

**SQL**  
KONFERENZ  
2025



# Mastering Spark Notebooks

A Guide for Data Engineers  
Version 2025



## Tillmann Eitelberg

CEO, oh22information services GmbH  
Microsoft Data Platform MVP



## Oliver Engels

CEO, oh22data AG  
Microsoft Data Platform MVP



**SQL**  
KONFERENZ  
2025

# Agenda

- Getting Started with Spark Notebooks
- Efficient coding with Python in Spark Notebooks
- Environments – Setting up your Spark workspace
- Mastering data manipulation in Spark
- Delta Tables – Optimizing data storage and performance
- Semantic Link – Connecting your data for a smarter workflow
- Simplifying data preparation with Data Wrangler
- Supercharging your Spark workflow with integrated tools



# What we can't cover today

- Fabric CI/CD

- [Notebook source control and deployment - Microsoft Fabric | Microsoft Learn](#)
- [Overview of Fabric Git integration - Microsoft Fabric | Microsoft Learn](#)

- Microsoft Spark Utilities

- [Introduction to Microsoft Spark utilities - Azure Synapse Analytics | Microsoft Learn](#)

- Python Workshop

- [Documentation \(python.org\)](#)
- [Python Tutorials – Real Python](#)



# Getting Started with Spark Notebooks



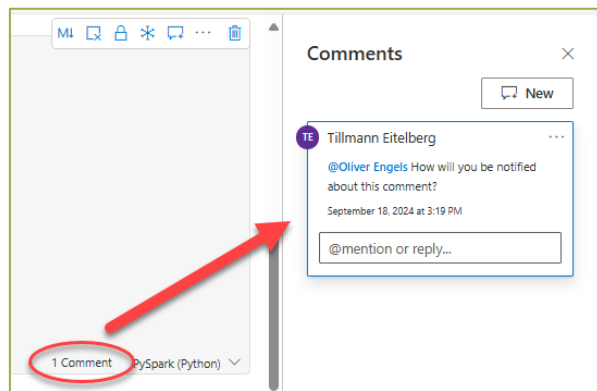
# Notebooks – The basics

- Add, delete and move rows
- Lock or freeze cells
- hide input or output
- Merge or split rows
- Shift + Enter / Strg + Enter
- Add Code and Markdown Cells
- Adding ipywidgets (IPython Widgets)
- Toggle parameter cell (used in pipeline activity)
- Fabric notebooks recognize the standard Jupyter .ipynb files
- [Fabric Notebook known limitation](#)




# Comments

- Creating comments / threads
- Comments are assigned to individual cells and can be accessed from there
- People can be mentioned in comments
- Comments can be closed as a thread



# Endorse Fabric Items

- Various objects within Fabric can be endorsed
- By default, users can promote notebooks (and other objects)
- With special permissions, they can also be set as certified
- Individual objects that contain data (currently only lakehouses and semantic models) can also be given a master data badge

Name	Type	Opened	Owner	Endorsement	Sensitivity	Location
#FABCONEUROPE	Workspace	now	—	—	—	Workspaces
03 - Snippets	Notebook	11 minutes ago	Tillmann Eitelberg	 Promoted	—	#FABCONEUROPE
02 - Magic Commands	Notebook	34 minutes ago	Tillmann Eitelberg	—	—	#FABCONEUROPE
01 - Basics	Notebook	an hour ago	Tillmann Eitelberg	—	—	#FABCONEUROPE





# Magic Commands

```
1 %%pyspark
2
3 result = spark.sql("SELECT * FROM random_numbers")
4
5 display(result)
```

[11] ✓ - Command executed in 1 sec 467 ms by Tillmann Eitelberg on 5:26:57 PM, 9/17/24

PySpark (Python) ▾

	12L Spalte1	12L Spalte2	12L Spalte3	12L Spalte4	12L Summe	
1	97	38	96	43	274	
2	78	23	34	72	207	
3	5	62	46	15	128	

```
1 %%sql
2 SELECT * FROM random_numbers
```

[12] ✓ - Command executed in 1 sec 451 ms by Tillmann Eitelberg on 5:26:59 PM, 9/17/24

Spark SQL ▾

	12L Spalte1	12L Spalte2	12L Spalte3	12L Spalte4	12L Summe	
1	97	38	96	43	274	
2	78	23	34	72	207	
3	5	62	46	15	128	

**SQL**  
KONFERENZ  
2025



# Magic Commands

```
1 # %who_ls
2 # %who
3 %whos
```

[13] ✓ - Command executed in 250 ms by Tillmann Eitelberg on 5:26:59 PM, 9/17/24 PySpark (Python) ▾

Variable	Type	Data/Info
HiveContext	type	<class 'pyspark.sql.context.HiveContext'>
streamingContext	type	<class 'pyspark.streaming.Context'>
col	function	<function col at 0x7809f7ec8d30>
data	list	n=3
df	DataFrame	DataFrame[Spalte1: bigint<...>4: bigint, Summe: bigint]
display	function	<function display at 0x7809eee4d360>
displayHTML	DisplayHTML	<notebookutils.visualizat...>object at 0x7809eee2b220>
exit	function	<function trap_exit at 0x7809ed30ae00>
initializeHContext	function	<function initializeHContext at 0x7809ed30b910>
msparkutils	module	<module 'notebookutils.ms...'>
notebookutils	module	<module 'notebookutils' f...>
prepare	function	<function prepare at 0x7809f6eb6290>
random	module	<module 'random' from '/h...>
result	DataFrame	DataFrame[Spalte1: bigint<...>4: bigint, Summe: bigint]
sc	SparkContext	<SparkContext master=yarn...>
spark	SparkSession	<pyspark.sql.session.Spark...>
spark_sum	function	<function sum at 0x7809f70c1990>
sqlContext	SQLContext	<pyspark.sql.context.SQLC...>

```
1 data
```

[14] ✓ - Command executed in 240 ms by Tillmann Eitelberg on 5:27:00 PM, 9/17/24 PySpark (Python) ▾

```
[(97, 38, 96, 43), (78, 23, 34, 72), (5, 62, 46, 15)]
```

```
1 %xdel data
```

[15] ✓ - Command executed in 283 ms by Tillmann Eitelberg on 5:27:01 PM, 9/17/24 PySpark (Python) ▾

```
1 data
```

[16] ✗ - Command executed in 276 ms by Tillmann Eitelberg on 5:27:03 PM, 9/17/24 PySpark (Python) ▾

```
NameError                                Traceback (most recent call last)
Cell In[47], line 1
----> 1 data
NameError: name 'data' is not defined
```

- List all variables with “%whos”
- Display variable
- Delete variables with “%xdel”
- %%capture - captures the output of the entire cell

