



# MAXIV

**Data Handling for ML and Analytics**

# Data Handling for ML and Analytics

Carla Takahashi

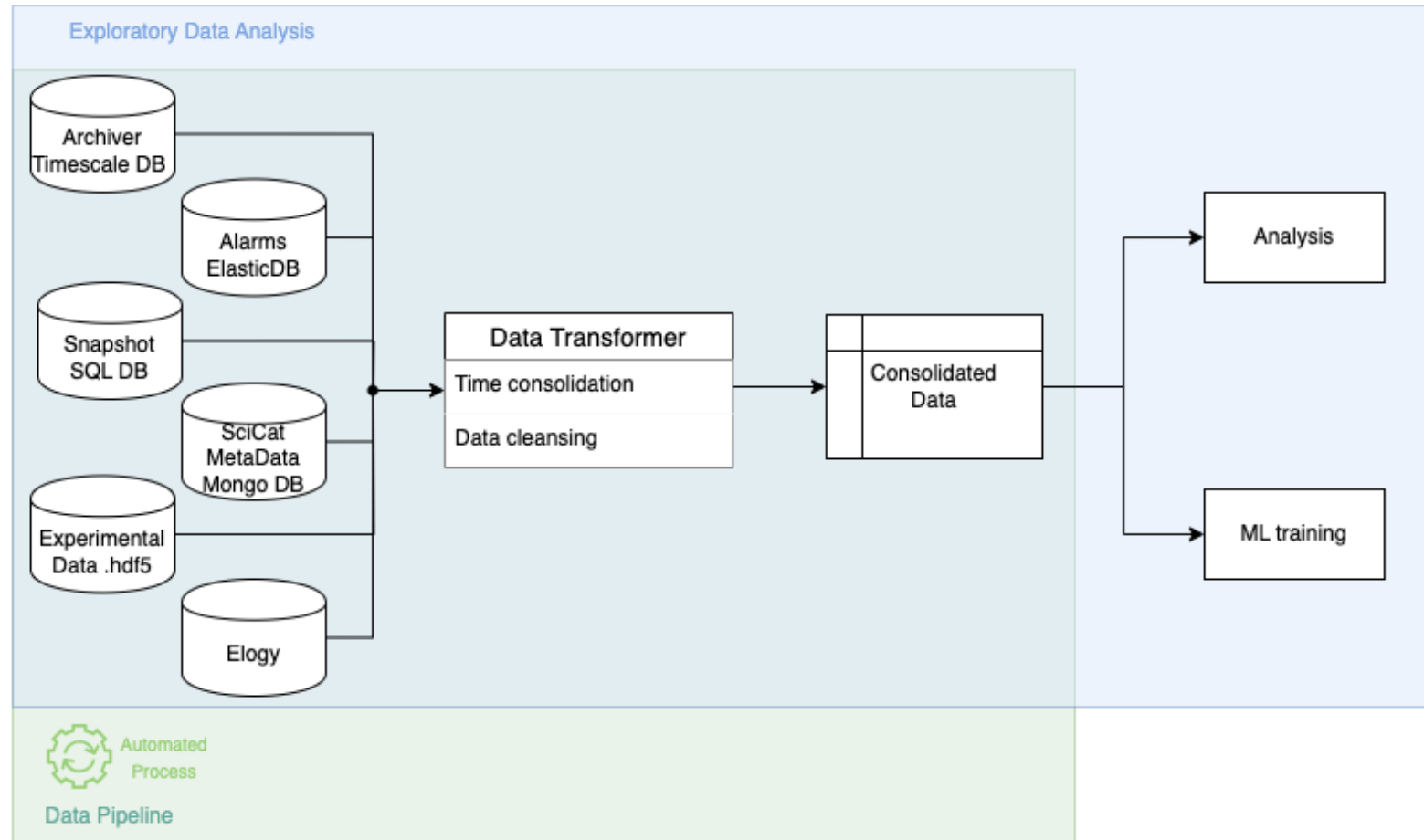
## Data Sources

- Archiver DB using Timescale.
- Alarms Logs.
- Snapshot Database.
- SciCat Database.
- Experimental Data.
- Electronic Logbook (possibly?).

