

# Enterprise-Ready Python for High-Performance Data Teams

Savannah (Ostrowski) Bailey, Doris Lee

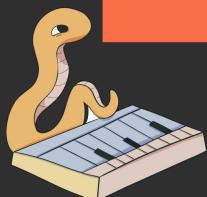


# Save the slides!

Slides for this talk can be found on GitHub



[savannahostrowski/pycon-2025](https://github.com/savannahostrowski/pycon-2025)



# Savannah (Ostrowski) Bailey



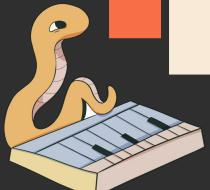
- 🐍 Python Core Developer
- 🚀 Faster CPython team
- 🧪 Jupyter Foundation Governing Board
- ❤️ 10+ years writing Python (3rd in-person PyCon!)
- 🔧 Prev: Pylance/Python in VS Code, Docker, Azure Developer CLI
- 💻 Now: Product Lead for Notebooks and Python



# Doris Lee

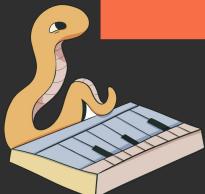


- 🎯 Build intelligent, easy-to-use tools for data science
- 🔭 Started career in astronomy (undergrad)
- 🐻 CS PhD @ UC Berkeley  
(Research in data viz, notebooks, Python - developed OSS projects)
- 🚀 Co-founder & CEO of Ponder (company behind Modin)
- 🤝 Acquired by Snowflake in 2023
- ❄️ Now: Building Python Product @ Snowflake

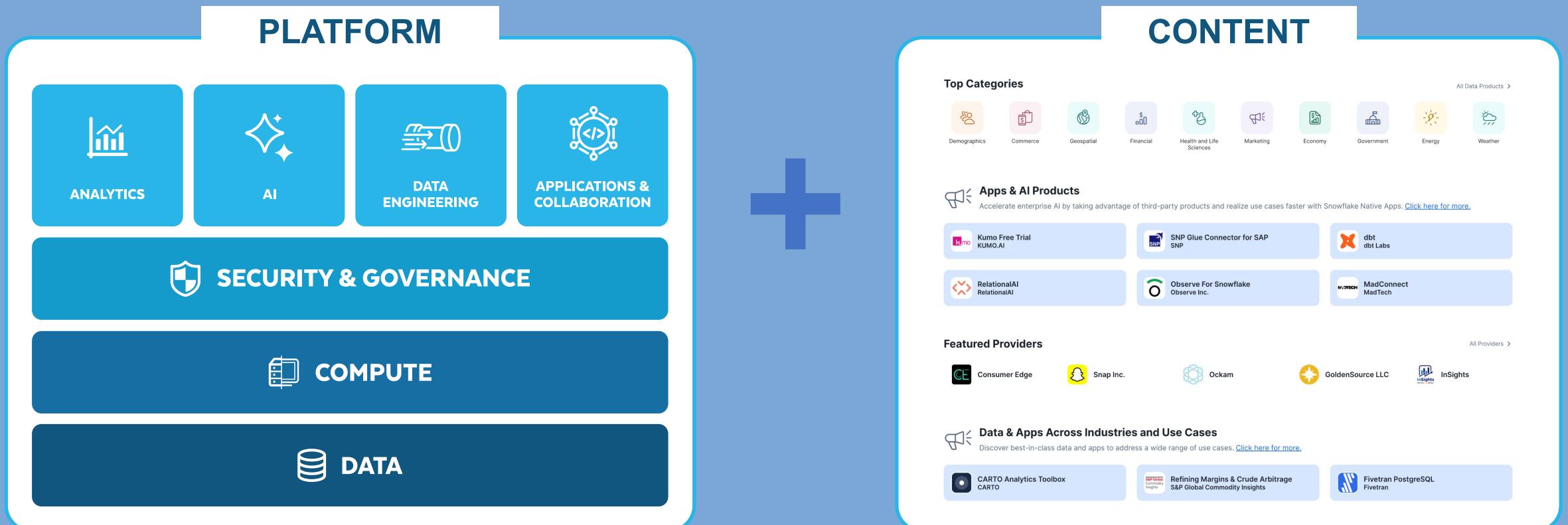


# Agenda

- Snowflake + PyData
- Lifecycle of data to model for enterprise team
- Enterprise-ready PyData APIs
- Enterprise-grade Runtime and Environments
- Best practices for PyData at Scale

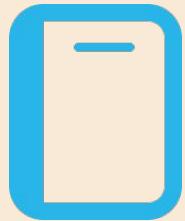


# Snowflake: Mobilize Data, Apps, and AI



# Python in Snowflake

Any data app or workload, all on Snowflake



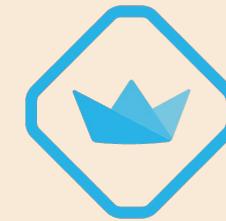
## Snowflake Notebooks

Explore, analyze and visualize data using Python and SQL



## Snowpark Python

Query and transform data with familiar pandas or Spark API



## Streamlit in Snowflake

Turn data and AI into interactive apps with Python

# PyData 🤝 Snowflake: Embracing Open Standards

Standards you know and love from PyData Ecosystem



IDE

From Snowflake Notebooks or any IDE



Visual Studio Code



Streamlit



APIs

Query and transform data with familiar APIs



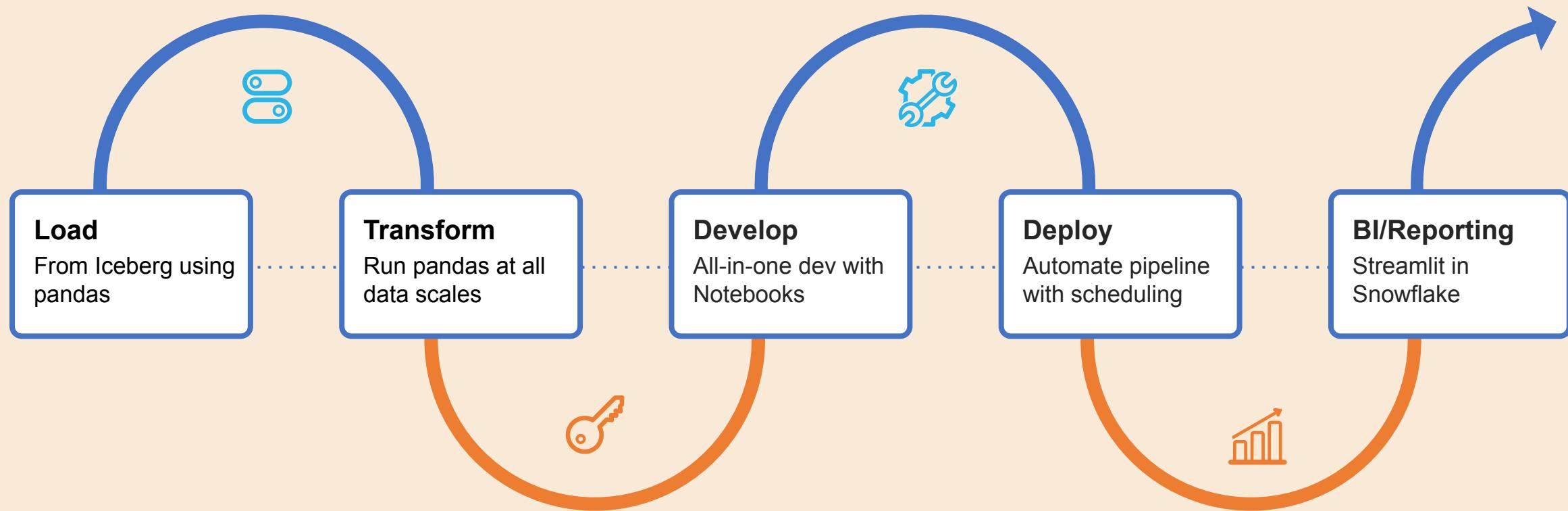
Packages

Bring your own packages and libraries



# From data to model

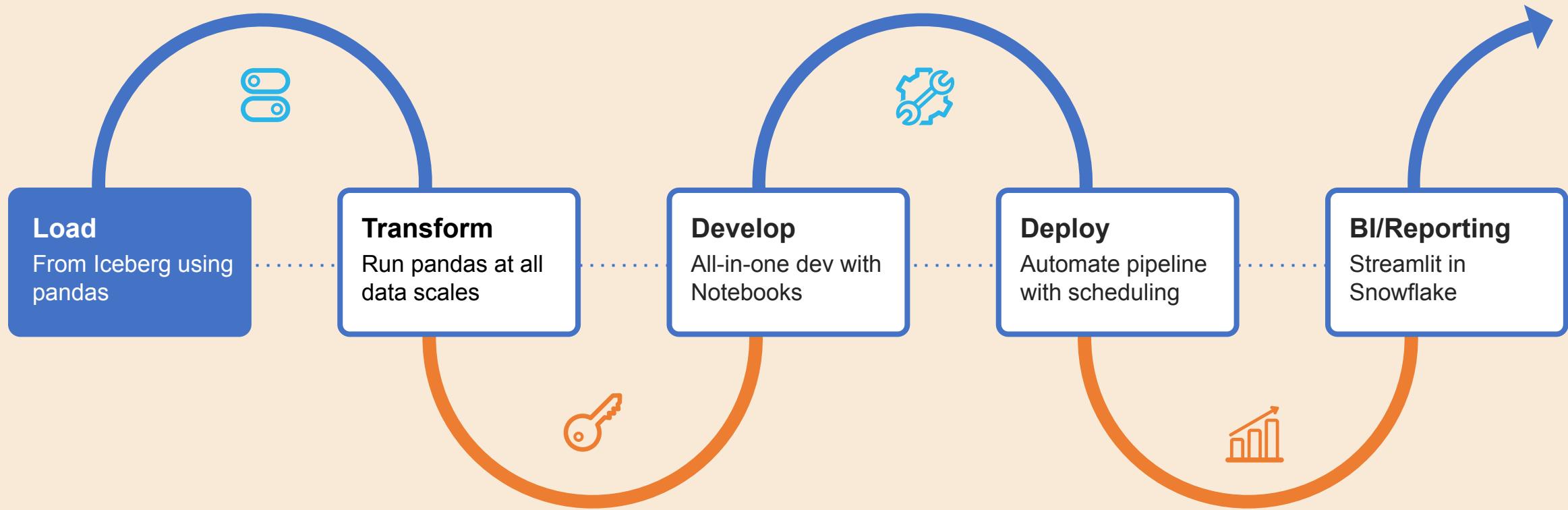
End-to-end development for high performance data teams



# Enterprise-ready PyData APIs

# From data to model

Load data from any source to work with using pandas



# What is pandas ?

The de-facto Python library for data cleaning, transformation, and analysis



#1 most popular Python library  
(used by 1 in 5 developers)



Flexible (600+ APIs!)



