# How to Guide

## Data Preparation Integration (DP)

**Version: Release 1.1**
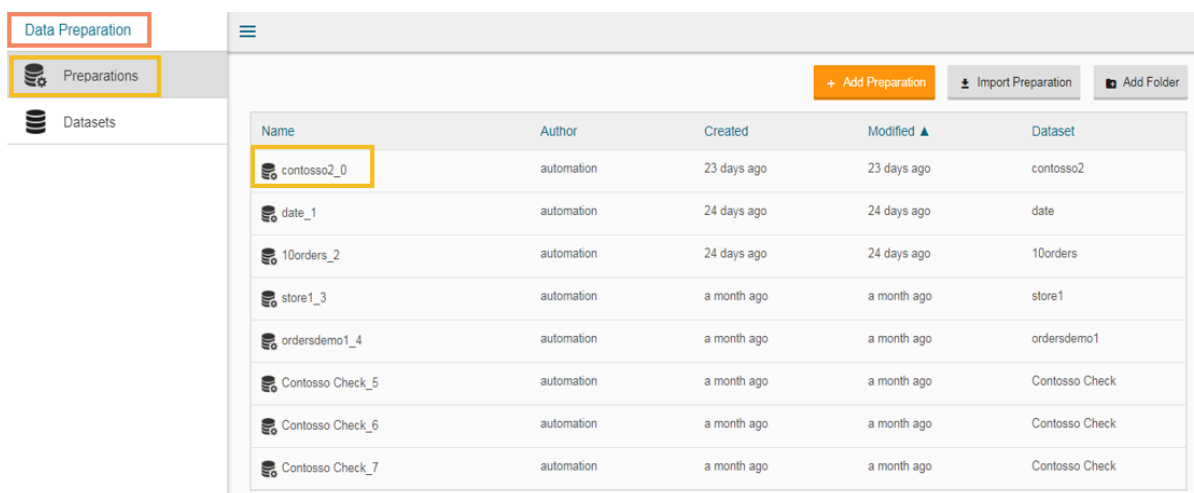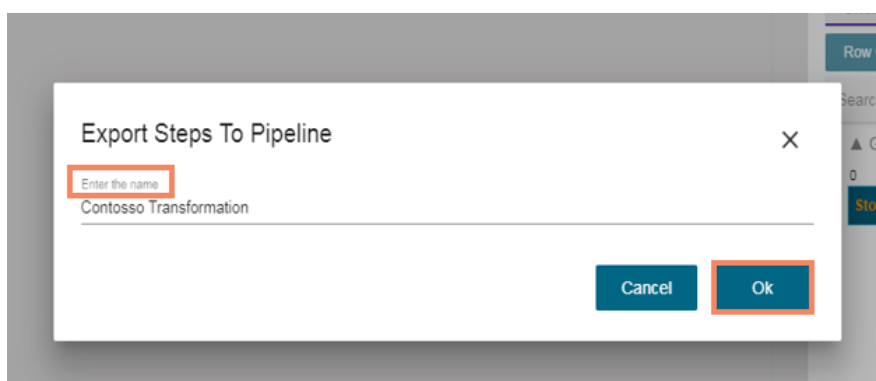
# Contents

# Data Preparation

Data Preparation is a BDB plugin that is used for cleansing data to prepare it for further analysis. The users can integrate the Data Preparation scripts with Data Pipeline to facilitate the data streaming.

## 1.1 Deploying a Data Preparation Model to Data Pipeline
1. Navigate to the Data Preparation landing page.
2. A list of all the preparation scripts appears under the '**Preparations**' option
3. Select a Preparation Model from the list



4. A new window '**Export Steps to Pipeline**' appears
   a. Enter a specific title for the Data Preparation model
   b. Click the '**Ok**' option to deploy the model to the Data Pipeline.



## 1.2 Using the Deployed Data Preparation Model in Data Pipeline
1. Navigate to the Pipeline Settings page
2. Select the Dataprep Scripts from the left panel
3. The '**Dataprep Scripts List**' appears with the deployed Data Preparation script.

4. The users can access the deployed Data Preparation script from the '**Script Name**' drop-down menu that appears inside the Dataprep Script Runner component of the Data Pipeline.