SQL  
KONFERENZ  
2025



# Magic Commands

`%%capture install_log`

- Suppresses or captures the output of a Jupyter cell
- Handles both stdout (regular output) and stderr (error messages)
- Can store the output in a variable for later use
- Useful for hiding long or unnecessary outputs (e.g., pip install)
- Does not affect the code execution

```
# Access the captured output later
print("Captured stdout:", captured_output.stdout)
print("Captured stderr:", captured_output.stderr)
```



SQL  
KONFERENZ  
2025

# Magic Commands

[Format code in Microsoft Fabric notebooks](#)  
[Microsoft Fabric | Microsoft Learn](#)

`%load_ext jupyter_black`

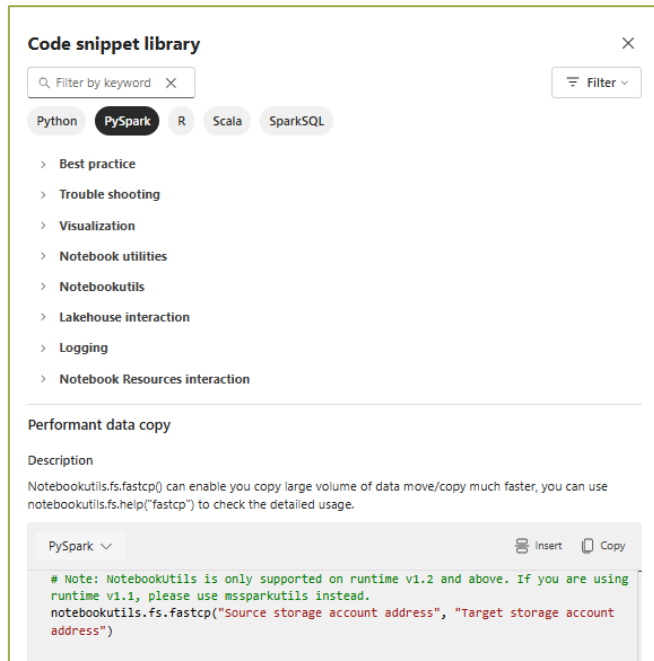
- extension for Jupyter Notebooks that auto-formats Python code using the Black code formatter
- enforces a consistent code style by reformatting code to follow best practices
- removes unnecessary whitespace, corrects indentation, and improves readability

```
data = {"A": np.random.randint(1, 100, 10), "B":  
np.random.randint(1, 100, 10), "C": np.random.randint(1, 100, 10), "D": np.random.randint(1, 100, 10)}  
df = pd.DataFrame(data)  
df = pd.DataFrame(data)
```



# Code snippet library

- Library with various example code on individual topics
- Can be filtered by language
- Code can be transferred directly to a notebook as a new line
- Individual variables / placeholders must be filled in the code
- Import of additional libraries is not always clearly defined
- Snippets can also be added via IntelliSense
- The examples provided do not appear to be identical



# Copilot

- Load a specific package for copilot
- Install Magic Commands for Copilot

#Run this cell to install the required packages for Copilot

```
%pip install https://aka.ms/chat_magics-0.0.0-py3-none-any.whl
```

```
%load_ext chat_magics
```



Mistakes are possible.  
Carefully review content created  
by AI before using it.



**SQL**  
KONFERENZ  
2025

# Copilot

No pandas imported in the current notebook  
File path and name is incorrect  
The context menu does a much better job

The screenshot shows the Databricks Copilot interface. On the left, a sidebar lists sources: All sources, Lakehouses, and sample\_datasets. The main area displays a notebook cell with the following content:

```
1 # Welcome to your new notebook
2 # Type here in the cell editor to add code!
3
```

Below the code, a status bar indicates: [2] ✓ 6 sec - Command executed in 272 ms by Tillmann Elteberg on 20:10 PM, 9/17/24 PySpark (Python) ✓

On the right, the Copilot panel shows a preview of generated code. A red box highlights the following code snippet:

```
df = pd.read_csv('/lakehouse/default/files/city_safety.csv')
```

Below the code, the Copilot panel explains the assumptions made:

While generating the code, I made the following assumptions:  
It is assumed that pandas is imported as pd

At the bottom of the Copilot panel, there are input fields for prompts like "What are the top 5 safest cities in the dataset?" and "Remove any duplicate entries in the DataFrame." and a "Type a request or click the spark icon for ideas." field.

SQL  
KONFERENZ  
2025



# Copilot

- Prompt "Analyze df and provide insights about the data"
- Output only in the Copilot window
- The summary is "OK"
- No code is copied into the notebook
- Further questions are a nice idea
- Other questions can also create tables, but also only in the Copilot window

## Data Insights:

### • Data Types:

- The 'dataType' column contains 100 unique values, indicating various types of incidents or responses.
- The 'category' column has 5189 unique values, representing different locations or addresses.

### • Location Information:

- Latitude ranges from 47.4956 to 47.7766, with a mean of approximately 47.6203.
- Longitude is constant at -122.331.

### • Temporal Information:

- The dataset includes entries from various dates and times, but the 'dateTime' column seems to have missing values.

### • Response Details:

- Most common response type is 'Aid Response', occurring 3549 times.

### • Missing Values:

- Columns like 'dataSubtype', 'dateTime', 'subcategory', and 'status' have missing values.

### • Source Information:

- All entries have the same source and extended properties values (2020 and 7 respectively).

These insights provide a basic understanding of the dataset's contents, including the types of incidents recorded, location details, temporal aspects, common response types, and presence of missing values in certain columns. Further analysis could involve exploring relationships between variables or trends over

What are the top 5 categories with the highest number of incidents?

Combine similar subcategories into broader categories for better analysis.

SQL  
KONFERENZ  
2025





## Copilot – using Chat-magics

- Chat Magics bring the magic of Copilot to your notebook
- The `%%chat` command allows you to ask questions about the state of your notebook
- The `%%code` enables code generation for data manipulation or visualization
- `%pin DataFrameName` to help the AI to focus on specific dataframes
- `%new_task` clear the AI to focus on a new task



SQL  
KONFERENZ  
2025

# Efficient coding with Python in Spark Notebooks



# Modularize Your Code

- Functions and modules allow code reuse, reducing duplication
- Improves readability and maintenance by breaking down complex tasks
- Modular code is easier to scale and modify as projects grow
- Isolate components for easier testing and bug fixing
- Centralize logic in functions/modules to avoid copy-pasting and reduce errors and thereby gain clean structure, better collaboration, version control, and easier maintenance
- Default Parameters: Flexibility in functions with minimal input, reducing repetitive code while allowing customization

- `%run <notebook>`

vs.

```
mssparkutils.notebook.run("notebook path", <timeoutSeconds>, <parameterMap>)
```

- The upper limit for notebook activities or concurrent notebooks in `notebookutils.notebook.runMultiple()` is 50.
- The statement depth for `%run` is up to 5, with a total of 1000 referenced cells.