Quick Prototyping



Easy-to-use



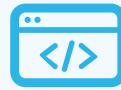
Widely adopted



`df.describe(...)`



`df.rename(...)`



`pd.concat(...)`



`df.dropna(...)`

`df.groupby(...)`



`pd.join(...)`

`pd.merge(...)`

`df.pivot(...)`

# Bringing best of pandas to Snowflake

## PANDAS ON SNOWFLAKE

### WHAT IS IT

Simple, unified interface for working with Snowflake data using familiar pandas API

### WHY USE IT

- Easily work with Snowflake data and bring in data from files
- Develop pandas pipeline at all data scales from prototype to production
- Powerful Snowflake analytics through convenient pandas semantics for flexible transformation

```
import modin.pandas as pd
import snowflake.snowpark.modin.plugin
# Load data from a Snowflake table (Supports Iceberg, views, etc.)
df = pd.read_snowflake('orders')

# Convert date column to datetime type for resampling
df['date'] = pd.to_datetime(df['date'])

# Classify text using Cortex LLM functions
from snowflake.cortex import ClassifyText
df['region'] = df["location"].apply(ClassifyText, /
                                      categories=['North America', 'Europe', 'Asia'])

# Resample to daily granularity and fill missing days with zero
daily_sales = df.resample('D').sum()
daily_sales_filled = daily_sales.fillna(0)

# Save result to a dynamic table in Snowflake
df.to_dynamic_table('processed_orders')
```

# Loading data with pandas on Snowflake

Seamlessly interact with Snowflake objects

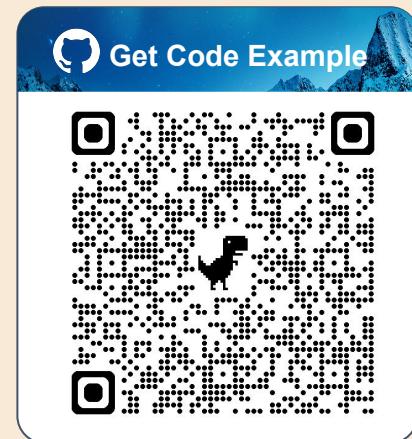
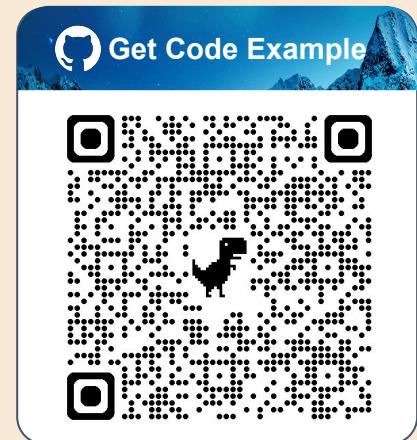
```
# Supports Tables, views, Iceberg, Dynamic Tables, etc.  
df = pd.read_snowflake("MY_TABLE")
```

<https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Read%20from%20Table/Reading%20from%20Table.ipynb>

Load files from local or stage paths across a range of formats

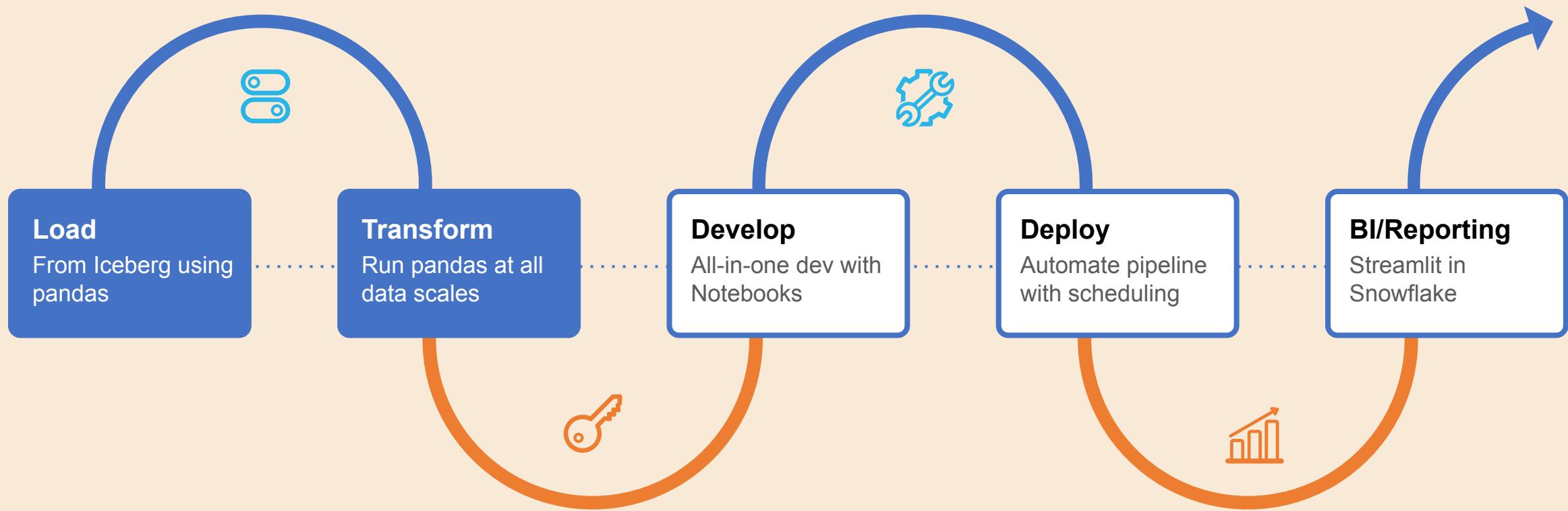
```
# Reading CSV from S3 and other remote source  
df = pd.read_csv("s3://mys3bucket/data.csv")  
# Reading Parquet from stage file  
df = pd.read_parquet("@MY_STAGE/data.parquet")
```

[https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Load%20files%20to%20Snowflake/load\\_files\\_snowflake.ipynb](https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Load%20files%20to%20Snowflake/load_files_snowflake.ipynb)



# From ingest to transformation

End-to-end development for high performance data teams



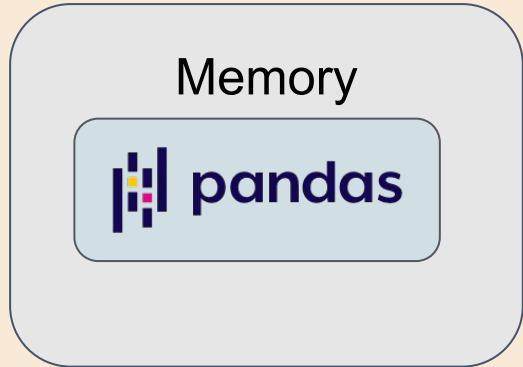
# Challenges of working with enterprise-scale data



