# Tombolo User Guide

HPCC Systems Solutions Lab

# Contents

# Introduction

Tombolo is a metadata tracking tool for HPCC Data Lake solution. It tracks the metadata around how every asset is used in a Data Lake, and the process flow as to how these assets evolve.

Tombolo helps you answer the following questions in a Data Lake environment.

"Who is the owner of xyz data?"

"What is the source of xyz data?"

"What does the data contain?"

"What are the compliance rules around xyz data?"

"Who approved the usage of this data?"

"When was this data last used?"

"Can you show me how this data is being used?"
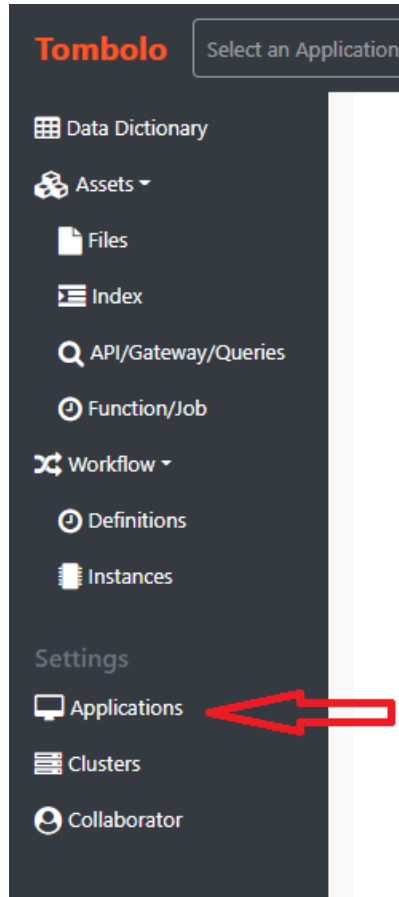
"Is this data being handled securely?"

"What is the impact of using this data?"
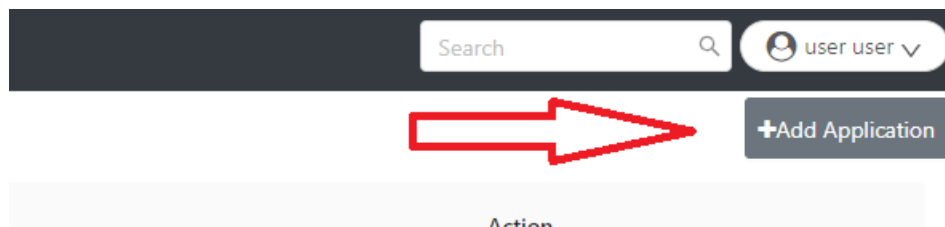
"What happens if this data does not arrive on time?"

"What happens if the data is not used on time?"
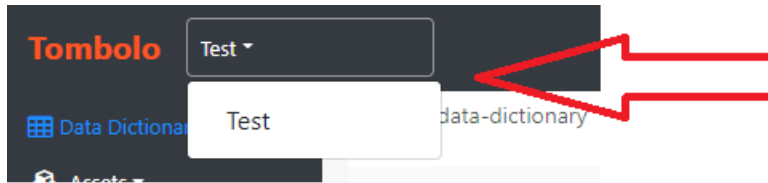
# Create an Application

In order to start using Tombolo, an "Application" has to be created. Application is a way of grouping your assets with in Tombolo. To create an application, click on the "Applications" link in the left nav. If you already have Applications, they will be listed in the Applications page



To create a new Application, click on Add Application button. Give the Application a meaningful name and description (optional)
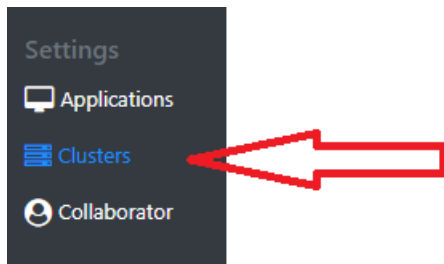


Click OK to create the Application. The Application should be now listed under the Applications dropdown.

## Add a Cluster

Tombolo gives you the ability to lookup your assets directly from an HPCC cluster. You can add Clusters through the Clusters options in the navigation.