SQL  
KONFERENZ  
2025



# Use Decorators for Reusable Logic

- Functions that modify the behavior of other functions or methods, without changing their code
- Allows code reuse for common patterns like logging, timing, authentication, caching, etc.
- Uses the `@decorator_name` syntax above the function to apply the decorator
- Decorators take a function as an argument and return a new function
- Multiple decorators can be stacked on top of a function



SQL  
KONFERENZ  
2025

# Use Decorators for Reusable Logic

```
1 def timeit(func):
2     def wrapper(*args, **kwargs):
3         start_time = time.time()
4         result = func(*args, **kwargs)
5         print(f"Execution time: {time.time() - start_time} seconds")
6         return result
7     return wrapper

[2] ✓ <1 sec - Command executed in 245 ms by Tillmann Eitelberg on 4:58:16 PM, 9/18/24 PySpark (Python) ✓
```

```
1 @timeit
2 def run_transformation(df):
3     return df.filter(df["age"] < 18)

[3] ✓ <1 sec - Command executed in 244 ms by Tillmann Eitelberg on 4:58:16 PM, 9/18/24 PySpark (Python) ✓
```

```
1 df = spark.sql("SELECT * FROM lakehouse.dbo.fabcon01 LIMIT 1000")
2 df_trd = run_transformation(df)
3

[4] ✓ 1 min 3 sec - Command executed in 1 min 2 sec 912 ms by Tillmann Eitelberg on 5:06:15 PM, 9/18/24 PySpark (Python) ✓
```

> Log ...

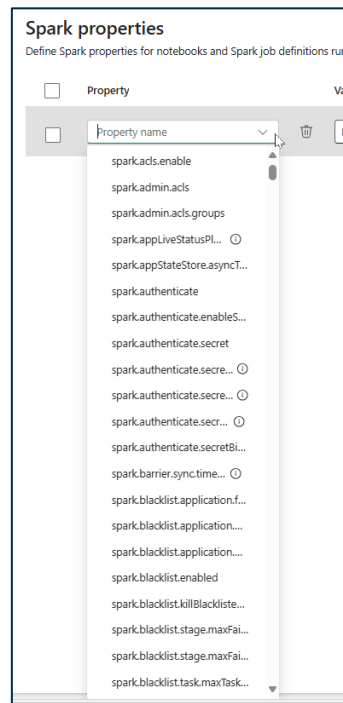
Execution time: 0.0027866363525390625 seconds

## **Environments – Setting up your Spark workspace**



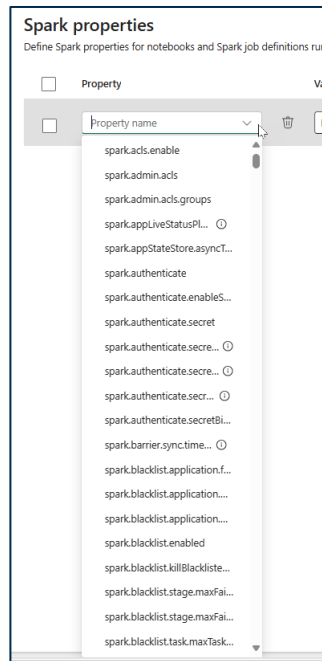
# Environments – Setting up your Spark workspace

- Enables the summarization of environmental parameters
- Add public libraries via PyPi in fixed versions (including older ones)
- Add your own non-publicly hosted packages
- Configuration of the Spark cluster from the number of driver cores to memory and executor instances
- Definition of various Spark properties
- Additional files can be added and accessed from the notebook.
- Environments must be saved and published
- Publishing takes several minutes
- Caution: a Spark Session with an environment that contains additional libraries takes longer to start  
(But the installation in the notebook otherwise also)



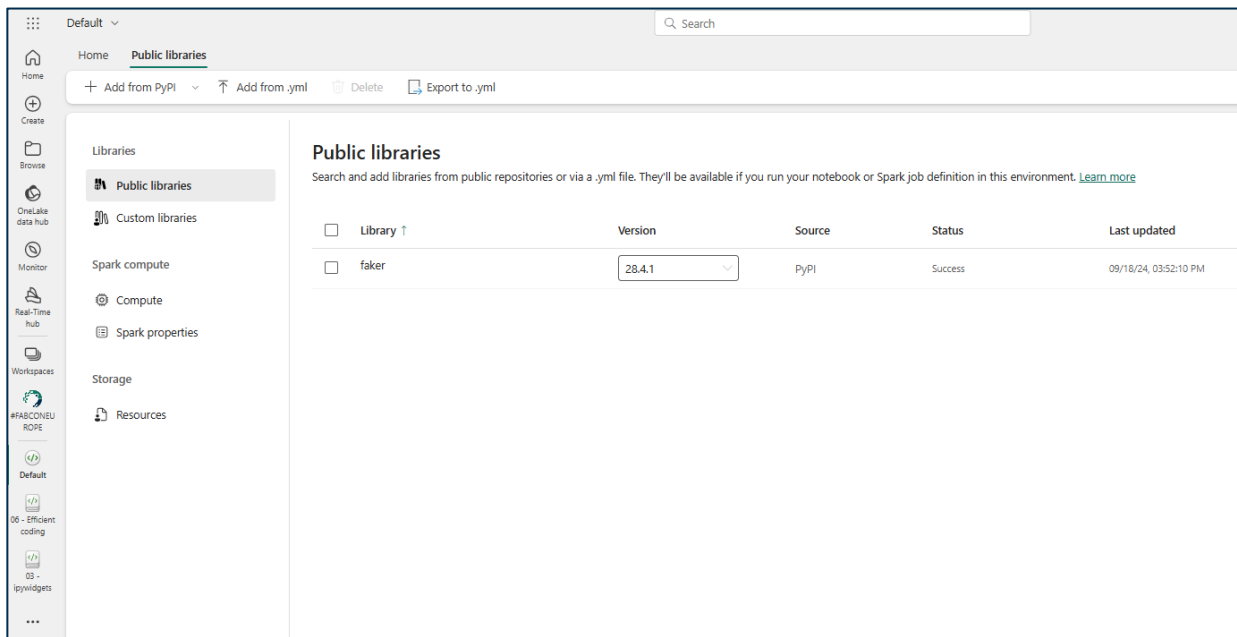
## Environments – Setting up your Spark workspace

- Enables the summarization of environmental parameters
- Add public libraries via PyPi in fixed versions (including older ones)
- Add your own non-publicly hosted packages
- Configuration of the Spark cluster from the number of driver cores to memory and executor instances
- Definition of various Spark properties
- Additional files can be added and accessed from the notebook.
- Environments must be saved and published
- Publishing takes several minutes
- Caution: a Spark Session with an environment that contains additional libraries takes longer to start (But the installation in the notebook otherwise also)





# Environments – Setting up your Spark workspace



The screenshot displays the Databricks workspace interface. The top navigation bar shows 'Default' and a search bar. The left sidebar contains various navigation icons and labels: Home, Create, Browse, OneLake data hub, Monitor, Real-time hub, Workspaces, #FABCONEU ROPE, Default, 06 - Efficient coding, 03 - ipynbwidgets, and a menu icon. The main content area is titled 'Public libraries' and includes a search bar and a table of libraries. The table has columns for Library, Version, Source, Status, and Last updated. A single library, 'faker', is listed with version '284.1' from PyPI, in a 'Success' status, and last updated on '09/18/24, 03:52:10 PM'.