	Characteristics	Challenges
Archiver	<ul style="list-style-type: none"> <li>From the control system side.</li> <li>Our main source of historical data.</li> <li>Stored on a Timescale database and allow some processing to be done on the db.</li> <li>Can be accessed from the white network.</li> </ul>	<ul style="list-style-type: none"> <li>From the machine side, there is a large amount of data available.</li> <li>From the beamlines side not much data is being archived.</li> <li>The data is not consolidated in time.</li> </ul>
Alarms	<ul style="list-style-type: none"> <li>Stores persistently all alarms.</li> <li>Stored in a Elastic Search.</li> </ul>	<ul style="list-style-type: none"> <li>It is a very large dataset.</li> <li>Data time resolution depends on the PLC time resolution.</li> </ul>
Snapshot	<ul style="list-style-type: none"> <li>Control system variables are stored triggered by different processes.</li> <li>For each snapshot data is consolidated in time.</li> </ul>	<ul style="list-style-type: none"> <li>The DB lives on the same host as the Tango DB and can not be accessed on the white network.</li> <li>It is uncertain the volume and quality of data stored there.</li> </ul>
SciCat	<ul style="list-style-type: none"> <li>SciCat database can also be used to metadata as input models.</li> <li>The quality of the experiment is one of the features available on the metadata.</li> </ul>	<ul style="list-style-type: none"> <li>Scientists are still on the process of learning scicat, thus the availability of data might be an issue.</li> </ul>
Experimental Data	<ul style="list-style-type: none"> <li>Experimental data, can be used to train models both for data analysis and for beamline alignment.</li> </ul>	<ul style="list-style-type: none"> <li>The use of experimental data might be subject to user/beamline restrictions.</li> <li>Experiments vary greatly, lead into data volume issues: too much or too little data.</li> </ul>
Elogy	<ul style="list-style-type: none"> <li>Electronic Logbook with human inputted data.</li> <li>Might contain calibration and operational information provided by scientists and operators.</li> </ul>	<ul style="list-style-type: none"> <li>Data is not curated or structured.</li> <li>Requires language processing to be used procedurally.</li> </ul>

# DATA CONSOLIDATION

- Even within the same data source, such as the Archiver, different attributes have different archiving periods, so we can say that our data is not consolidated in time.
- Different experiments and setups have different sets of metadata and snapshots stored. Thus, even considering a long period of time, the availability of data might be inconsistent.
- For experimental data, the scope tends to be better defined, both in time and which variables need to be stored, however, depending on the experiment the volume of data might be a challenge.



# Data Pipeline Challenges

- Currently, there is no standardized location for automated data processing. As a result, the ETL process is performed directly on the high-performance cluster and manually triggered through Jupyter notebooks.
- Similarly, there is no defined storage location for final training data blobs, which is important for auditing the final models. For now, these data are stored either locally or in Active Directory storage.
- Alternatively, the process used to generate the data could be stored, allowing for data regeneration if needed.
- Implementing an Airflow-based ETL solution is advisable.

# HIGH PERFORMANCE CLUSTER

- Jupyter Hub runs on the HPC.
- Has GPU nodes available.
- Can connect to the Archiver Timescale DB.
- Can connect to ML Flow to track experiments.
- Can create specific Kernels for different experiments.

