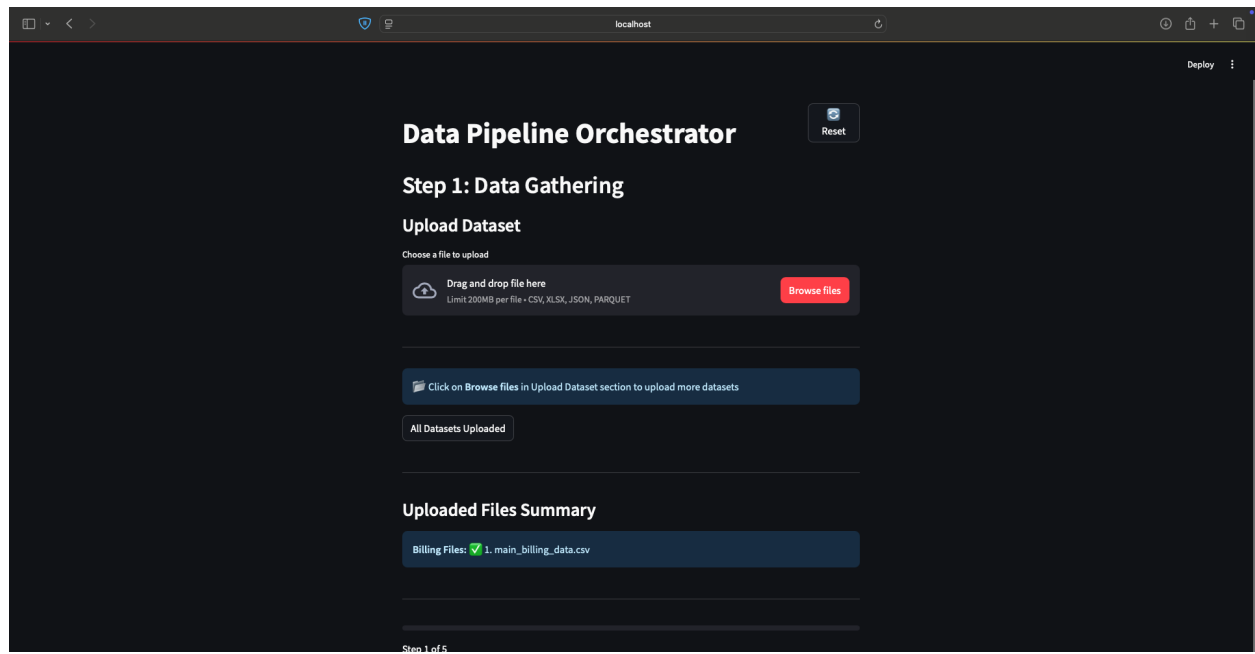# Data Prep Orchestrator Tutorial

## 1. Introduction

The Data Prep Orchestrator will help you unify your business data into a single dataset without code. Start by uploading your datasets one by one.