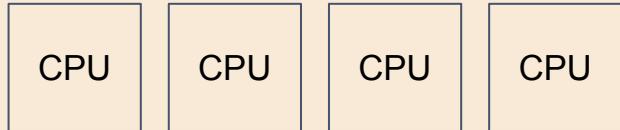
## 6+ months

Time spent to bring pandas workloads to production

# How pandas works?



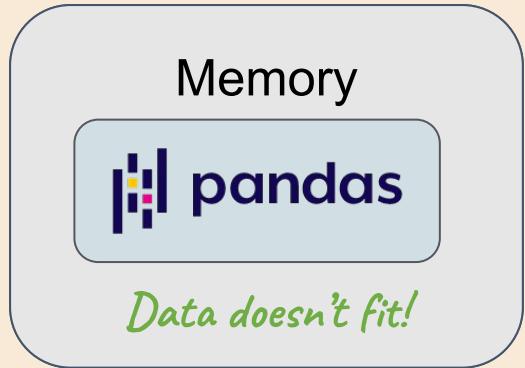
pandas operates on data **in-memory**



pandas is **single threaded**



# Enter Enterprise-scale data...



pandas operates on data **in-memory**



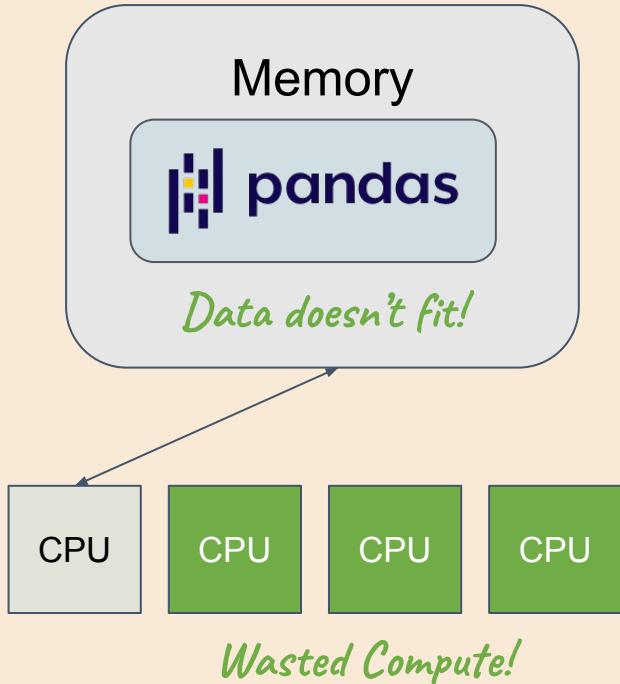
**Out of memory errors!**



pandas is **single threaded**



# Enter Enterprise-scale data...

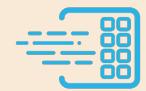


pandas operates on data **in-memory**

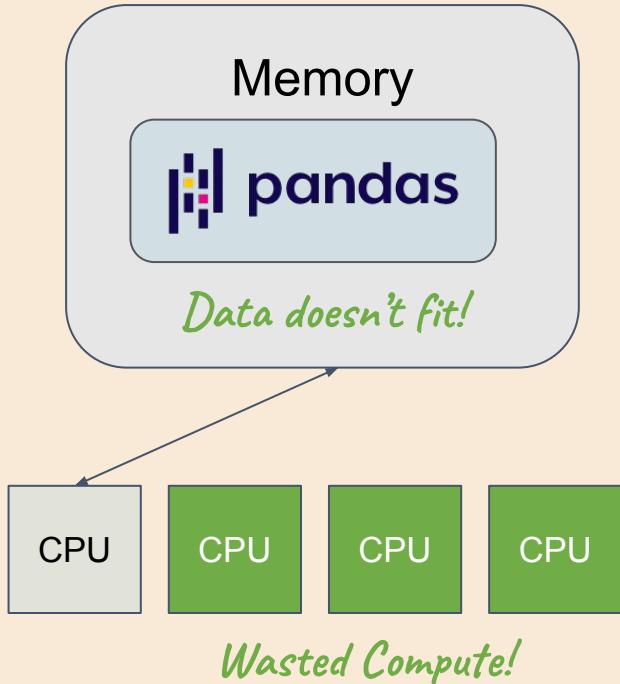
🚫 **Out of memory errors!**

pandas is **single threaded**

🚫 **Slow – No distributed processing!**



# Enter Enterprise-scale data...



pandas operates on data **in-memory**

🚫 **Out of memory errors!**

pandas is **single threaded**

🚫 **Slow – No distributed processing!**

...  **pandas** **doesn't scale** (even on ~ tens GB/millions rows)

# What this means for data teams

- ➡️ **Costly Re-translation**
- 🔧 **Difficult to experiment & debug**
- ⬇️ **Loss of productivity**





Scalable “drop-in” replacement for pandas

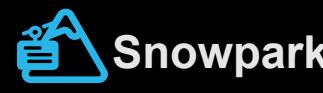
 modin\_project/modin

Star

10k+

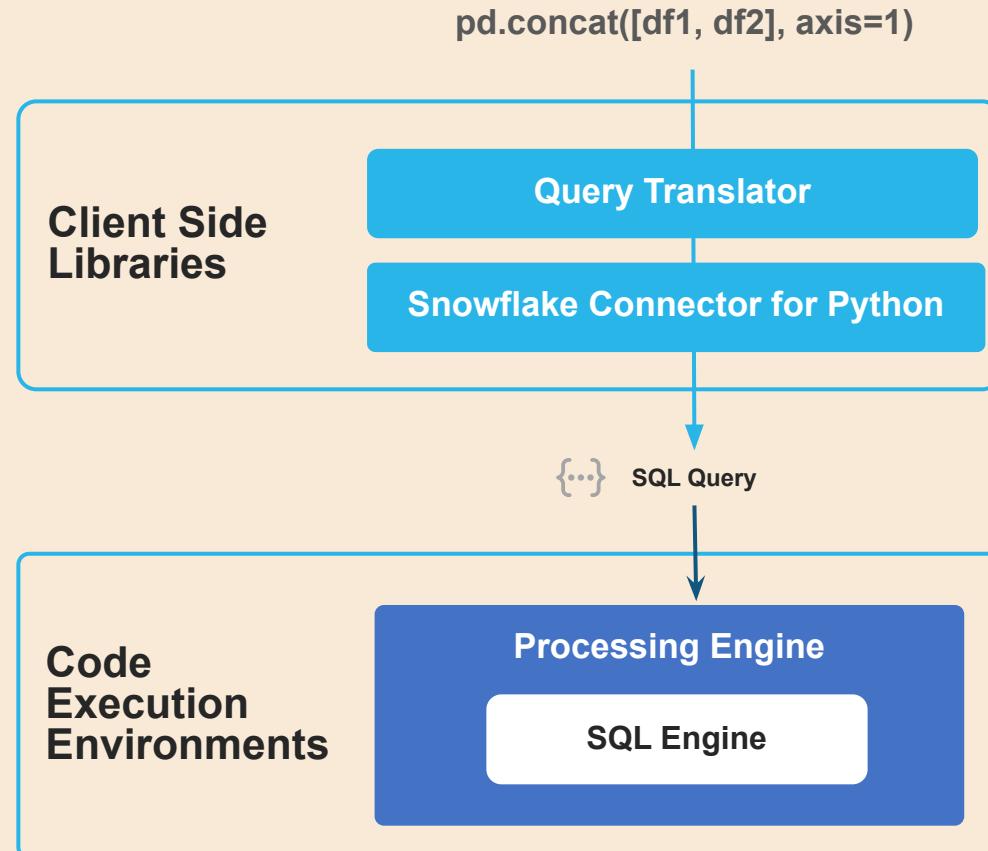
-  **1.3M+ monthly** downloads
-  **100+ contributors**
-  **Used by 1K+ organizations**

Runs on scalable OSS backends:



# Distributed Pandas at Scale in Snowflake

Automatic translation from pandas to SQL - no code rewrite required



## Seamless Prototype to Production

No need to manage or tune compute infrastructure



## Zero Data Movement

Compute happens directly in Snowflake securely

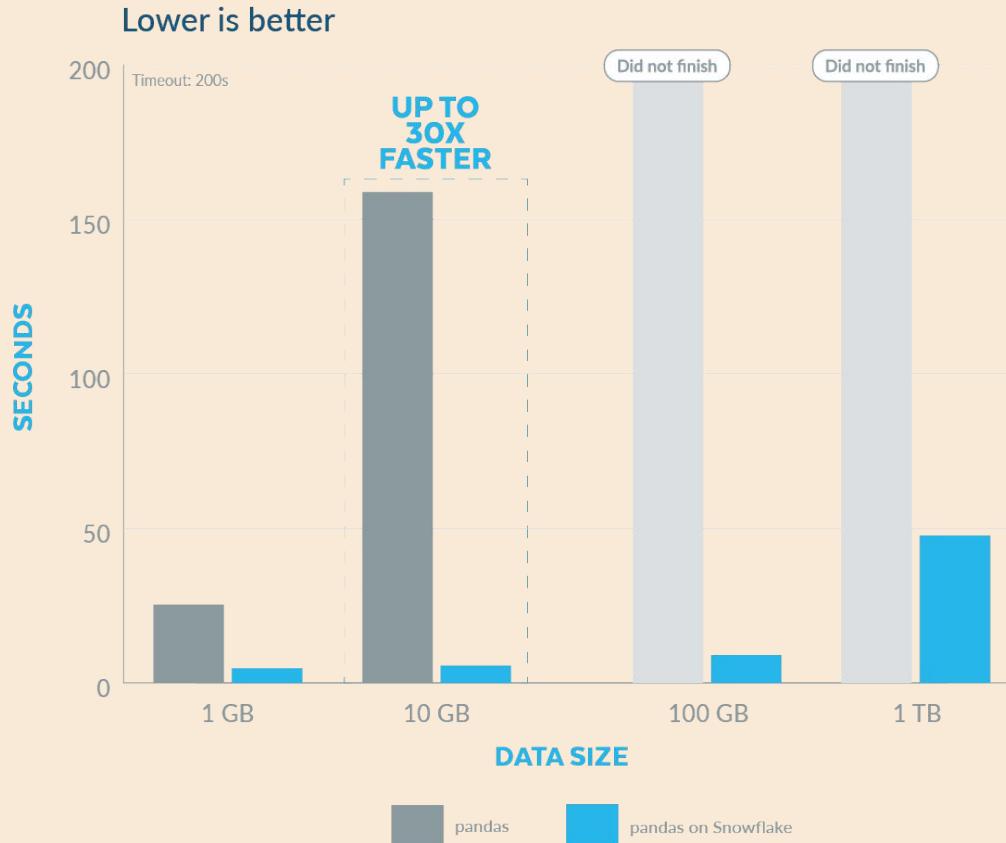


## Write with pandas

Leverage familiar API. No rewrites necessary.



# Run pandas on 1TB+ enterprise data



**30X+ faster performance  
than pandas**

- Address painful out of memory issues to scale pandas
- Secure access within Snowflake removes sensitive data risks on local machines
- Accelerate prototype to production and scale without tuning or managing infrastructure

Benchmark testing across different data sizes shows how pandas on Snowflake can scale up to 1TB and is up to 30x faster than pandas at the 10GB scale.