PS: The system will allow you to add only the pre-configured clusters. If you need other clusters to be added, please let us know.
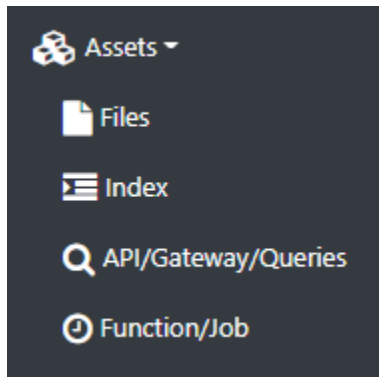


## Assets

Tombolo currently supports tracking metadata for the following Asset types:

- Files (Thor, CSV, JSON, XML)
- Index (HPCC)
- API/Queries (Roxie queries/other API's)
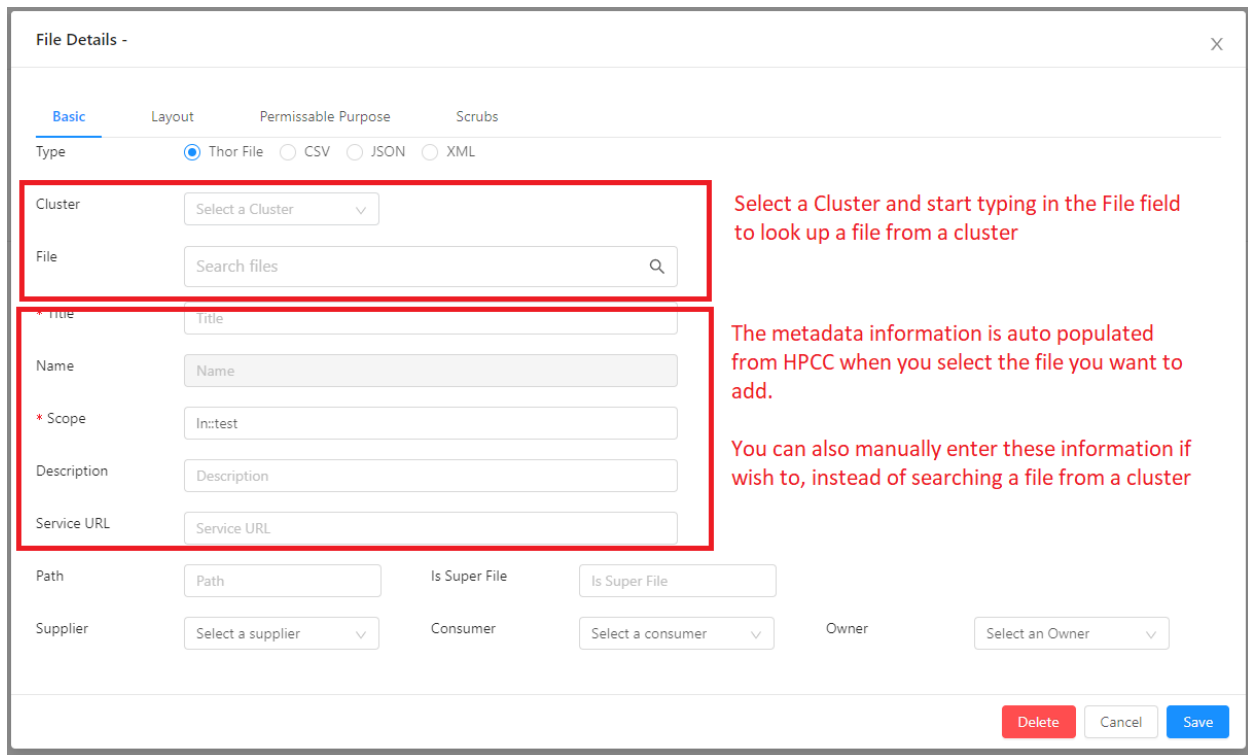- Function/Job (HPCC Jobs/other jobs)

# Files

Files can be added through Files option under Assets in the navigation.



Click Add button under each asset type to add respective asset.

## Files Details

## File Layouts:

Layouts for any files that is looked up directly from cluster will be auto populated. But you can also manually add Layout information for a file using 'Add a row' option.



## License Restrictions for files.

If you have any licensing restrictions for your files, record them here. The list of licenses are configurable in the system.

## Scrubs (WIP)

You can configure Rules for each fields in your files to be then consumed by a downstream application/job.