## 2. Confirming Dataset Type

Upon upload, the Orchestrator will categorize your data so it can (in the background) prepare to map your data to an internal schema. If you agree that your data fits in the category suggested by the AI, select yes. Otherwise, saying no will allow you to manually select the category the data belongs to. Repeat this step for each dataset you upload.
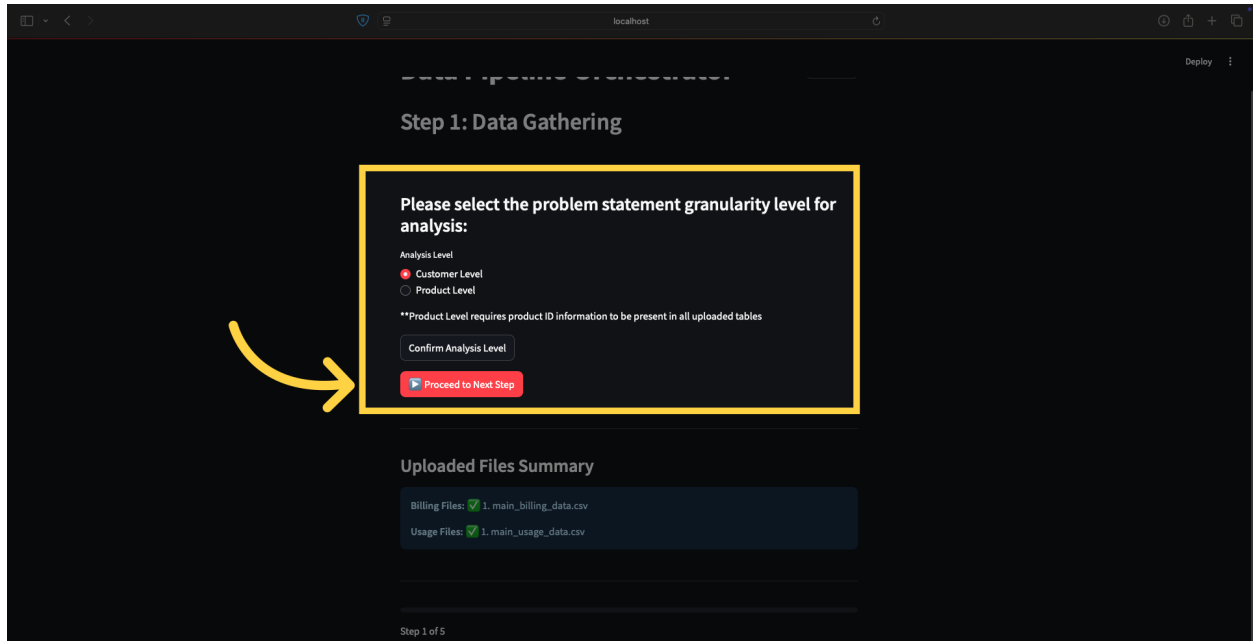
# 3. Unified Data Level

There are two levels at which your data can be transformed: Product and Customer level. Transforming at the Product level would lead to your data being unified at the ID/Date/Product level while the Customer level would lead to unification at the ID/Date level. Select the best level for your use case.
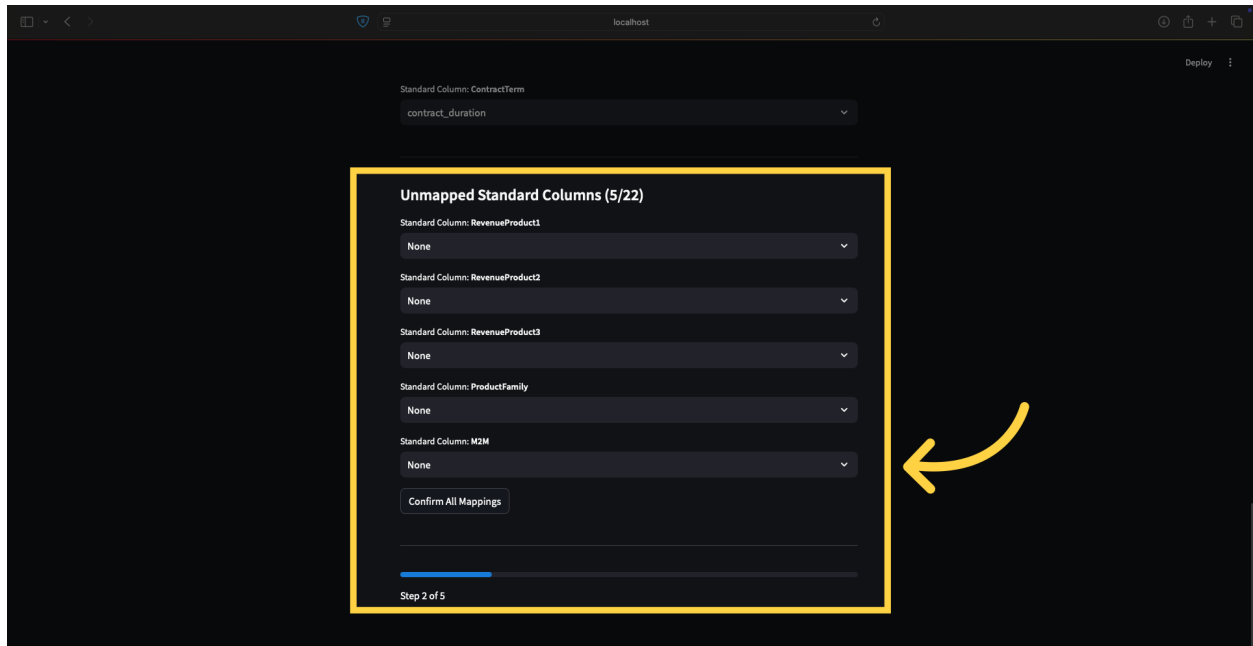
## 4. Mandatory Mappings

After selecting the target level, you will be asked to map your data to the agent's internal schema. This mapping allows the Orchestrator to leverage different strategies for cleaning + transforming your data in later steps. The Orchestrator will suggest some mappings for you, but if you want to change them, click on the respective dropdown. At minimum, your data must map to the mandatory fields (denoted by # next to the column name).

