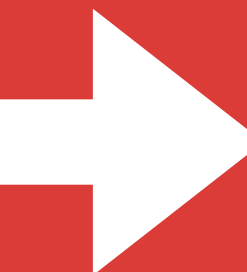


HAVE YOU EVER WANTED TO

BUILD YOUR OWN DATA SOURCE

○ ○ ○

```
1 # Register my Data Source
2 spark.dataSource.register(MyNewDataSource)
3
4 # Load data from your own data source
5 spark.read.format('MyNewDataSource').load()
6
7 # Write data to your own data source
8 spark.write.format("MyNewDataSource").save()
```



THE BASICS

Sign Up to the newsletter
@ databricks.news

Find more databricks tips @
dailydatabricks.tips

LOAD A CUSTOM DATA SOURCE AND REGISTER IT IN YOUR SESSION

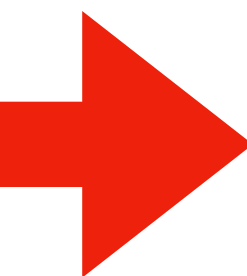
There are some great examples already available including ones that support the Faker library in the `pyspark-data-sources` library

Once you've loaded the community library, you can easily load a custom data source e.g. Google Sheets

These custom data sources can also leverage Sparks parallelism and file handling capabilities

〇〇〇

```
1 #install the community library with pip install pyspark-data-sources
2 from pyspark_datasources import GoogleSheetsDataSource #Import Google Sheets
3 spark.dataSource.register(GoogleSheetsDataSource)#Register Google Sheets DataSource
4
5
6 df = spark.read \
7     .format("googlesheets") \
8     .options(url="<your_google_sheets_document_url>") \
9     .load()
```



WHAT DOES THIS MEAN FOR ME?

Sign Up to the newsletter
@ [databricks.news](#)

Find more databricks tips @
[dailydatabricks.tips](#)

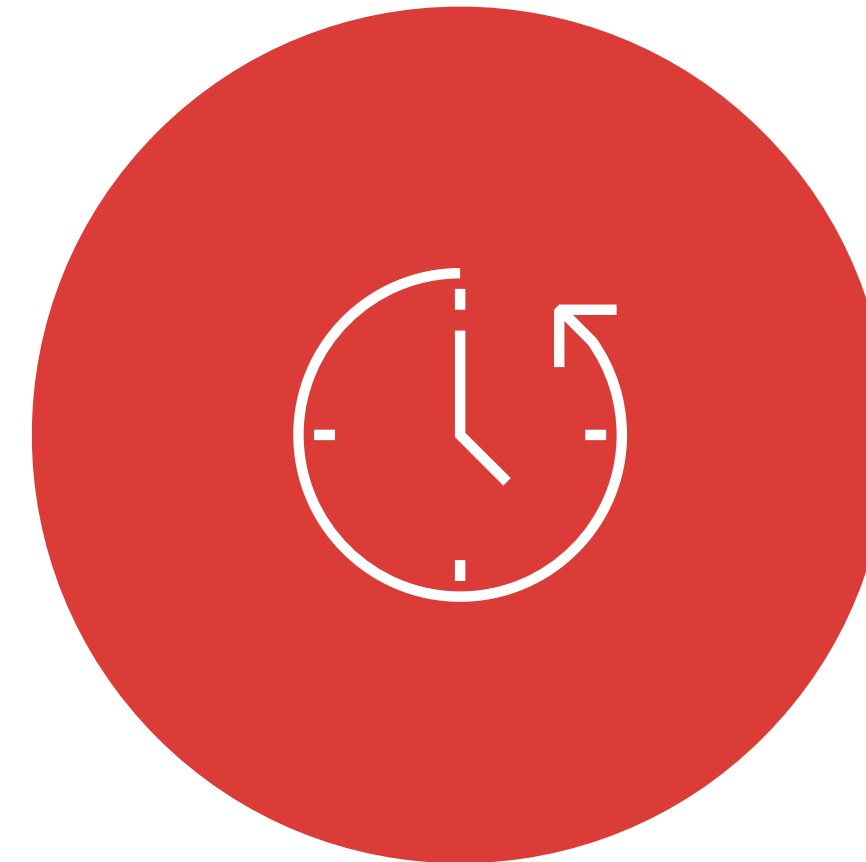


CAPABILITIES

Build Custom Spark connectors
in Python to connect to any data
source.