PS - This feature is currently a work in Progress



## File Preview

A preview of data. This tab will be shown only if your Tombolo Role has access to see the file data



## Dataflows – Shows the Tombolo Dataflows this file belongs to

# Indexes

Click on the Index option on the left nav to view the Indexes that are already added to Tombolo. New Indexes can be added using Add button.

## Basic Info



## Source File

## Index fields

**Index Details**                                                                       ✕

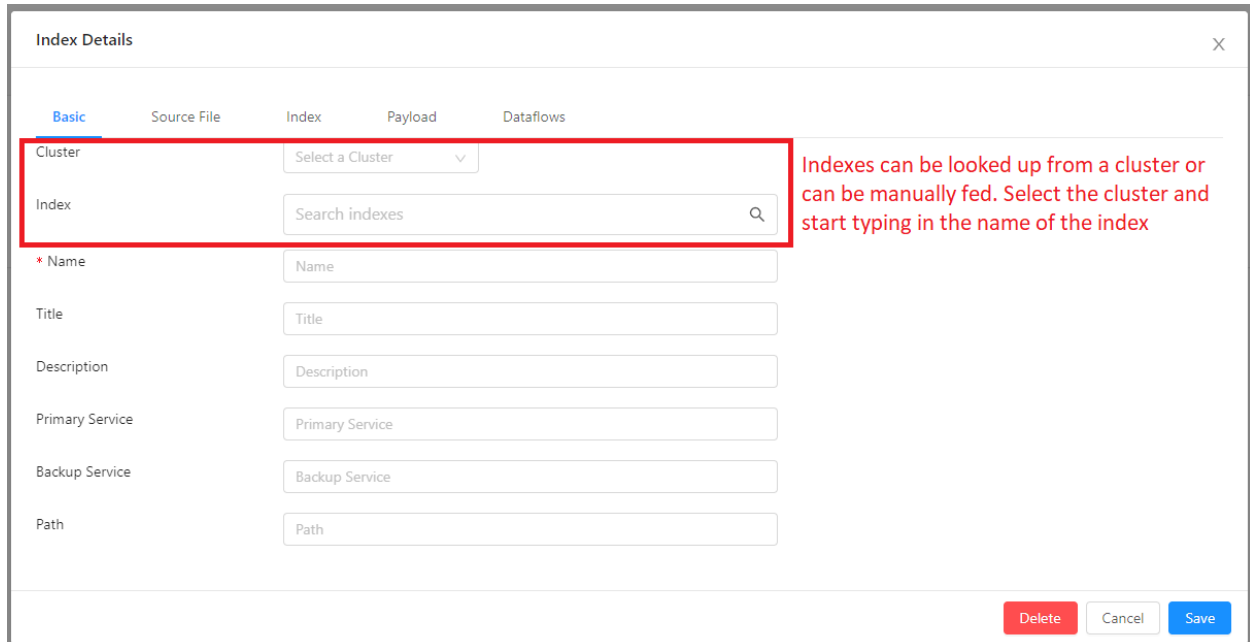| Basic | Source File | **Index** | Payload | Dataflows |

| Name | Type | Action |
|------|------|--------|
| imported_seconds_utc | Integer | 🗑 |
| userid | String | 🗑 |

*Key fields for the Index - auto populated from cluster*

Add a row

                                                        Delete   Cancel   Save

## Payload fields

**Index Details**                                                                       ✕

| Basic | Source File | Index | **Payload** | Dataflows |

| Name | Type | Action |
|------|------|--------|
| elevation | String | 🗑 |
| activitycalories | String | 🗑 |
| imported_time_utc | String | 🗑 |
| sedentaryminutes | String | 🗑 |
| floors | String | 🗑 |
| marginalcalories | String | 🗑 |
| lightlyactiveminutes | String | 🗑 |
| veryactiveminutes | String | 🗑 |

*Payload fields auto populated from cluster*

Add a row

                                                        Delete   Cancel   Save

**Dataflows** – Shows the Tombolo Dataflows this Index belongs to

# API/Gateway/Queries



## Input Fields

## Output Fields



**API/Gateway/Query Details**                                               ✕

| Basic | Input Fields | **Output Fields** |

| Name | Type | Possible Value | Value Description | Action |
|------|------|----------------|-------------------|--------|
| result_count | number | | | 🗑 |