Library ↑	Version	Source	Status	Last updated
<input type="checkbox"/> faker	284.1	PyPI	Success	09/18/24, 03:52:10 PM



**SQL**  
KONFERENZ  
2025

# Mastering data manipulation in Spark



# DataFrames



**SQL**  
KONFERENZ  
2025

# DataFrames

- Two-dimensional, tabular data structures with labeled axes (rows and columns), similar to a spreadsheet or SQL table
- Widely used for data manipulation, analysis, and processing tasks
- The primary framework in Spark is, unsurprisingly, Spark DataFrames
- Nevertheless, Pandas is the most widely used DataFrame framework
  - Pandas can be used in both Python and pySpark
  - Widest support through various libraries
- Developers have to be careful which DataFrame they use
  - Reconversion costs a lot of time and resources
  - Not all functions are directly supported by all data frames.



SQL  
KONFERENZ  
2025

# Mastering data manipulation in Spark

- Each DataFrame has a schema that defines the structure, meaning the column names and data types are known
- Load data from various sources (CSV, JSON, Parquet, Hive, etc.)
- Perform operations like filtering, selecting, joining, and aggregation to transform raw data into useful insights
- Save transformed data in formats like Parquet, ORC, CSV, or to databases like Hive or relational databases via JDBC



SQL  
KONFERENZ  
2025

# Mastering data manipulation in Spark

- **Spark DataFrames:** DataFrames are processed in parallel across a cluster, making them highly scalable for big data
- **Pandas:** DataFrames are processed locally on a single machine, suitable for small to medium data workloads
- Both have more or less the same functions, but different syntax
- Due to different functions and different internal structures, not all pySpark / Python functions and modules can always handle both DataFrames identically



SQL  
KONFERENZ  
2025



# Mastering data manipulation in Spark

- When converting data between Spark DataFrames and Pandas, several issues can arise:
  - Memory limits: Pandas operates in-memory, so converting a large Spark DataFrame into Pandas may exceed memory limits
  - Type mismatches: Data types supported by Spark may not map directly to Pandas (e.g., null handling, categorical data)
  - Performance: Spark is optimized for distributed systems, while Pandas is not. Converting large datasets to Pandas can slow down performance significantly



SQL  
KONFERENZ  
2025



# Mastering data manipulation in Spark

Description	Pandas Commands	Spark Commands
Select specific columns	<code>df[['column']]</code>	<code>df.select('column')</code>
Filter rows based on conditions	<code>df[df['column'] &gt; value]</code>	<code>df.filter(df['column'] &gt; value)</code> <code>df.where(df['column'] &gt; value)</code>
Group data and perform aggregations	<code>df.groupby('column').agg()</code>	<code>df.groupBy('column').agg()</code>
Add or modify a column	<code>df['new_column'] = df['existing_column'] * 2</code>	<code>df.withColumn('new_column', df['existing_column'] * 2)</code>
Remove columns	<code>df.drop(columns=['column'])</code>	<code>df.drop('column')</code>
Join two DataFrames	<code>df.merge(other_df, on='key')</code>	<code>df.join(other_df, on='key')</code>
Sort data	<code>df.sort_values(by='column')</code>	<code>df.orderBy('column')</code> or <code>df.sort('column')</code>
Remove duplicate rows	<code>df.drop_duplicates()</code>	<code>df.dropDuplicates()</code>
Select unique rows	<code>df['column'].unique()</code>	<code>df.distinct()</code>
Limit the number of rows	<code>df.head(n)</code>	<code>df.limit(n)</code>
Replace NULL/NaN values	<code>df.fillna(value)</code>	<code>df.fillna(value)</code>
Remove rows with NULL/NaN values	<code>df.dropna()</code>	<code>df.na.drop()</code>
Rename a column	<code>df.rename(columns={'old_name': 'new_name'})</code>	<code>df.withColumnRenamed('old_name', 'new_name')</code>
Convert DataFrame to Pandas (Spark only)	N/A	<code>df.toPandas()</code>





# DataFrames



SQL  
KONFERENZ  
2025

# DataFrames - Pandas-on-Spark, Dask, Polars

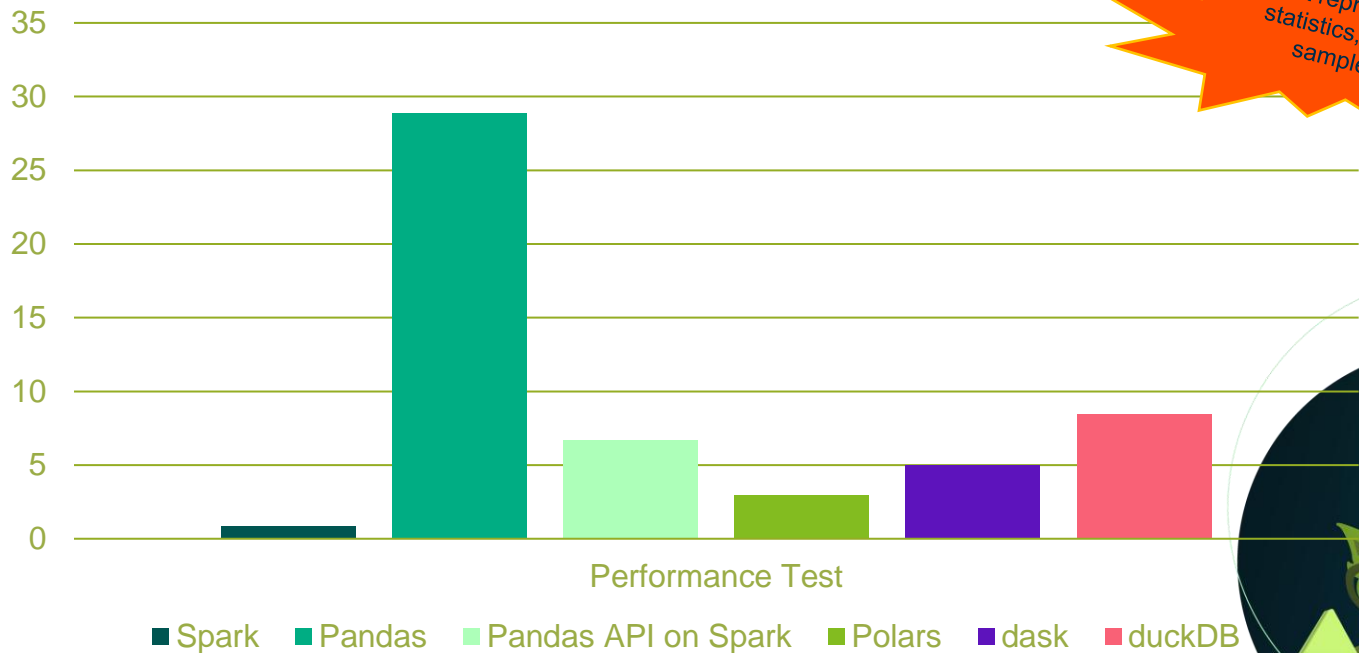
- **Pandas-on-Spark:** Combines the simplicity of Pandas with Spark's scalability, enabling Pandas-like operations on large, distributed data
- **Dask:** Scales Pandas-like operations across multiple cores or machines, handling datasets larger than memory
- **Polars:** A high-performance DataFrame library optimized for speed and multi-threading, handling both small and large datasets efficiently
- **DuckDB:** A lightweight, in-process SQL database optimized for fast analytical queries on large datasets, especially columnar formats like Parquet