Connect your own APIs, file
formats or databases

Leverages Apache Arrow
For performance



ANY CADENCE, SAME CODE

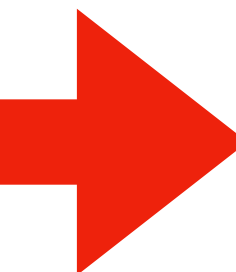
Your Custom Datasource Class
Works for batch +
streaming reads/writes.



UNIFIED EXECUTION

Simple Integration with
External Services

Integration with Declarative
Pipelines or Pyspark
Notebooks



MAKE YOUR OWN DATA SOURCE

Sign Up to the newsletter
@ databricks.news

Find more databricks tips @
dailydatabricks.tips

```
class OpenMeteoDataSource(DataSource):  
    """  
    Open-Meteo weather data source for PySpark.  
    """  
    @classmethod  
> def name(cls) -> str: ...  
  
> def schema(self) -> str: ...  
  
> def reader(self, schema: StructType) -> DataSourceReader: ...
```

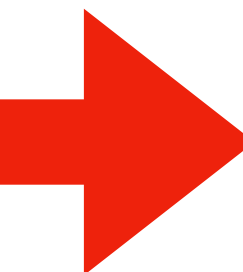
Create a Datasource for your custom data source like above.

LOAD THE DATA SOURCE LIBRARY

This Example is the Open-Meteo Weather API that provides access to high-resolution weather data (1-11km resolution) from national weather services

○○○

```
1 import openmeteo_datasource # You might want to move your datasource to it's own py file  
2 from pyspark.sql.datasource import DataSource, DataSourceReader #Import to register data sources  
3  
4 spark.dataSource.register(openmeteo_datasource.OpenMeteoDataSource) # Register Data Source  
5  
6 df = spark.read \  
7     .format("openmeteo") \  
8     .option("latitude", "52.52") \ # add your own data source specific options e.g. co-ordinates  
9     .option("longitude", "13.41") \  
10    .option("start_date", "2024-01-01") \ # Options that can support features like predicate  
    pushdown or pagination  
11    .option("end_date", "2024-01-31") \  
12    .load()
```



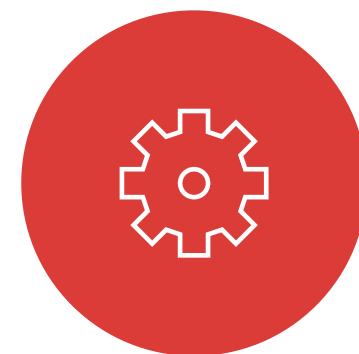
FUTURE CAPABILITIES



COMMUNITY CONNECTORS



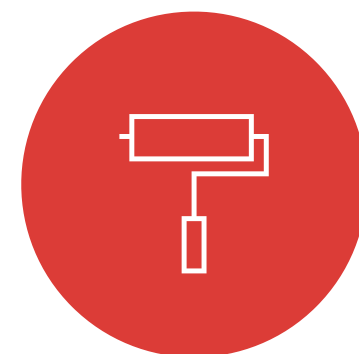
FILTER PUSHDOWN



PERFORMANCE OPTIMISATIONS



COLUMN PRUNING



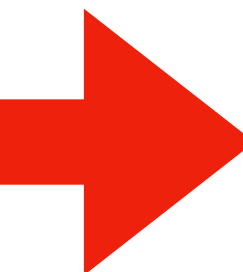
CUSTOM STATISTICS



OBSERVABILITY IMPROVEMENTS

Sign Up to the newsletter
@ databricks.news

Find more databricks tips @
dailydatabricks.tips



GOT QUESTIONS?

JOIN ME ON
NOV 20TH

FOR PYSPARK DATASOURCE
WEBINAR

SIGN UP TO MY DATABRICKS NEWSLETTER
FOR DETAILS



Databricks.**news**