Add a row

*Output fields of a query is identified automatically from cluster. User can also add custom fields by clicking Add a row*

## Job/Functions



**Job Details**                                                             ✕

| Basic | Input Params | Input Files | Output Files | Dataflows |

Cluster          Select a Cluster ⌄

Job              Search jobs 🔍

Name             Ingest_JH_data

Title            Ingest_JH_data

Description      Description

Git Repo         Git Repo

Entry BWR        Ingest_JH_data

Contact          Contact                          Author        Author

Job Type         ⌄

☐ Automatically create dependant files            Delete   Cancel   Save

*Search for a job from a cluster to retrieve some metadata*

*Capture contact info, Author of jobs here*

*If the job source resides in a GitHub repo, you can configure that as well.*

# Input Files

## Job Details      ✕

| Basic | Input Params | **Input Files** | Output Files | Dataflows |

Input Files:   Select Input Files ⌄   **Add**

| Name | Description |
| --- | --- |
| hpccsystems::covid19::file::public::johnhopkins::world.flat | |
| hpccsystems::covid19::file::public::johnhopkins::us.flat | Input files for HPCC Jobs are auto populated |
| hpccsystems::covid19::file::public::uscountypopulation::population.flat | |
| hpccsystems::covid19::file::public::worldpopulation::population_gender.flat | |
| hpccsystems::covid19::file::public::uspopulation::population.flat | |

‹   1   ›

☐ Automatically create dependant files     **Delete**   Cancel   **Save**

# Output files

## Job Details      ✕

| Basic | Input Params | Input Files | **Output Files** | Dataflows |

Output Files:   Select Output Files ⌄   **Add**

| Name | Description |
| --- | --- |
| hpccsystems::covid19::file::input::source::jh::level3 | Output files for HPCC jobs auto populated from |
| hpccsystems::covid19::file::input::source::jh::level2 | cluster. Files already existing in Tombolo can |
| hpccsystems::covid19::file::input::source::jh::level1 | also be added here using Output Files dropdown |

‹   1   ›

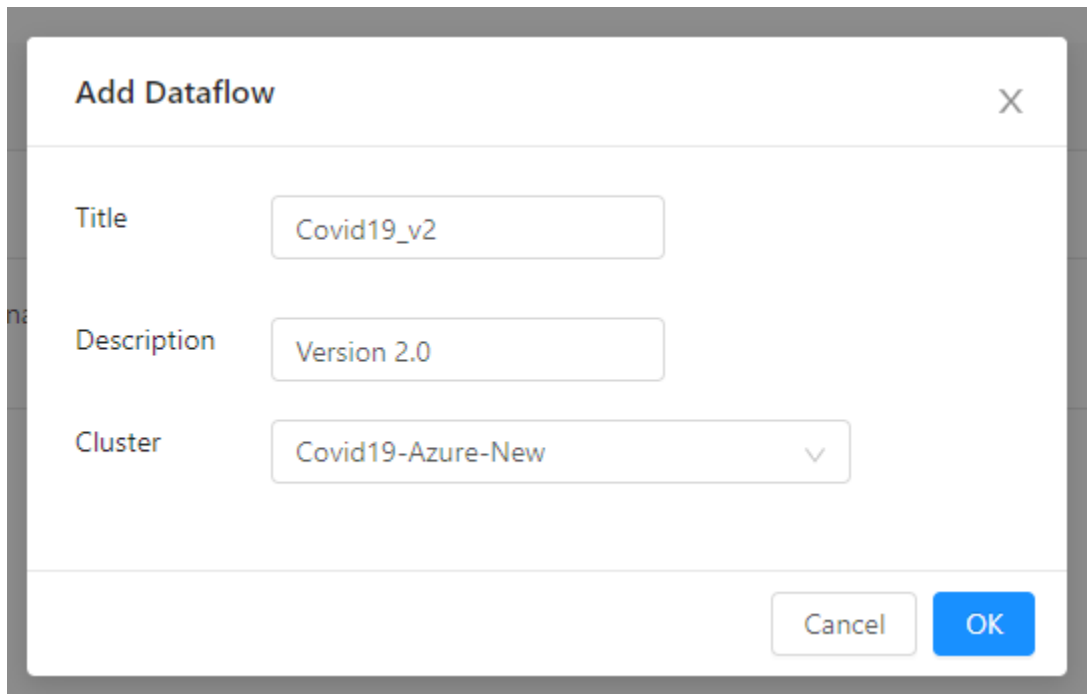☐ Automatically create dependant files     **Delete**   Cancel   **Save**

# Workflow Definitions

Capturing Data Lineage of a Data Lake is a key feature of Tombolo.