# Embracing Open Standards

Integration with PyData Libraries through Dataframe Interchange Protocol

- pandas is often used with third-party libraries for visualization, ML.
- Added support for dataframe interchange protocol
- Increase interoperability with other PyData libraries



# Get Started Today!

## Installation

Use with Snowflake Notebook or in your own IDE

```
$ pip install "snowflake-snowpark-python[modin]"
```

```
import modin.pandas as pd  
import snowflake.snowpark.modin.plugin
```

## Resources

### Developer Guide



<https://docs.snowflake.com/en/developer-guide/snowpark/python/pandas-on-snowflake>

### Quickstart



[https://quickstarts.snowflake.com/guide/getting\\_started\\_with\\_pandas\\_on\\_snowflake](https://quickstarts.snowflake.com/guide/getting_started_with_pandas_on_snowflake)

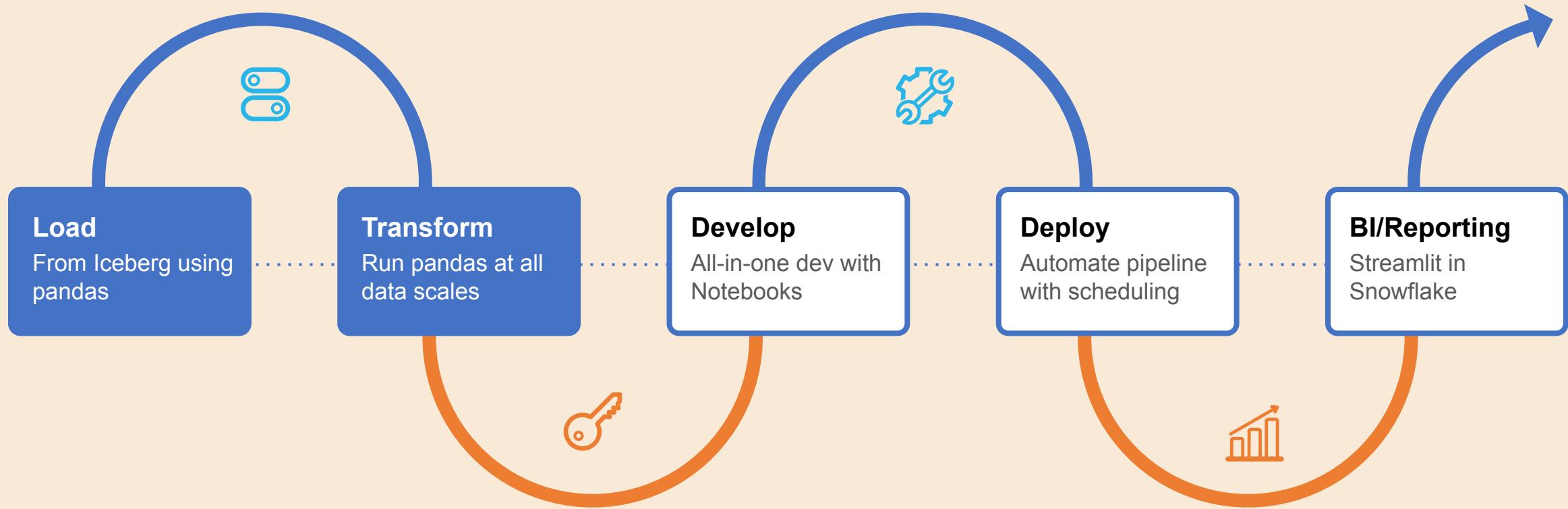


# Enterprise-grade runtimes and environments



# From exploration to model

End-to-end development for high performance data teams



# Snowflake Notebooks

The screenshot shows the Snowflake Notebook interface. At the top, there's a navigation bar with 'Notebooks' and 'My First Notebook Project'. Below it are buttons for 'Packages', 'Active', and 'Run all'. On the left, there's a sidebar with icons for file operations like 'New', 'Search', 'Copy', 'Edit', 'Delete', and 'Ask Copilot'. The main area has two code cells:

- Python v as cell3**:

```
1 # Import Python packages used in this notebook
2 import streamlit as st
3 import altair as alt
4
5 # Pre-installed libraries that comes with the notebook
6 import pandas as pd
7 import numpy as np
8
9 # Package that we just added
10 import matplotlib.pyplot as plt
```
- Markdown v as cell4**:

**SQL Querying at your fingertips 🎉**

We can easily switch between Python and SQL in the same worksheet.

Let's write some SQL to generate sample data to play with.
- SQL v as cell5**:

```
1 -- Generating a synthetic dataset of Snowboard products, along with their price and rating
2 v SELECT CONCAT('SNOW-',UNIFORM(1000,9999,RANDOM())) AS PRODUCT_ID,
3          ABS(NORMAL(5,.3,RANDOM())) AS RATING,
4          ABS(NORMAL(750,.200::FLOAT,RANDOM())) AS PRICE
5 FROM TABLE(GENERATOR(ROWCOUNT => 100));
```

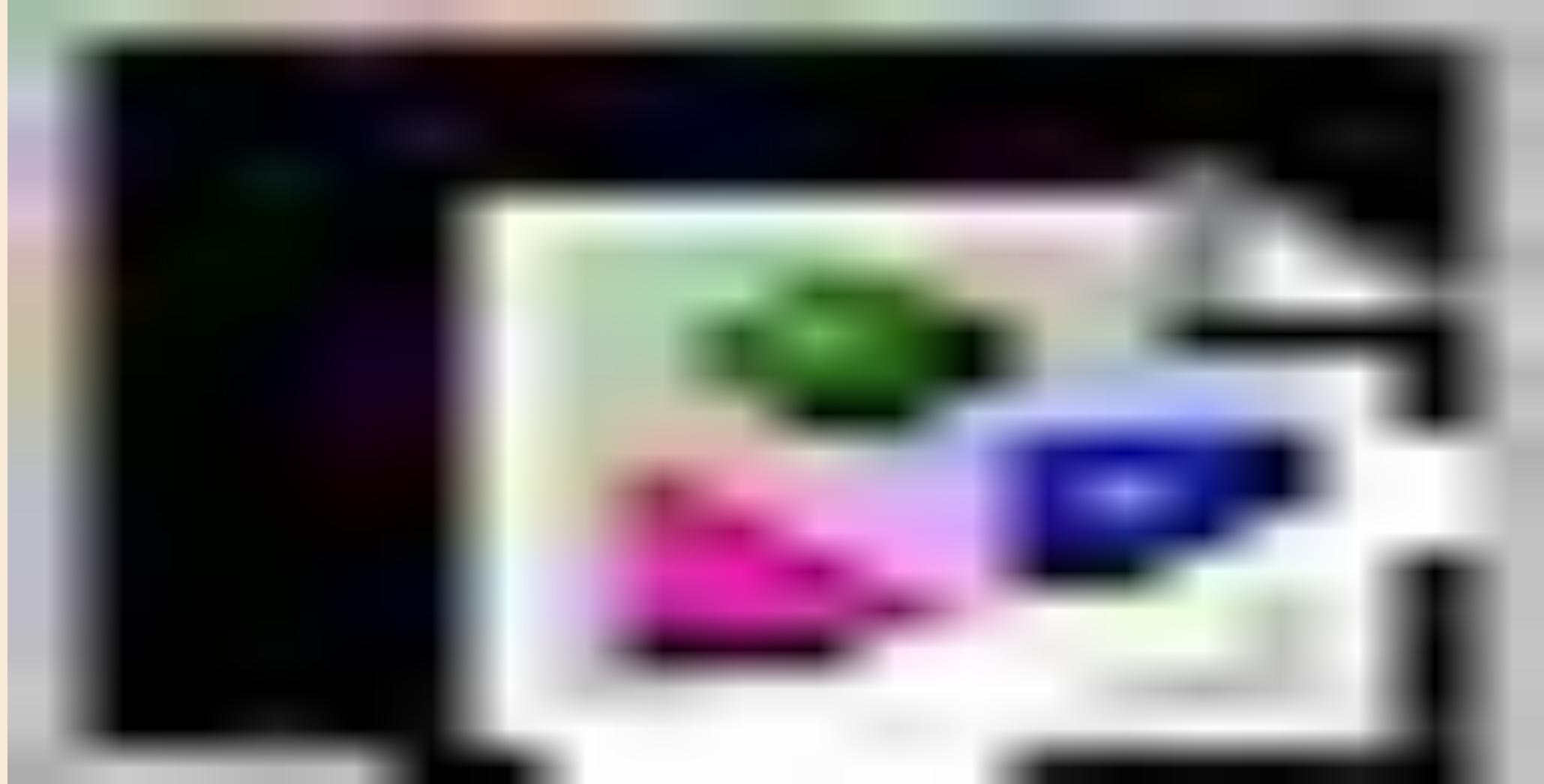
Below the SQL cell, there's a table output:

	PRODUCT_ID	RATING	PRICE
0	SNOW-9681	2.2102	879.2882
1	SNOW-9657	3.3833	835.3502
2	SNOW-3551	3.511	864.5094
3	SNOW-8226	0.9605	673.7927
4	SNOW-8473	3.1038	496.0473

A unified cell-based, interactive development surface that makes it faster and easier than ever to leverage your Snowflake data in Data Science, ML, and Data Engineering workflows

