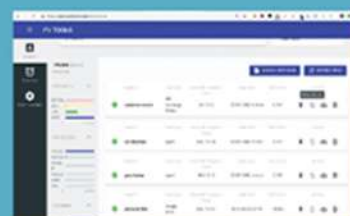
# Data Protection platform integration with Hadoop Hive



Hadoop Hive returns the requested data for the scan

Connector app creates an ODBC connection to connect to Hive database

## Data Protection Platform



Connector app retrieves the requested data for a specific timeline, processes the response data, sample it based on predefined sampling techniques, parse the response and sends it to the Customer Platform



Connector App

Connector triggers a call to the connector app to initiate a scan on the Hive database



## Customer

Customer is a leading Personal Data Privacy and Protection provider.

It enables organizations to discover and map all types of data from all enterprise data sources; automatically classify, correlate, and catalogue identity & entity data into profiles; manage and protect enterprise data with advanced data intelligence; and automate data privacy and protection.

It identifies all PII across structured, unstructured, cloud & Big Data.

Customer requested to build a Connector app to integrate their platform with Hadoop Hive to scan the data present in the Hadoop Hive for finding the PII information.



## Requirement



## Technology Solution

- Apache Hive is a data warehouse software project built on top of Apache Hadoop for providing data query and analysis. Hive gives a SQL-like interface to query data stored in various databases and file systems that integrate with Hadoop.
- Sacumen developed the Connector app to integrate Hadoop Hive using C# 8.0 (.NET Core 3.0). The Connector app performs the following actions:
  - Creates an ODBC driver connection.
  - Following are the important parameters to get-set Connection String. Rest of the parameters can be set as required by ones application.
    - DRIVER={Microsoft Hive ODBC Driver} ▫ Host=server\_name ▫ DRIVER={Microsoft Hive ODBC Driver} ▫ Host=server\_name
  - Connects to Hive cluster using the connection string.
  - Retrieves the data of all the tables from the database.
  - Gets the requested data within a certain timeline.
  - Samples the fetched data using predefined sampling techniques.
  - Formats the received data in the required format and pass it to the customer.