To create a Dataflow, click on Definitions under Workflow in the navigation. Dataflows that are already created will be listed. Click on Add and select a Cluster to which you want to point the dataflow. The cluster selection will be used later for automating tracking of workflows.
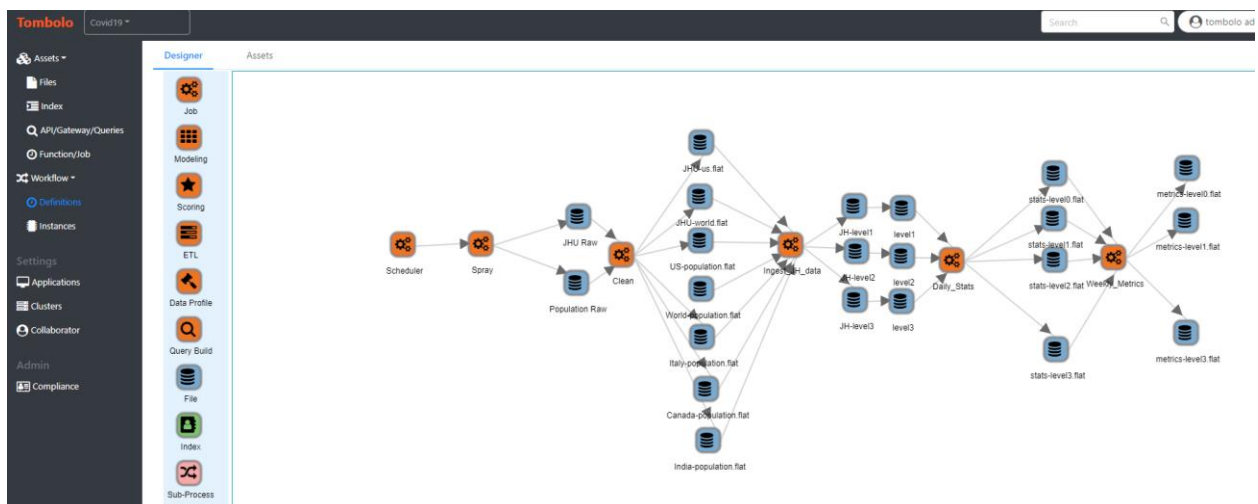


Once the Dataflow is created, click on the Dataflow name to view the Designer.

The Palette contains various nodes that are supported currently. Even though all the Jobs captures the same metadata, the idea is to capture job specific metadata in the future.

- Job – Any ECL Job
- Modelling – ML Modelling job
- Scoring – ML Scoring Job
- ETL – Any ETL job
- Data Profile – To run a Data Profile job
- Query Build – A job that builds and publishes roxie query
- File – Logical File/CSV/JSON/XML
- Index – An HPCC Index
- Sub-Process – A sub-process (child Dataflows within the main dataflow)

To use a node in the Dataflow, click on the node in the left pallet and drag it to the Designer.

The nodes can be associated with any of the asset (File/Index/Job/Query) by double clicking on it. It will then open the same Details dialog where you can either lookup an asset from a cluster or manually add the metadata.

## Designer Controls

Add a node to the designer – select the node from palette and drop to the designer

Add node details – Double click on a node

Connecting nodes – Keep holding Shift key and drag mouse from Source node to target node

Delete a node – Click on the node and press Delete button or use the delete icon that shows up on the node

Delete a connection – select the connection and press Delete button

Move a node – select the node and drag the mouse to where the node needs to be moved.

Zoom in/Zoom out – Place mouse on the designer and roll the scroll control on the mouse up/down

## Dataflow Instances

Tombolo has live workflow support to track what is happening in your workflow. Workflow tracking is done using Kafka as the integrator. This would mean that your ECL jobs will have to integrate with Kafka.