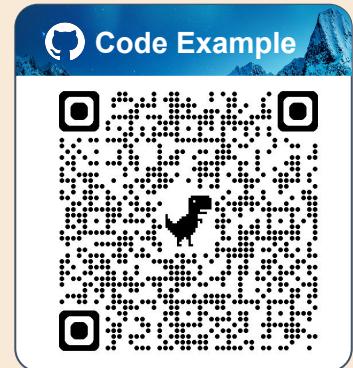
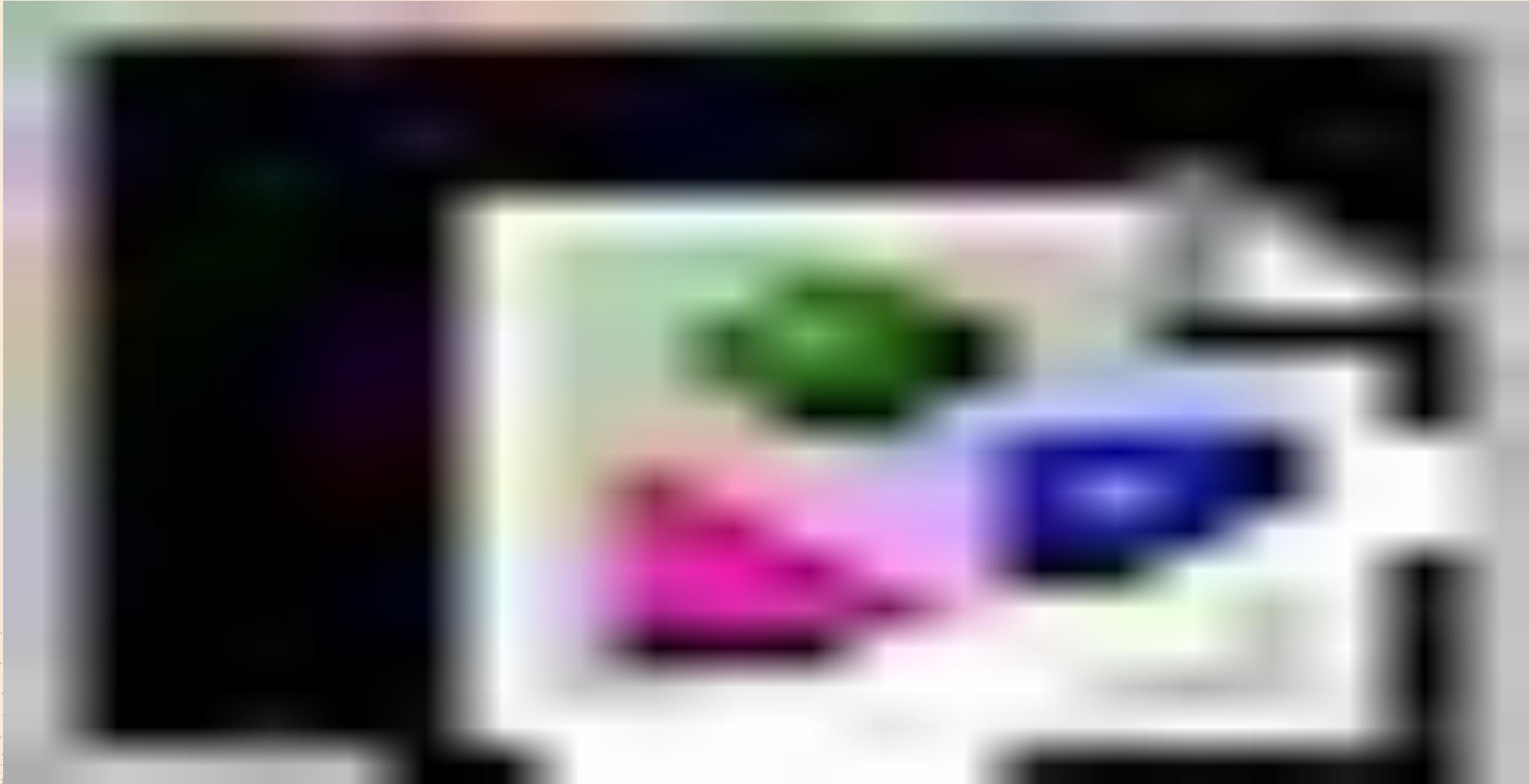
# Instant access to data from your Notebook

Focus on code, not configuring access to your data



# Seamless interop between languages

## Blend SQL and Python with cell and variable referencing



<https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Python%20SQL%20cells/python-sql-cells.ipynb>

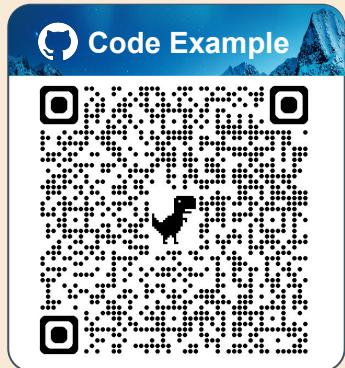
# Interactive data visualization in < 100 lines of code

# Embed interactive Streamlit apps directly in your Notebook for deeper insights

The screenshot shows a Jupyter Notebook interface. The top navigation bar includes 'Notebooks' and 'SAVANNAH\_PYCON\_CR' tabs, along with 'Packages', 'Share', and 'Active' buttons. The main area displays Python code in a cell labeled 'cell4'. The code imports Streamlit and Altair, converts a DataFrame to a pandas DataFrame, sets a title, creates an Altair chart with 'CARAT' on the x-axis, 'PRICE' on the y-axis, and 'CUT' as a color variable, and finally displays the chart using Streamlit's altair\_chart function. A sidebar on the right lists other cells in the notebook.

```
Python as cell4 • 0.0s ▶
```

```
1 import streamlit as st
2 import altair as alt
3
4 df = diamonds_df.to_pandas()
5
6 st.title("Diamond Carat vs Price Analysis")
7
8 chart = alt.Chart(df).mark_circle().encode(
9     x='CARAT',
10    y='PRICE',
11    color='CUT',
12    tooltip=['CARAT', 'PRICE', 'CUT', 'COLOR']
13 ).properties(
14     width=600,
15     height=400
16 )
17
18 st.altair_chart(chart, use_container_width=True)
```



[https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Data%20Visualization%20with%20Streamlit/data\\_viz\\_with\\_streamlit.ipynb](https://github.com/Snowflake-Labs/snowflake-python-recipes/blob/main/Data%20Visualization%20with%20Streamlit/data_viz_with_streamlit.ipynb)

# Challenges in production ML development



## Cost vs performance

Meeting the performance and quality standards while attempting to keep costs under control is challenging



## Scalability

With complex model schemas and large datasets, training models require sophisticated infrastructure capable of that scale



## Governance

Maintainability and reproducibility requires consistency in build environments and effective isolation – often a challenging DevOps problem