# 5. Optional Mappings

In addition to the suggested mappings, you can map other fields that the Orchestrator did not recognize. This step is optional, but any additional fields that are mapped will help the Orchestrator understand your data for later steps.
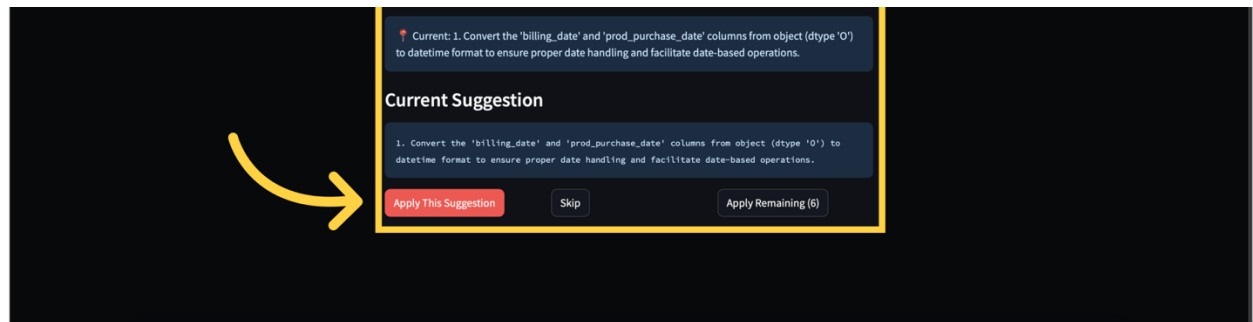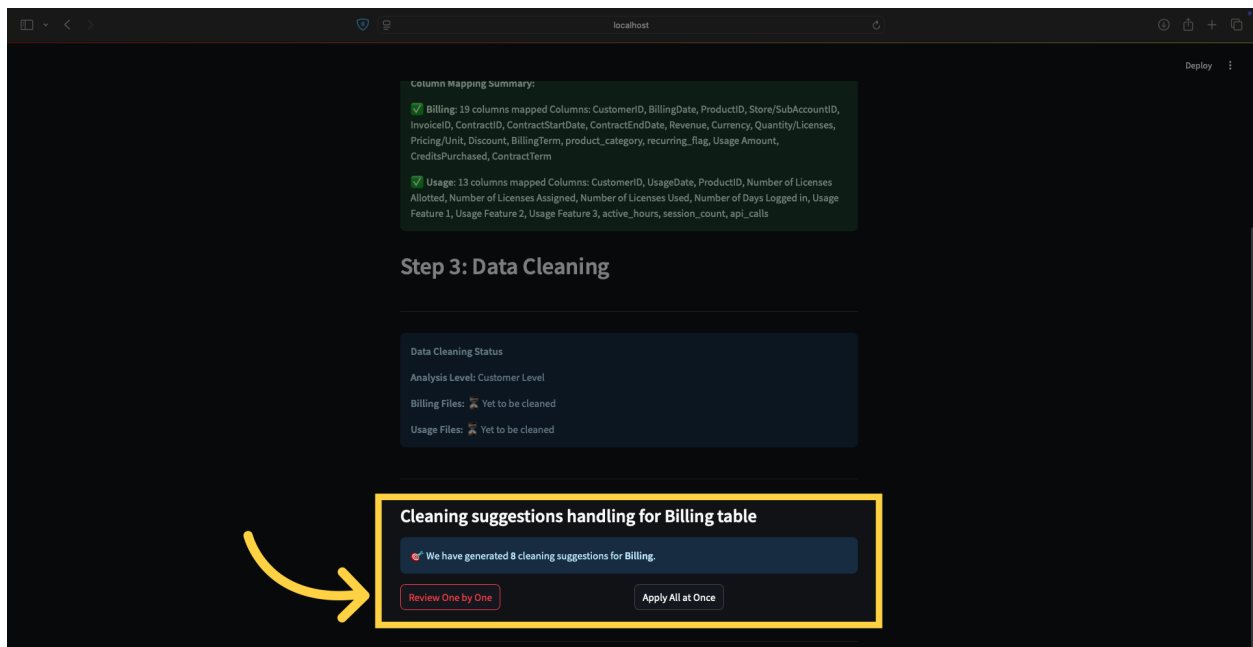
# 6. Data Cleaning