SQL  
KONFERENZ  
2025



# Mastering data manipulation in Spark



Not representative  
statistics, just a few  
sample tests

SQL  
KONFERENZ  
2025



# Mastering data manipulation in Spark

- When using multiple DataFrame frameworks in the same notebook (e.g., Pandas, Spark, Dask, Polars), issues include:
  - Inconsistent APIs: Each framework has its own methods and syntax, leading to confusion or errors when switching
  - Data transfer overhead: Converting data between frameworks (e.g., Spark to Polars) can be time-consuming and resource-intensive
  - Memory management: Different frameworks handle memory differently, potentially leading to inefficiencies or crashes if not managed carefully
  - Library version conflicts: Different versions of libraries may be incompatible, leading to errors or issues that are sometimes difficult or impossible to resolve



# **Delta Tables – Optimizing data storage and performance**



# Delta Tables – Optimizing data storage and performance

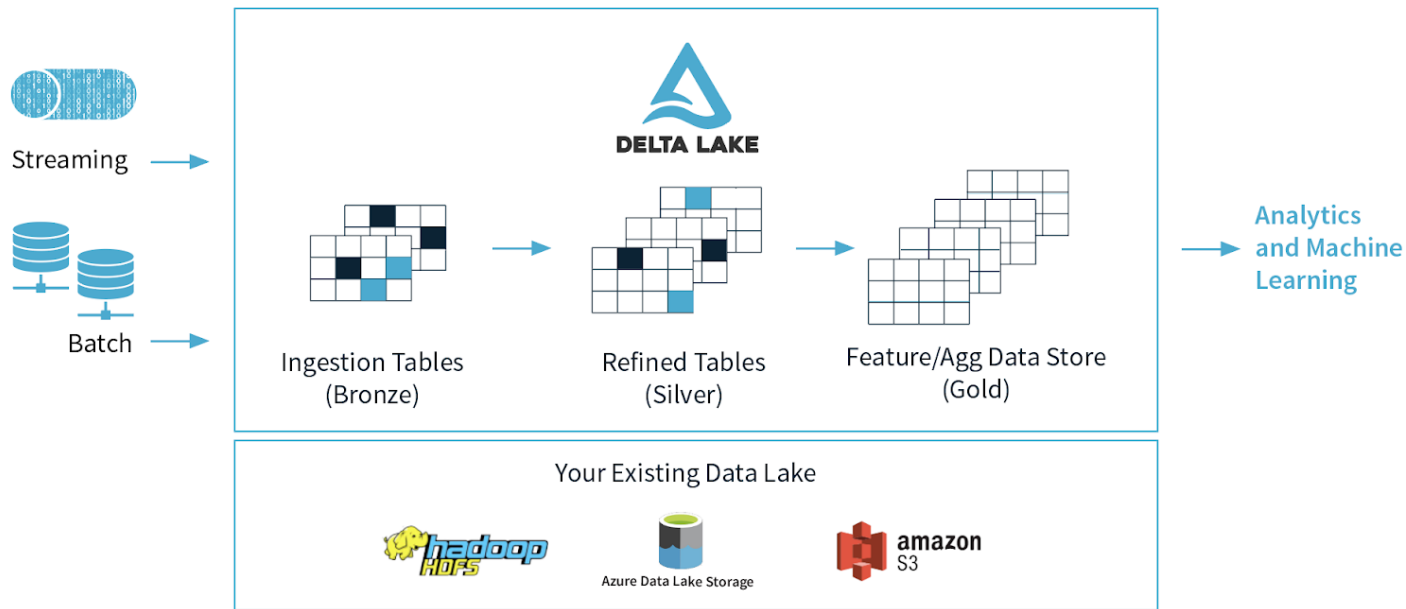


- Delta Tables combine the scalability of data lakes with ACID transactions for reliable data management
- Maintain historical versions of data, enabling time travel and auditability
- Support for upserts (MERGE), deletes, and updates directly in the table
- Can automatically evolve to accommodate new data without manual schema adjustments
- Use data indexing, caching, and compaction techniques to improve query speed
- Integrates seamlessly with Spark, Microsoft Fabric, Databricks and various cloud platforms for easy adoption



SQL  
KONFERENZ  
2025

# Delta Tables – Optimizing data storage and performance



## Delta Tables – Example commands

SQL Command	Description
CREATE TABLE AS (Delta)	Creates a new Delta table from the results of a query.
MERGE INTO	Performs upserts by merging new data into an existing Delta table based on a condition.
DESCRIBE DETAIL	Provides detailed information about the Delta table, such as version history and storage details.
DESCRIBE HISTORY	Shows the version history and metadata of a Delta table, useful for time travel and audit trails.
OPTIMIZE	Optimizes the performance by combining smaller files into larger ones in a Delta table.
VACUUM	Removes old versions of data files that are no longer needed to reclaim storage space.
RESTORE	Restores a Delta table to an earlier version.
TIMESTAMP AS OF	Queries a Delta table as of a specific timestamp.
VERSION AS OF	Queries a Delta table as of a specific version number.





