# Trino and Starburst Training Series:
## Transformation processing with PyStarburst

v1.0.0

## Session 5 of 5:  [Full Series Information](#)

Prior sessions

1. https://www.starburst.io/resources/creating-querying-data-lake-tables-on-demand/
2. https://www.starburst.io/resources/modern-table-formats-apache-iceberg-on-demand/
3. https://www.starburst.io/resources/data-pipelines-views-data-products/
4. https://www.starburst.io/resources/experience-warp-speed-in-action/

This workshop is focused on introducing & exercising the PyStarburst DataFrame API.  These are the goals for this session.

- Differential SQL-based data engineering from programming-based.
- Understand how PyStarburst implements lazy execution with Starburst Galaxy.
- Explore the DataFrame API.
- Write Python code to perform analytical questions and transformation processing.
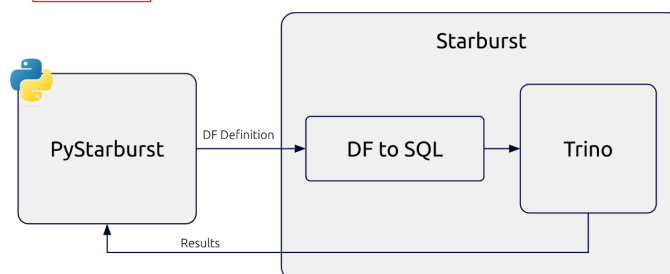
# Table of Contents

## PyStarburst Overview



```
df_missions = df_missions.with_column("date", f.sql_expr("COALESCE(TRY(date_parse(\"date\", '%a %b %d, %Y %H:%i UTC')), NULL)"))

print(df_missions.schema)

df_missions = df_missions\
    .filter(col("date") > datetime(2000, 1, 1))\
    .sort(col("date"), ascending=True)

df_missions.show()
```

- PySpark-like Syntax
- Lazy execution
- Python gets converted to SQL
- Heavy lifting done by Trino

Links: [API Docs](#) & [Example Code](#)

Currently supported on Starburst Galaxy and Starburst Enterprise.
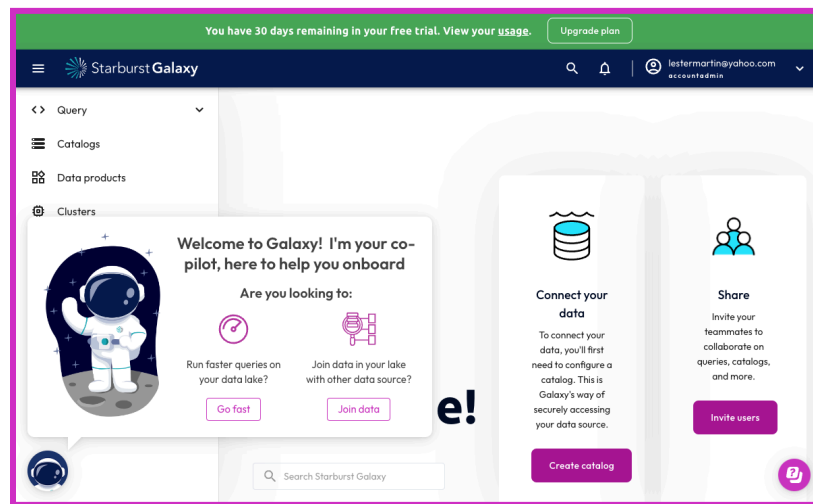
# Create Starburst Galaxy account

To sign up for Starburst Galaxy, follow the instructions on the free registration page at https://www.starburst.io/platform/starburst-galaxy/start/.

**Note**: You will receive an "invitation" email.  Please check your spam or junk folder if it does not immediately arrive in your inbox.
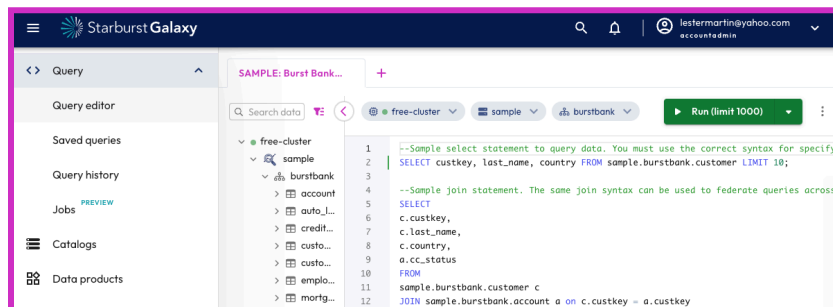
After you have entered your confirmation code, set a password, and selected your new domain name, you will be presented with a series of questions about your desired usage for Starburst Galaxy. Complete these with whatever you choose to share.

Eventually, you will likely be presented with a page similar to the following screenshot.



Click on the astronaut helmet icon in the lower-left to silence the pop-up coming from it.

At this point, you should see something similar to this to indicate you are fully configured.