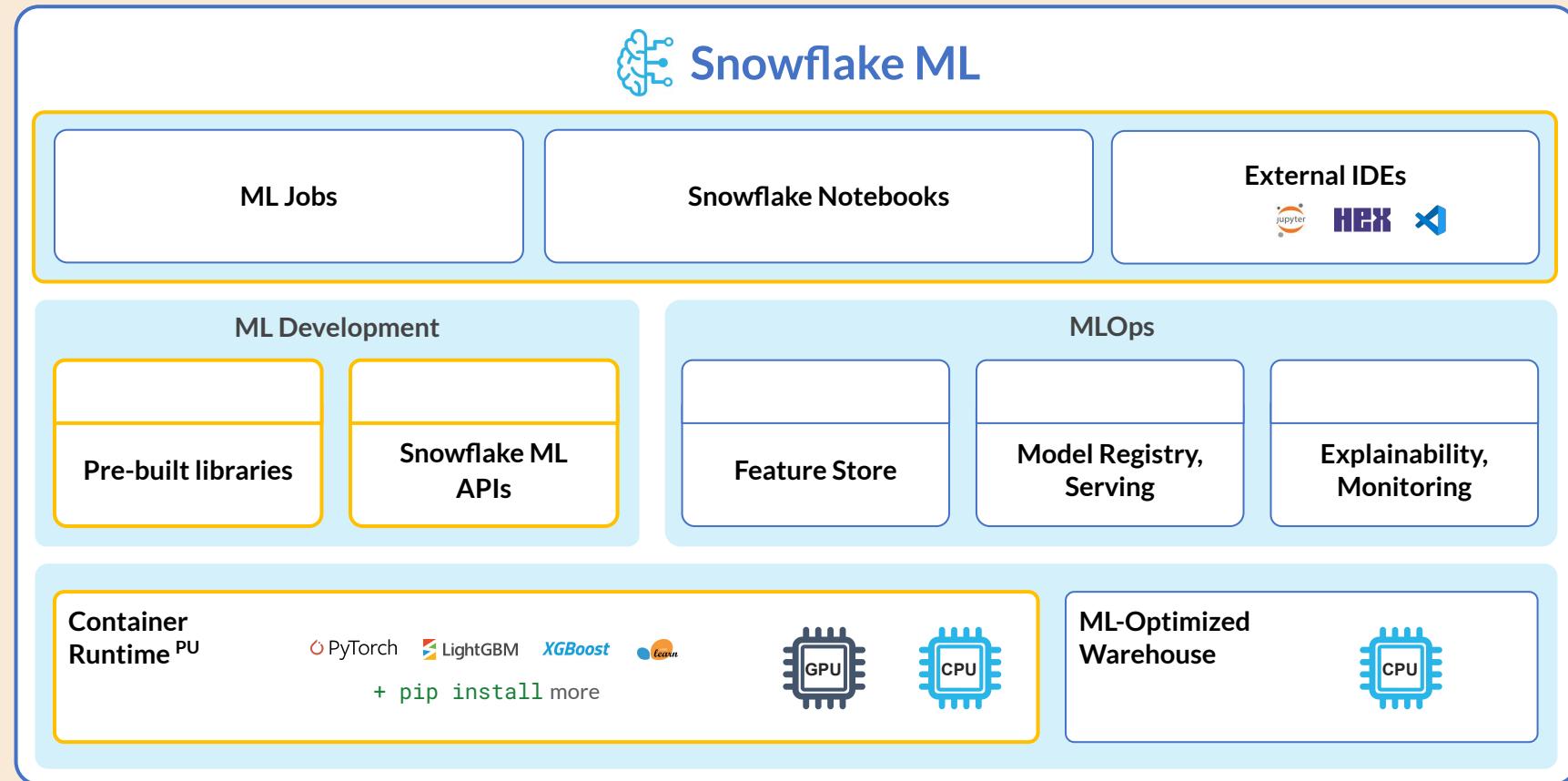


## Access to GPUs

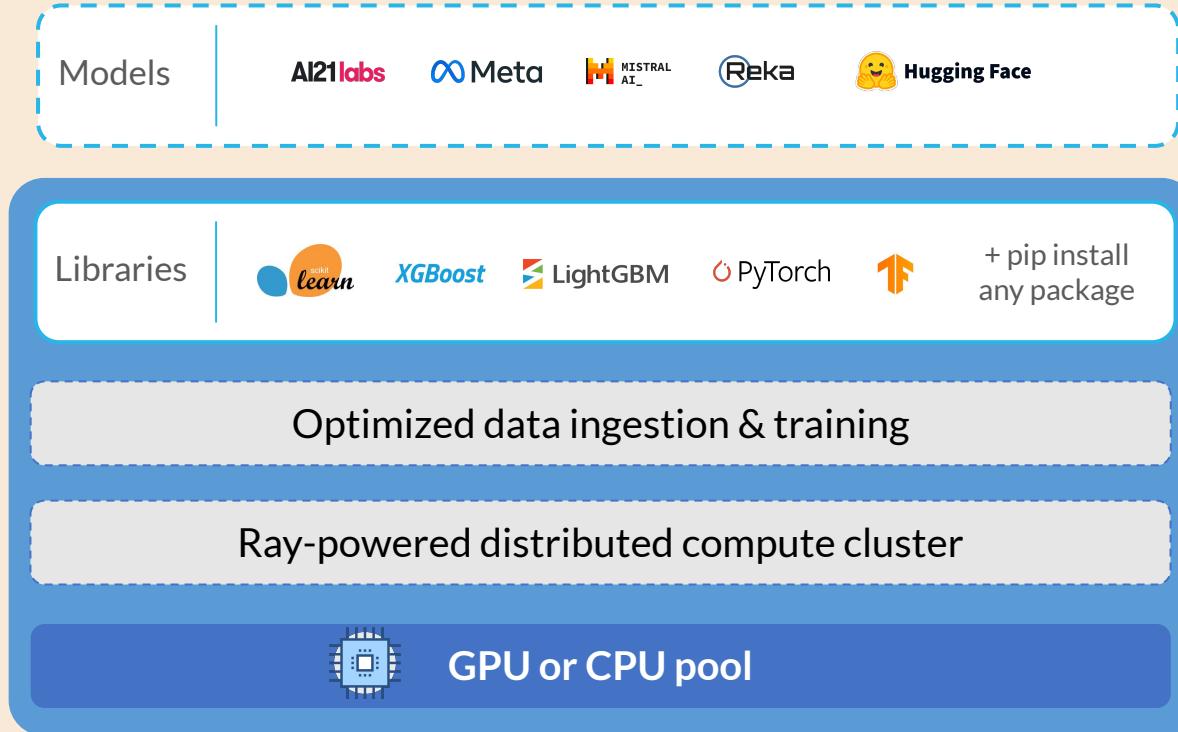
Deep learning and LLM related tasks require processing on GPUs – often hard to access without administrative setup or DevOps knowledge

# Snowflake ML

Integrated set of capabilities for end-to-end ML in Snowflake



# Container Runtime for ML

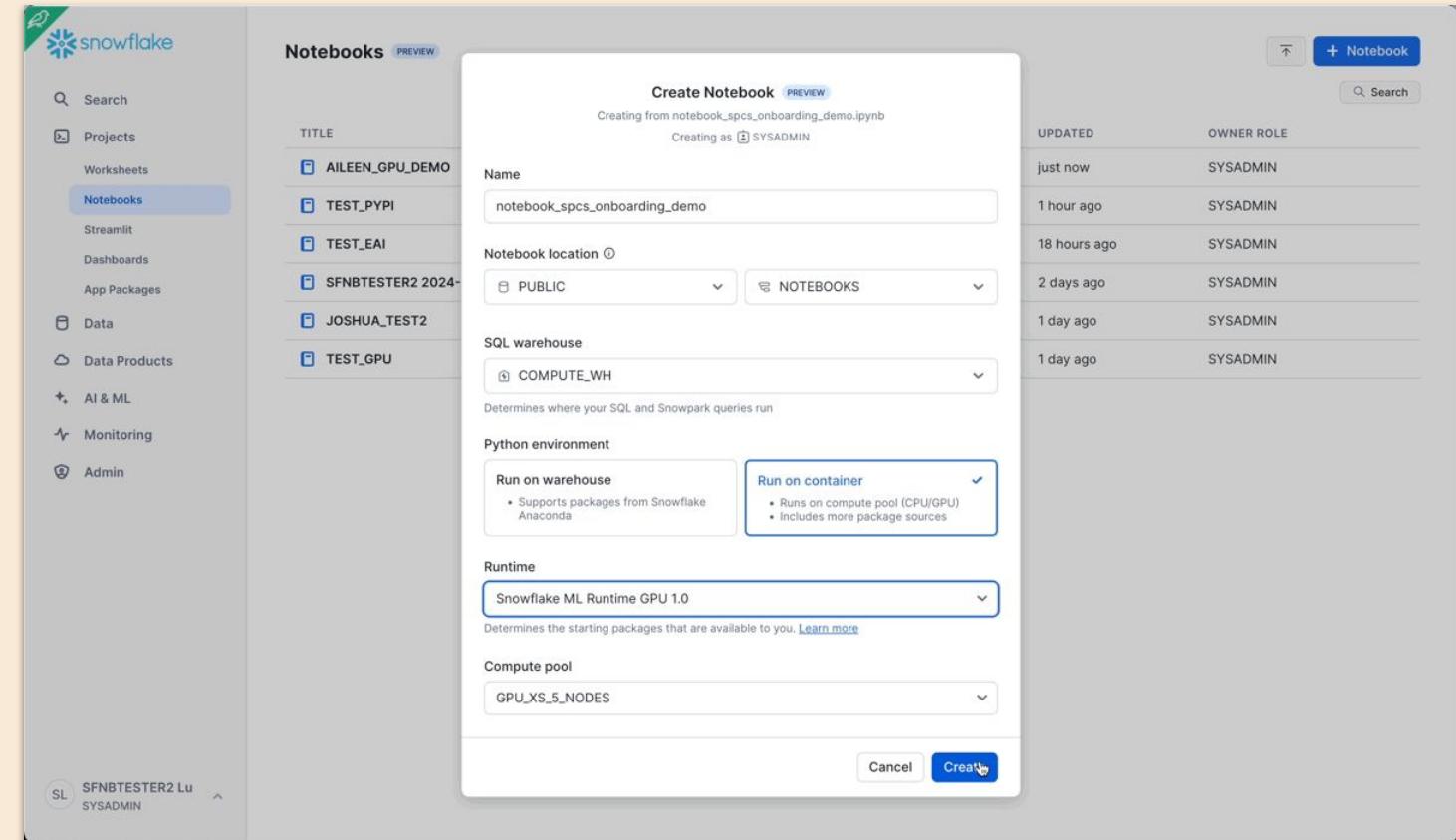


- Pre-built container environment with **latest ML packages** already installed
- Data connector APIs for **optimized data ingestion**
- Optimized and **scalable training APIs** for frameworks like XGBoost, LightGBM, PyTorch
- **Distributed hyperparameter optimization** using Bayesian strategy
- Ray-powered compute cluster for efficient multi-core and multi-node **distributed training**
- Easy access to **CPU and GPU** compute pools

PyTorch, the PyTorch logo and any related marks are trademarks of The Linux Foundation.  
TensorFlow, the TensorFlow logo and any related marks are trademarks of Google Inc.

# Snowflake Notebooks on Container Runtime

- Notebooks that run on a fully managed container environment with access to CPUs and GPUs, optimized data loading from Snowflake, automatic lineage capture and Model Registry integration
- Flexibility to run ML workloads at scale without data movement
- Leverage pre-installed ML packaged or pip to install any custom package



# Optimized ML APIs

## Data Loading with DataConnector API

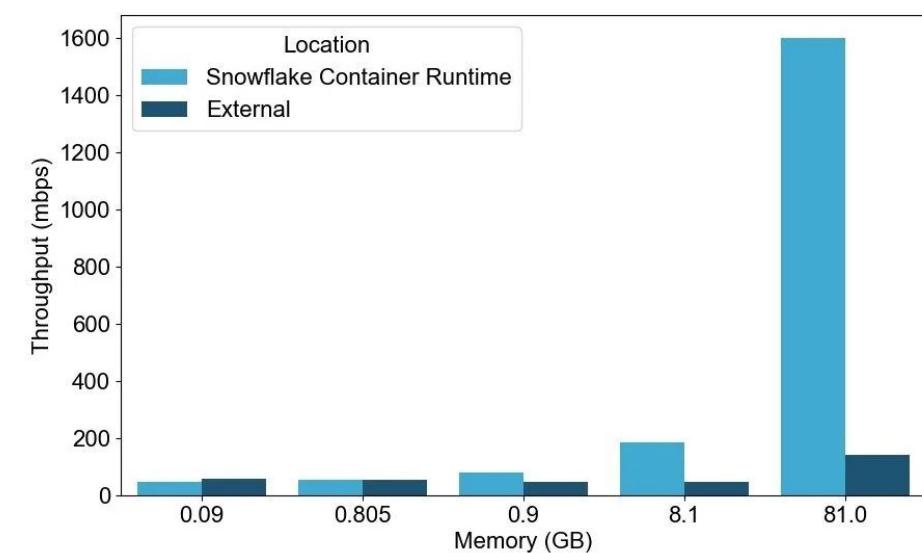
### With DataConnector API

*81 GB of data into pandas DataFrame in  
50 seconds vs. ~9.5 minutes with OSS  
equivalent outside Snowflake*

```
Python v as cell19

1  from snowflake.ml.data.data_connector import DataConnector
2
3  # Retrieve data from a snowflake table
4  table_name = 'MLR_BENCHMARK_DB.PUBLIC.REGRESSION_DATA_12GB' # 12GB table
5  snowpark_df = session.table(table_name)
6
7  # Materialized it into a pandas dataframe using DataConnector
8  t0 = time.time()
9  pandas_df = DataConnector.from_dataframe(snowpark_df).to_pandas()
10 duration = time.time() - t0
11 print(f'Materialization of 12GB data uses {duration:.2f} secs')

Materialization of 12GB data uses 7.36 secs
```