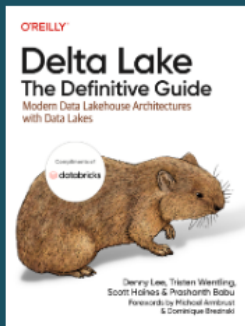
## Delta Tables – V-Order

- The Lakehouse and the Delta Lake table format are central to Microsoft Fabric
- V-Order is a write time optimization to the parquet file format that enables lightning-fast reads under the Microsoft Fabric compute engines, such as Power BI, SQL, Spark, and others
- V-Order sorting has a 15% impact on average write times but provides up to 50% more compression
- 100% open-source parquet format compliant
- Enabled by default in Microsoft Fabric

[Delta Lake table optimization and V-Order - Microsoft Fabric | Microsoft Learn](#)



# Delta Tables – Optimizing data storage and performance



## Delta Lake: The Definitive Guide

Building modern data lakehouse architectures with Delta Lake with forewords by Michael Armbrust and Dominique Brezinski

[https://delta.io/pdfs/dldg\\_databricks.pdf](https://delta.io/pdfs/dldg_databricks.pdf)

# Semantic Link – Connecting your data for a smarter workflow



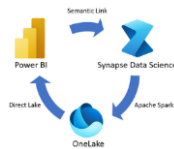
# Semantic Link

- The semantic link feature connects semantic models with Synapse Data Science in Microsoft Fabric
- Enhance data connectivity
- Propagate semantic information
- Bridges the gap between Power BI and the Data Science experience
- **SemPy** is the corresponding Python package that provides all the functions within Fabric and is Lib is preinstalled!!
- Semantic Link Labs – Additional library with more functions to simplify working with Microsoft Fabric

[microsoft/semantic-link-labs](https://github.com/microsoft/semantic-link-labs): Early access to new features for Microsoft Fabric's Semantic Link. ([github.com](https://github.com))



# Semantic Link [-Labs]



## • From your Fabric Notebook you can:

Home Edit Run View

Run all Connect PySpark (Python) env\_sampy Data Wrangler Copilot

### Get Workspaces combined with Semantic Models

This sample code produces an overview on all your workspaces including the semantic models

```
1 # List your accessible workspaces
2 wspaces = fabric.list_workspaces().sort_values(by='Name', ascending=True)
3
4 # Filter the list only for workspaces on a dedicated capacity
5 filtered_ws = pd.DataFrame(wspaces).query('["Is On Dedicated Capacity" == True]')[["Name", "Id"]]
6
7 # Loop the workspace list and build extract the underlying semantic models (incl. additional XMLA Props)
8 datasets = [
9     fabric.list_datasets(ws, additional_xmla_properties=['Model.DefaultMode', 'Model.DirectLakeBehavior', 'CompatibilityLevel']).assign(workspace=ws)
10     for ws in filtered_ws['Name']
11 ]
12
13 # Build a simple data catalog of all the datasets in your workspaces
14 ds_catalog = pd.concat(datasets, ignore_index=True)
15 display(ds_catalog)
```

[9] ✓ - Command executed in 51 sec 792 ms by Oliver Engels on 8:27:22 PM, 8/13/24

ABC Dataset Name	ABC Dataset ID	Created Timestamp	Last Update	ABC Model Default Mode	ABC Model Direct Lake Behavior	Compatibility Level	ABC workspace
Adventure Works ...	7410f999-0f8d-...	2019-09-17 10:50:29.000	NULL	Import	Automatic	1567.0	Market_LP_UDEV
Adventure Works ...	8294215f-a511-...	2019-09-17 05:50:29.000	NULL	Import	Automatic	1567.0	Market_DP_PRD
Adventure Works ...	5e8413d5-4c35-...	2019-09-17 05:50:29.000	NULL	Import	Automatic	1567.0	Market_DP_QAS

Functional Dependencies in your data

Insides in your Fabric landscape

Gover/Administer Fabric landscape

Manage Semantic Models at Scale

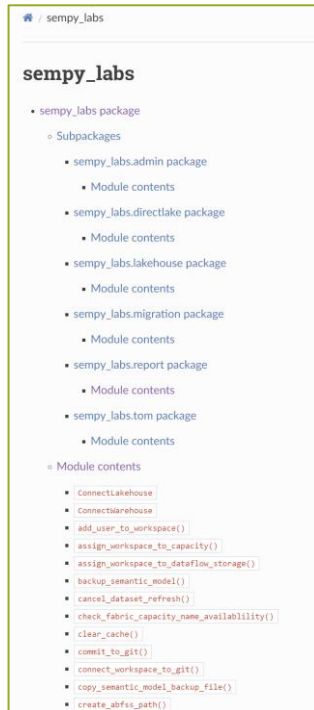
Monitoring your Fabric landscape

Quality Gates at scale

- <https://github.com/microsoft/semantic-link-labs/tree/main/notebooks>

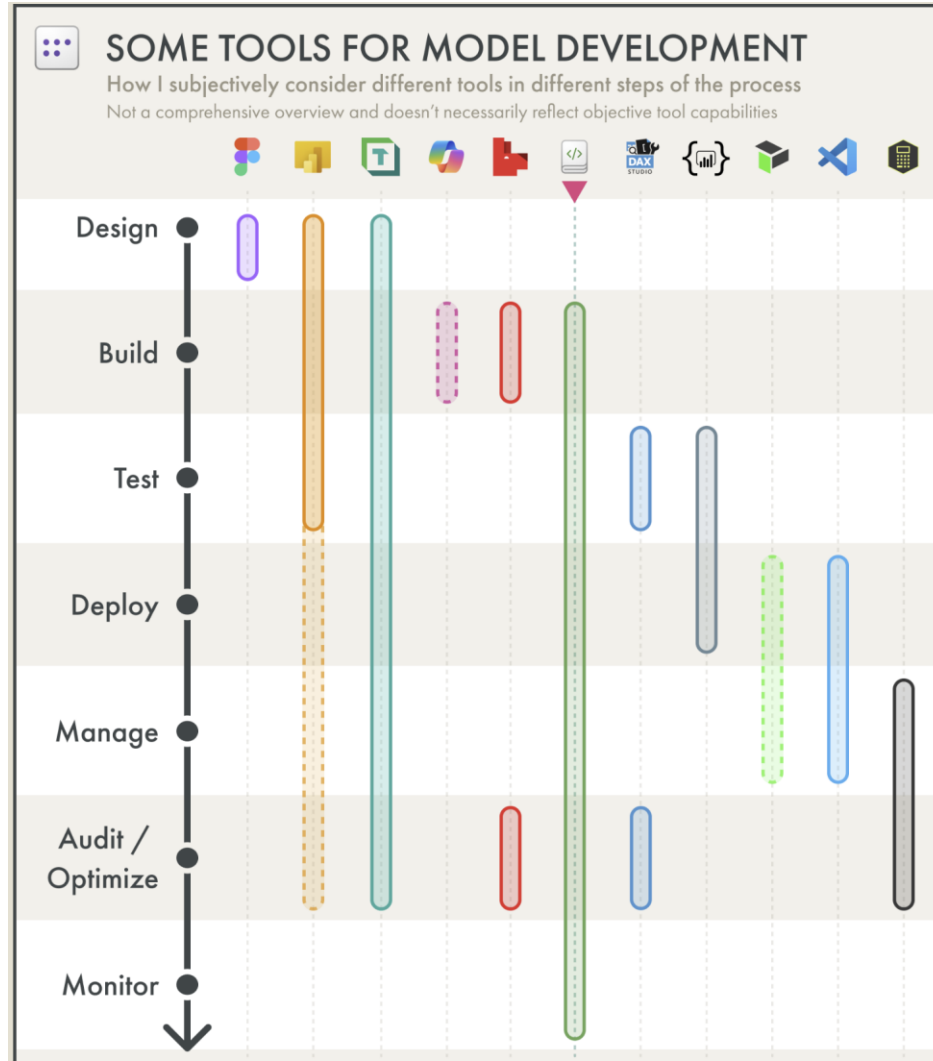
# Semantic Link Labs

- Innovator: Michael Kovalsky, Microsoft CAT team
- Library has > 300 and more functions
- You don't need to be a python expert
- You don't need to deal with complex API
- The swiss knife for all Fabricators