# PyStarburst lab exercises

The code used in this webinar is in a Jupyter notebook file that is available at
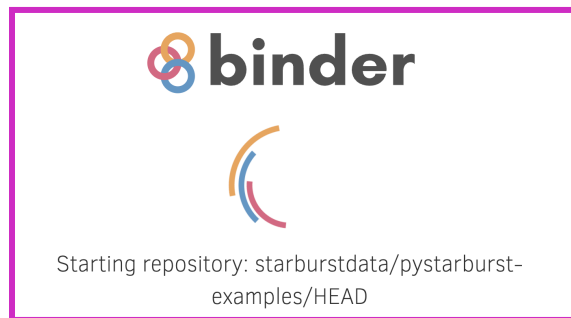https://github.com/starburstdata/pystarburst-examples/blob/main/notebooks/tpch.ipynb.

If you already have access to a jupyter instance, you can simply upload this notebook into it.  If you do not, visit https://github.com/starburstdata/pystarburst-examples and click on the button labeled "**launch binder**".



Alternatively, click on https://mybinder.org/v2/gh/starburstdata/pystarburst-examples/HEAD if you cannot access GitHub directly.
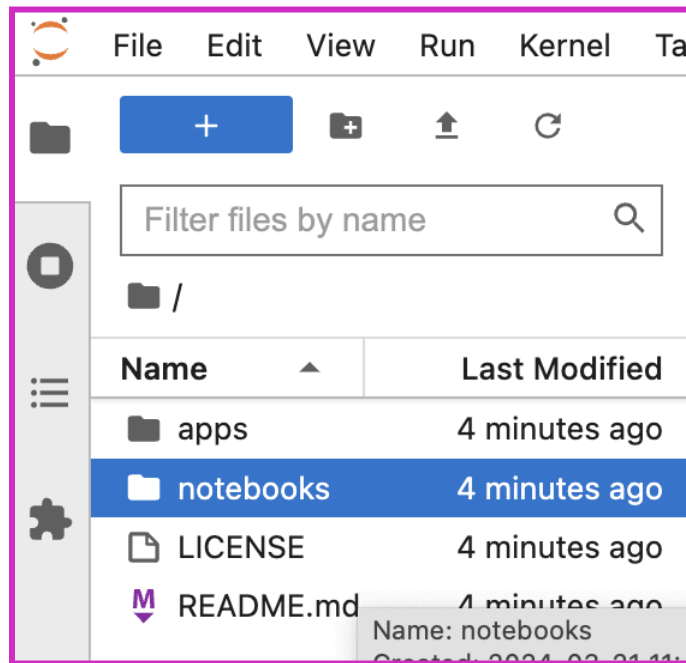
Regardless of which route you choose, your browser will look like this once launched.
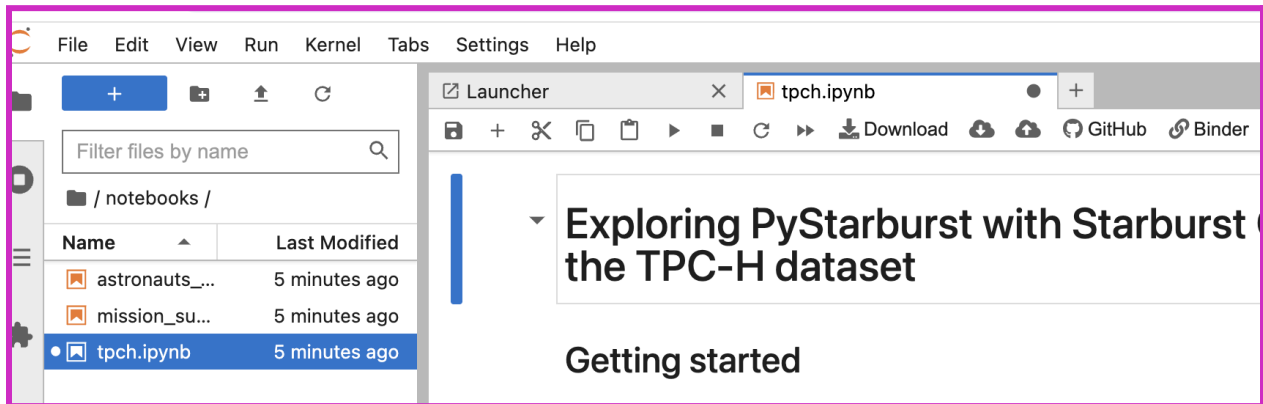


After a short time, a jupyter web-based notebook system will load.

Now, double-click on **notebooks** in the explorer pane on the left.



Then double-click on **tpch.ipynb** in the same explore pane on the left which will render the notebook that contains the code for this lab.



The instructor will walk through executing this notebook, cell by cell, and will provide additional information and discussions on each one.  If reviewing this document after the live webinar, please feel free to view the following YouTube video with the same discussions.

▶ Using PyStarburst on Starburst Galaxy with the TPC-H dataset