# Optimized ML APIs

## Distributed Training APIs

### With Snowflake ML - XGBoost

*3-7x speed up compared to open source XGBoost by utilizing all four available GPUs*

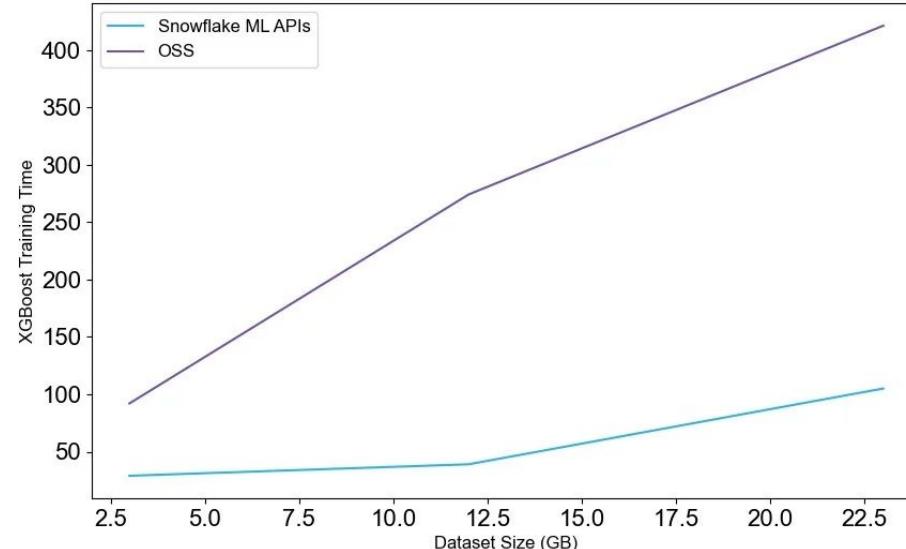
Distributed APIs for PyTorch, XGBoost & LightGBM available!

```
1  from snowflake.ml.modeling.distributors.xgboost import XGBEstimator, XGBScalingConfig
2  from snowflake.ml.data.data_connector import DataConnector
3
4  table_name = 'MLR_BENCHMARK_DB.PUBLIC.REGRESSION_DATA_12GB' # 12GB table
5  snowpark_df = session.table(table_name)
6  data = DataConnector.from_dataframe(snowpark_df)
7  input_cols = ["FEATURE1", "FEATURE2", "FEATURE3", "FEATURE4", "FEATURE5", "FEATURE6", "FEATURE7"]
8  label_col = 'TARGET'
9
10 # Leverage GPU, and let system figure out the rest by passing in -1
11 scaling_config = XGBScalingConfig(use_gpu=True, num_workers=-1, num_cpu_per_worker=-1)
12 estimator = XGBEstimator(n_estimators=10, scaling_config=scaling_config)
13
14 model = estimator.fit(data, input_cols=input_cols, label_col=label_col)
```

2024-09-10 23:24:23,436 INFO tune.py:622 -- [output] This will use the new output engine with verbosity

Training started without custom configuration.

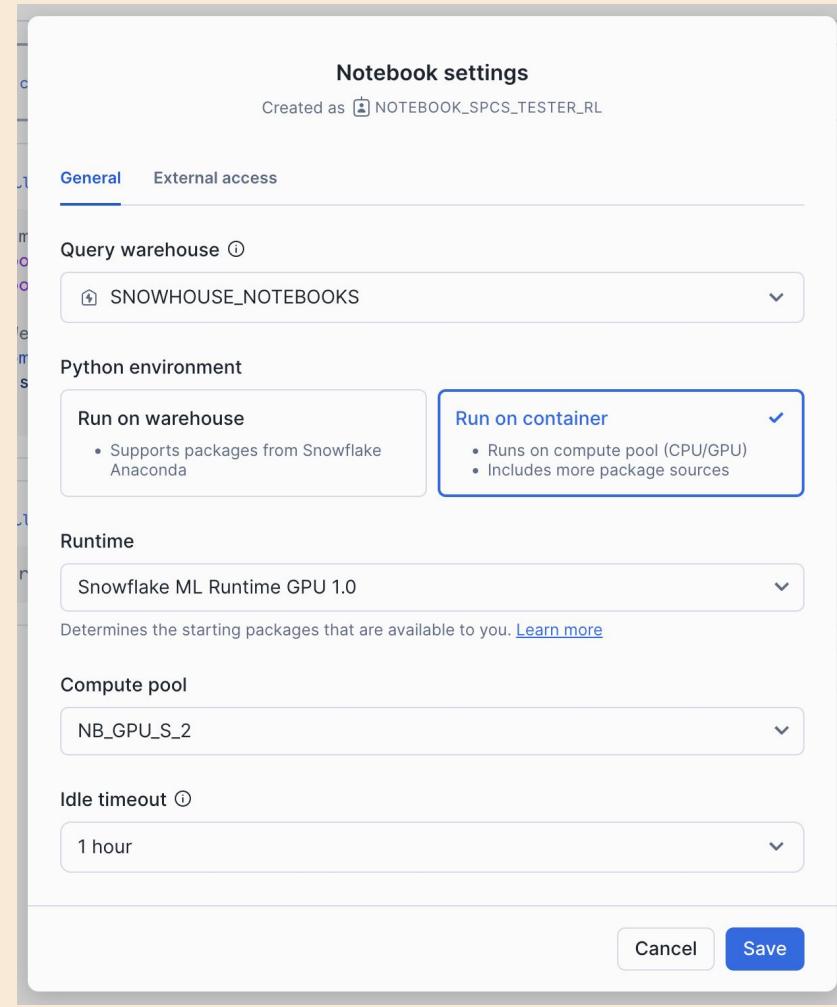
[36m(SplitCoordinator pid=6502) [0m Starting execution of Dataset. Full log is in /tmp/ray/session 2024



# Managed Access to GPUs

With admin-led governance

- ✓ Admin configured GPU pools easily accessible via Notebook config
- ✓ Easily scale up and down *during a NB session* to only use resources as needed
- ✓ Snowflake ML APIs configured to automatically detect and utilize all available GPUs



# Containerized ML Environment

## CPU and GPU-specific runtime images

- ✓ Comprehensive list of pre-installed packages curated for ML development
- ✓ Flexibility to supplement with any open source package with [pip](#)
- ✓ Pre-configured CUDA libraries to easily leverage GPUs

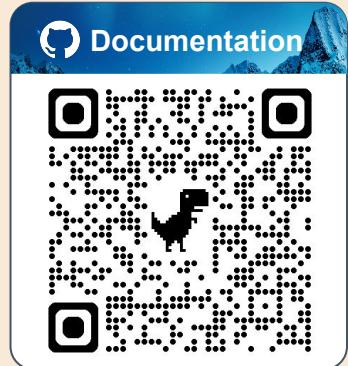
### Full list for GPU v1 image

This table is a full list of Python packages pre-installed on the GPU v1 image:

Package	Version	Package	Version
Babel	2.11.0	Bottleneck	1.3.7
Brotli	1.0.9	Brotli	1.0.9
GitPython	3.1.41	Jinja2	3.1.3
MarkupSafe	2.1.3	PyJWT	2.4.0
PySocks	1.7.1	PySocks	1.7.1
PyYAML	6.0.1	Pygments	2.15.1
Pympler	0.9	Send2Trash	1.8.2
absl-py	1.4.0	accelerate	0.28.0
aiobotocore	2.7.0	aiohttp	3.9.1
aiohttp	3.9.3	aiohttp-cors	0.7.0
aioitertools	0.7.1	aioprometheus	23.12.0
aiosignal	1.2.0	aiosignal	1.3.1
altair	4.2.2	annotated-types	0.6.0