# Semantic Link [-Labs]

- If you are a “Semantic”  
– Data Engineer or modeler
- Kurt Buhler has an excellent  
blog on tools for model  
development  
<https://data-goblins.com/power-bi/semantic-link-labs>
- The only place where notebooks  
can't support you is design



# Simplifying data preparation with Data Wrangler





# Simplifying data preparation with Data Wrangler

- Notebook-based tool for exploratory data analysis
- Code Accelerator

The screenshot displays the Data Wrangler interface, which is a notebook-based tool for exploratory data analysis. The interface is divided into several sections:

- Notebook:** Contains a code cell with the following Python code:


```
1 df = spark.read.format("csv").option("header", "true").load("Files/sample_datasets/ImmoScout24_CH_Rental_Property_Dataset.csv")
2 # df now is a Spark DataFrame containing CSV data from "Files/sample_datasets/ImmoScout24_CH_Rental_Property_Dataset.csv".
3 display(df)
```

The code cell shows a success message: "4 sec - Command executed in 2 sec 482 ms by PySpark (Python)".
- Data Preview:** A table showing the first 1000 rows of the dataset. The columns are: ID, type, address, city, postcode, lat, lon, rentNet, rentGross, currency, and livingSpace. The table is titled "Data Wrangler: df" and shows a summary of the data.
- Summary Panel:** A panel on the right side of the interface showing the data shape (1,000 rows x 51 columns) and a list of columns. It also shows the number of rows with missing values (0) and the number of missing values (23,259).

A large green arrow points from the code cell to the Data Preview table, indicating the flow of data from the notebook to the visualization.

# Simplifying data preparation with Data Wrangler

- Works with Pandas **and** Spark DataFrames
- Automatically converts Spark DataFrames to pandas samples for performance reason

 Your Spark DataFrame has been converted into a pandas sample for performance but all the code generated by Data Wrangler will be converted to PySpark when you add it to your notebook.

- Custom Sample possibility

## Launch Data Wrangler

DataFrame \*

df

Sample method and size \*

Random

5000

rows

Open

Cancel

# Simplifying data preparation with Data Wrangler

- No Code Data Prep Operations for your DataFrame
- Code will be generated in the background and can be applied afterwards to your notebook
- Code is generated as function to be easily reused and operated

## Add code to notebook

The code generated by Data Wrangler will be converted to PySpark when you add it to your notebook.

☐ Include pandas code

☒ Preview PySpark code

```
# Code generated by Data Wrangler for PySpark DataFrame
from pyspark.sql import functions as F
from pyspark.sql import types as T

def clean_data(df):
    # Drop column: 'ID'
    df = df.drop('ID')
    # Derive column 'NewType' from column: 'type'
    # Transform based on the following examples:
```

Add Cancel

Operations

Search for operations...

- > Find and replace (4)
- > Format (7)
- > Formulas (3)
- > Numeric (4)
- > Schema (5)
- > Sort and filter (2)
- Group by and aggregate
- New column by example

Operations

← Back to current operation

New column by example

Automatically create a column when a pattern is detected from the examples you provide. Powered by Microsoft Flash Fill.

Target columns

type

Derived column name

NewType

Update Cancel

Cleaning steps

1 Load data from variable

4 New operation

#	index	type	NewType	address	city
		Missing: 0 (0%) Distinct: 19 (2%)	Missing: 0 (0%) Distinct: 4 (<1%)	Missing: 116 (12%) Distinct: 766 (77%)	Missing: 0 (0%) Distinct: 449 (45%)
		[APARTMENT, 'FLAT']: 65% [APARTMENT]: 13% [APARTMENT, 'DUPLEX']: 5% Other: 17%	Apartment: 96% House: 3% Hobby_Room: <7% Other: <1%	766 Distinct values	449 Distinct values
0	[HOUSE, 'SINGLE_HOUSE']	House	Altenstrasse 4	Filzbach	
1	[APARTMENT, 'STUDIO']	Apartment	Hauptstr. 30 (Nord)	Döttingen	
2	[APARTMENT, 'FLAT']	Apartment	Falmenstrasse 2d	Uster	
3	[APARTMENT, 'DUPLEX']	Apartment	Am Mattenhof 2b	Kriens	
4	[APARTMENT]	Apartment	Dorfstrasse 31	Benzenschwil	
5	[APARTMENT, 'FLAT']	Apartment	Preyenstr. 21	Wetzikon ZH	
6	[SINGLE_HOUSE]	Single_House	Beggengrassstrasse 102	Schliethelm	
7	[APARTMENT, 'LOFT']	Apartment	alte Spinnerei 16	Murg	
8	[APARTMENT, 'FLAT']	Apartment	Allmannstrasse 55, 8052 Zürich	Zürich	
9	[APARTMENT, 'FLAT']	Apartment	Unterer Rollring 20	Hägendorf	
10	[APARTMENT]	Apartment	Winzerhalde 109	Zürich	
11	[APARTMENT, 'SINGLE_ROOM']	Apartment	Kreuzstrasse 11	Bülach	
12	[APARTMENT, 'FLAT']	Apartment	Hauptstr. 2	Biberist	
13	[APARTMENT, 'DUPLEX']	Apartment	Untere Kräzern 2b	St. Gallen	

## 3 New column by example

```
1 # Derive column 'NewType' from column: 'type'
2 # Transform based on the following examples:
3 #   type      Output
4 # 1: "[HOUSE, 'SINGLE_HOUSE']" => "House"
5 pandas_df.insert(1, "NewType", pandas_df.apply(lambda row: row["type"].split("'")[1].title(), axis=1))
```

# Supercharging your Spark workflow with integrated tools



# Working with VSCode

- 2 different MS Fabric extensions are available
- Data Wrangler also exists as a standalone VSCode extension
- In addition to Microsoft extensions, extensions from the community are also available
- To best use the add-in locally, devcontainers are used  
[devcontainers \(github.com\)](https://github.com/devcontainers)
- Container image must be built locally
  - [VS Code extension with Docker support - Microsoft Fabric | Microsoft Learn](#)
  - [SynapseVSCode/samples/.devcontainer at main · microsoft/SynapseVSCode \(github.com\)](#)
- VSCode on the web from Fabric currently with version bugs