The Orchestrator will display cleaning suggestions and the reasons it wants to perform them on your data. If you feel that the suggestion is both necessary and beneficial to your data, you can apply it. If you don't feel that you need to apply a given suggestion, you can always skip it. You can always apply the Orchestrator's suggestions at once by clicking the Apply Remaining button.

# 7. Aggregations

The Aggregation agent will provide you with a set of recommended methods to aggregate your data. You can manually tune them by clicking on their respective text boxes. If you want to see why the agent chose what it did, click on the Show Aggregation Explanations button.

## 8. Aggregation Explanations

The Orchestrator will then show, column by column, the aggregations it selected and why. If you are satisfied with the reasons, confirm them, and the Orchestrator will then aggregate your data. Otherwise, you can go back and manually edit them via the checkboxes.
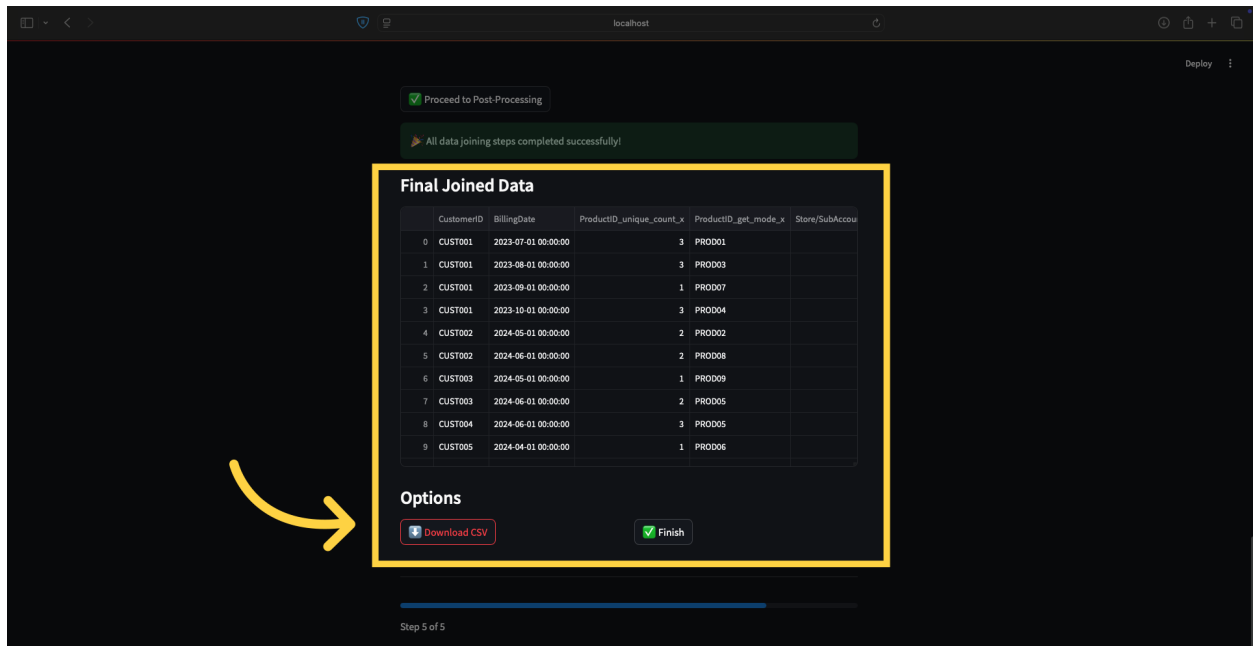
# 9. Joining your Data

The Orchestrator will determine which columns are best for joining your datasets together based on the level you selected in Step 3. Click proceed, and your data will be joined for you.

# 10. Final

Download your data, and you have finished your end-to-end data unification!