# Run your Notebook on a schedule

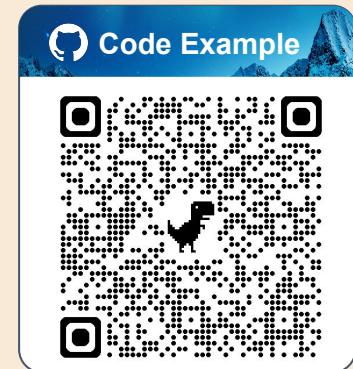
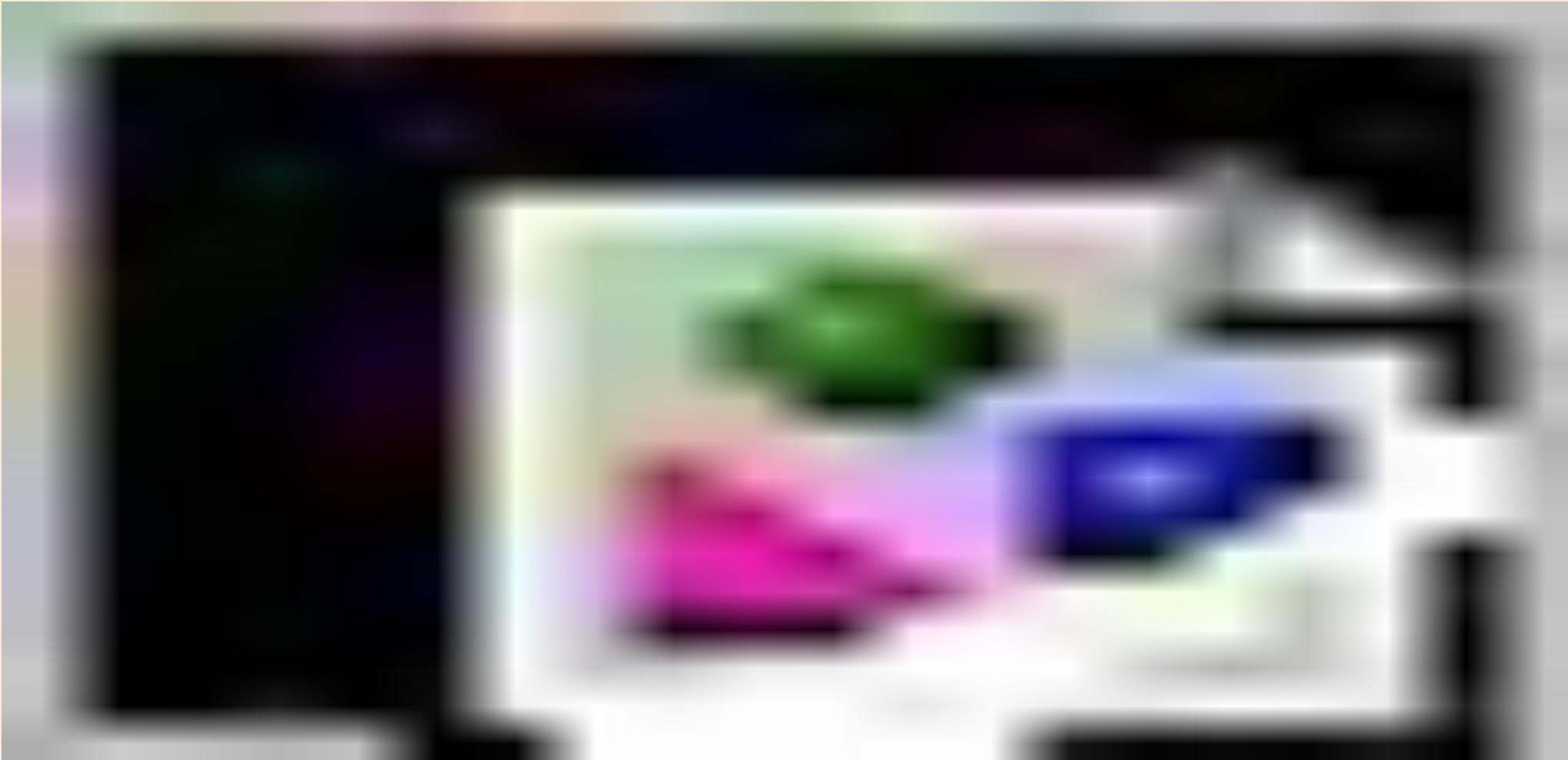
Keep your analysis or model up to date using scheduling



<https://docs.snowflake.com/en/user-guide/ui-snowsight/notebooks-schedule>

# Easily bring your model into Streamlit for reporting

## Share your work with stakeholders



<https://github.com/Snowflake-Labs/snowflake-python-recipes/tree/main/Snowflake%20ML%20Model%20in%20Streamlit>

# Learn about Snowflake Notebooks for ML

Documentation



<https://docs.snowflake.com/developer-guide/snowflake-ml/notebooks-on-spcs>

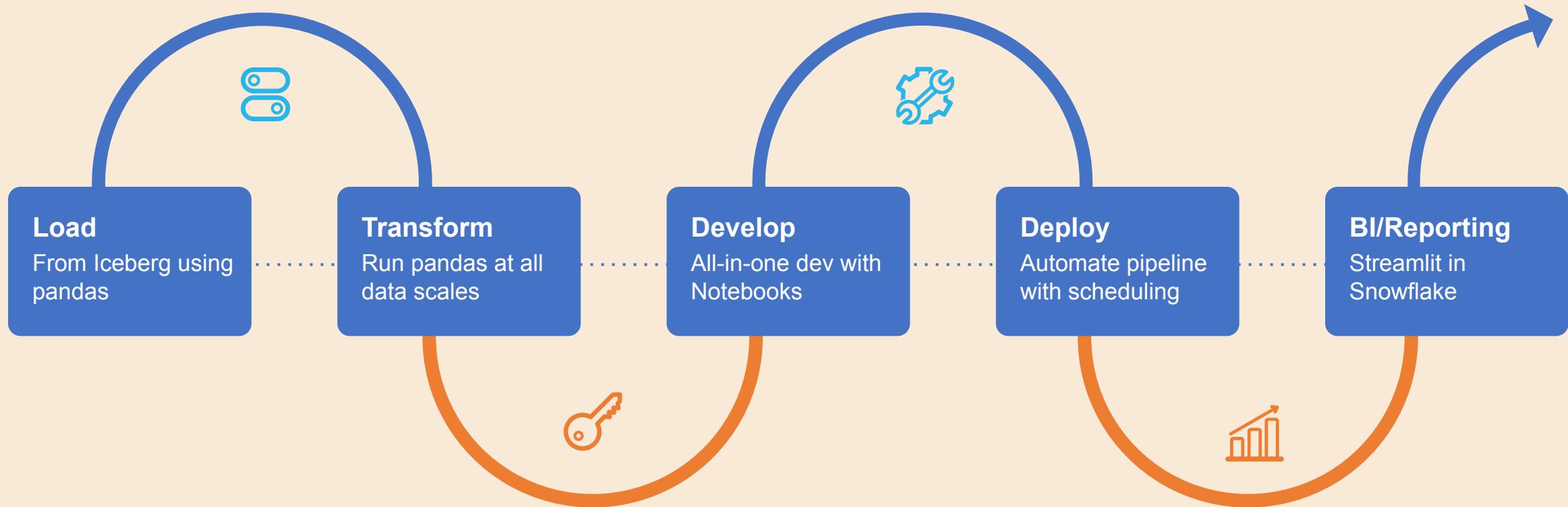
Container Runtime +  
Notebooks Quickstart



<https://quickstarts.snowflake.com/guide/notebook-container-runtime/>

# From data to model to insights

End-to-end development for high performance data teams



# Snowflake ❤️ PyData



## Easy

Built for developer productivity.  
Streamline workflow practices  
across data teams.



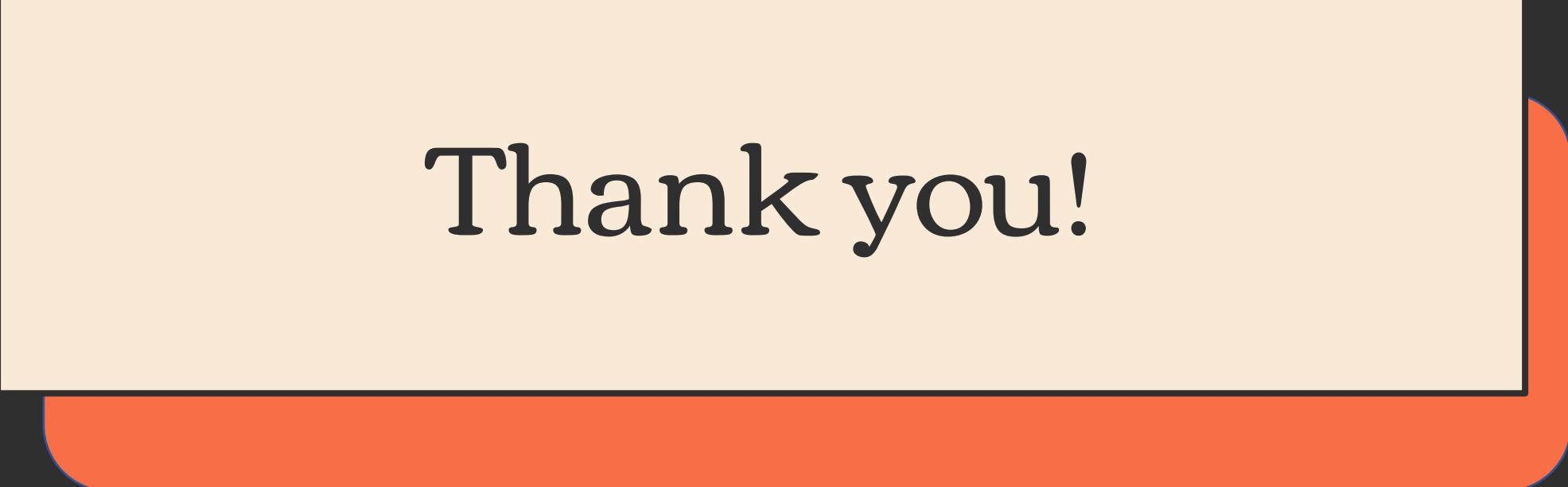
## Open

Embrace open formats and  
standards. Bring your familiar  
libraries and frameworks.



## Trusted

Interoperability of data,  
across clouds, with universal  
governance at scale



Thank you!

# Snowflake Python Recipes

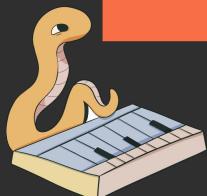
Collection of code examples, snippets, best practices for Python in Snowflake



**TRY IT ON YOUR OWN TODAY!**



Snowflake-Labs/snowflake-python-recipes



# We'd love to hear from you!

Fill out this 5 minute survey to share your experience



[snowflake.co1.qualtrics.com/jfe/form/SV\\_036qhfYM7YwEK7s](https://snowflake.co1.qualtrics.com/jfe/form/SV_036qhfYM7YwEK7s)



Come by Booth 309 for SWAG