SQL  
KONFERENZ  
2025

# Working with VSCode



Synapse VSCode

[Synapse VS Code - Visual Studio Marketplace](#)



Synapse VS Code - Remote

[Synapse VS Code - Remote - Visual Studio Marketplace](#)



Fabric Studio

[Fabric Studio - Visual Studio Marketplace](#)



OneLake VSCode

[OneLake VSCode - Visual Studio Marketplace](#)



Data Wrangler

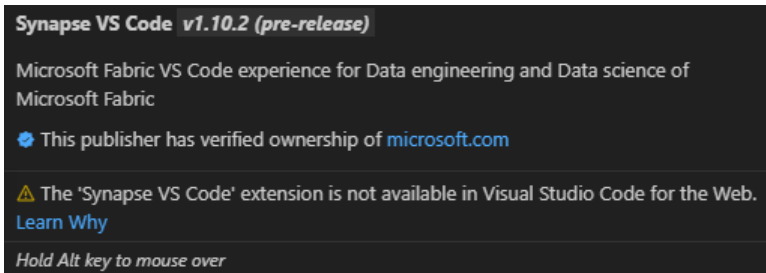
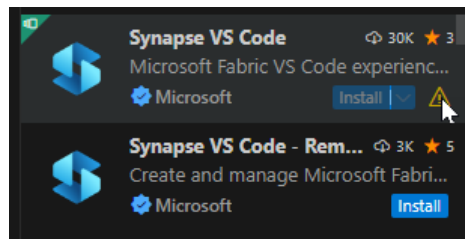
[Data Wrangler - Visual Studio Marketplace](#)



SQL  
KONFERENZ  
2025

# Working with VSCode

- Different versions for the desktop and web versions
- Synapse VS Code – Remote is intended as the extension for the web



SQL  
KONFERENZ  
2025



# Working with VSCode

The screenshot displays the VS Code interface with the Extensions Marketplace open. The left sidebar shows the 'EXTENSIONS' view with a search bar and a list of installed and popular extensions. The 'Jupyter' extension by Microsoft is highlighted. The main panel shows the details for the 'Jupyter' extension (v2024.9.0, Built-in). A red circle highlights the compatibility warning: 'Extension is not compatible with Code 1.93.1. Extension requires: ^1.94.0.' Below this, the 'Extension Pack (4)' section lists four related extensions: 'Jupyter Keymap', 'Jupyter Notebook Renderers', 'Jupyter Slide Show', and 'Jupyter Cell Tags'. The bottom of the page features the text 'Jupyter Extension for Visual Studio Code'.

**EXTENSIONS**

Search Extensions in Marketplace

**INSTALLED** 0

**POPULAR**

- Python**  
Python language support with exte...  
Microsoft
- Pylance**  
A performant, feature-rich languag...  
Microsoft
- Jupyter**  
Jupyter notebook support, interacti...  
Microsoft (Incompatible)
- Jupyter Keymap**  
Jupyter keymaps for notebooks  
Microsoft
- Jupyter Notebook Renderers**  
Renderers for Jupyter Notebooks (w...  
Microsoft
- Prettier - Code for...** 49.1M ★ 3.5  
Code formatter using prettier  
Prettier

**Extension: Jupyter**

**Jupyter** v2024.9.0 Built-in  
Microsoft | microsoft.com | 82,669,550 | ★★★★★ (323)  
Jupyter notebook support, interactive programming and computing that supports Intellisense, debugging and more.

Extension is not compatible with Code 1.93.1. Extension requires: ^1.94.0.

**DETAILS** FEATURES CHANGELOG

**Extension Pack (4)**

- Jupyter Keymap**  
Jupyter Keymaps for notebooks  
Microsoft
- Jupyter Notebook Renderers**  
Renderers for Jupyter Notebooks (with plotly...  
Microsoft
- Jupyter Slide Show**  
Jupyter Slide Show support for VS Code  
Microsoft
- Jupyter Cell Tags**  
Jupyter Cell Tags support for VS Code  
Microsoft

**Jupyter Extension for Visual Studio Code**





□

e2c2dcf58b18 

latest







In use

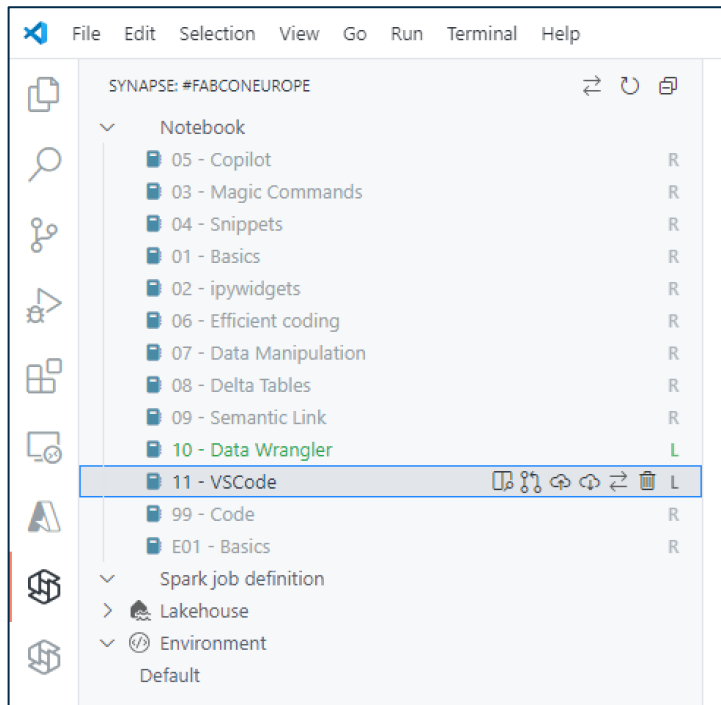
12 minutes ago

25.26 GB



# Working with VSCode

-  Open Notebook Folder
-  Update Notebook
-  Publish Resource Folder
-  Update Resource Folder
-  Switch Notebook Environment
-  Delete Notebook



## Links

- GitHub Repository  
[TEitelberg/MasteringSparkNotebooks](https://github.com/TEitelberg/MasteringSparkNotebooks)
- Semantic Link Labs  
[microsoft/semantic-link-labs](https://github.com/microsoft/semantic-link-labs)



## Session Feedback



**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025







**SQL**  
KONFERENZ  
2025

**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025





**SQL**  
KONFERENZ  
2025

**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025





**SQL**  
KONFERENZ  
2025



**SQL**  
KONFERENZ  
2025





**SQL**  
KONFERENZ  
2025

**SQL**  
KONFERENZ  
2025