The screenshot shows a Jupyter Notebook titled "Data Analysis.ipynb" running on a kernel. The notebook contains the following code:

```
[3]: import polars as pl
import pandas as pd
import matplotlib.pyplot as plt

import pickle
import time

from ydata_profiling import ProfileReport

[27]: f_path = "./Data"
f_name_prefix = "test_aug_8_polars"
data = pl.read_csv(f"{f_path}/{f_name_prefix}_cleaned", use_pyarrow=True)

[28]: data
```

The output of the code is a table with 12 columns and 10 rows of data. The columns are labeled with timestamps and various identifiers. The data is as follows:

Datetime	301L/DIA/BPM-01/xpossa	R3-301L/DIA/BPM-01/ypossa	R3-301M1/DIA/BPM-01/xpossa	R3-301M1/DIA/BPM-01/ypossa	R3-301M1/DIA/BPM-02/xpossa	R3-301M1/DIA/BPM-02/ypossa	R3-301U1/DIA/BPM-01/xpossa	R3-301U1/DIA/BPM-01/ypossa	R3-301U2/DIA/BPM-01/xpossa	R3-301U2/DIA/BPM-01/yj
2024-06-14 21:01:23	664	39822.0	305637.0	-534.0	75.0	45.0	-1511.0	491.0	1192.0	68.0
2024-06-14 21:01:24	664	37512.0	328994.0	-6157.0	238.0	334.0	-1205.0	4475.0	1330.0	808.0
2024-06-14 21:01:25	664	34332.0	322329.0	-4291.0	211.0	228.0	-1253.0	3375.0	1314.0	620.0
2024-06-14 21:01:26	664	31322.0	296317.0	-3692.0	174.0	147.0	-1326.0	2179.0	1294.0	436.0
2024-06-14 21:01:27	664	29716.0	316058.0	-1969.0	69.0	151.0	-1509.0	1866.0	1211.0	423.0
...	...	...	...	...	...	...	...	...	...	...
2024-06-14 21:58:56	664	30527.0	293541.0	353.0	7.0	-47.0	-1621.0	-467.0	1159.0	-131.0



The screenshot shows a Jupyter Notebook titled "MLFlow.ipynb" with the following content:

### Log experiment and model to ML Flow

#### Experiment Log

The ML Experiments of this notebook are stored in the [Experiments Tab On MLFlow Quickstart](#). Specifically on the Run named `bedecked-elik-947`

#### Model Register

You can see the registered model on the right side of the screen with the name `tracking-quickstart, v#`. The model is also available in the [Models Tab](#).

```
[29]: # Create a new MLflow Experiment
mflow.set_experiment("MLflow Quickstart")

# Start an MLflow run
with mflow.start_run():

    # Log the hyperparameters
    mflow.log_params(params)

    # Log the loss metric
    mflow.log_metric("accuracy", accuracy)

    # Set a tag that we can use to remind ourselves what this run was for
    mflow.set_tag("Training Info", "Basic LR model for iris data")

    # Infer the model signature
    signature = infer_signature(X_train, lr.predict(X_train))

    # Log the model
    model_info = mflow.sklearn.log_model(
        sk_model=lr,
        artifact_path="iris_model",
        signature=signature,
        input_example=X_train,
        registered_model_name="tracking-quickstart",
    )

Successfully registered model 'tracking-quickstart'.
2024/03/09 23:28:31 INFO mflow.store.model_registry.abstract_store: Waiting up to 300 seconds for model version to finish creation. Model name: tracking-quicksta
```

At the bottom, the status bar shows: "ML Flow v1 | Idle Mem: 348.60 MB Time remaining: 11:48 Mode: Command Ln 6, Col 59 MLFlow.ipynb"

# ML Flow

- Can be accessed from the HPC Jupyter Hub.
- Using the MLflow api within the jupyter notebook, it is possible to compare different runs for the same experiment.
- MLflow can publish models, in order to create a traceable repository of stable models.
- The Web Interface can also be used to compare models and publish them.

The screenshot shows the MLflow web interface for the "MLflow Quickstart" experiment. The page includes a search bar, a list of experiments, and a table of runs.

### MLflow Quickstart

Experiment ID: 1 Artifact Location: `mflow-artifacts/1`

Search Experiments:

Time created: [dropdown] State: Active [dropdown]

Datasets: [dropdown] Sort: Created [dropdown] Columns: [dropdown]

Table | Chart | Evaluation | Experimental

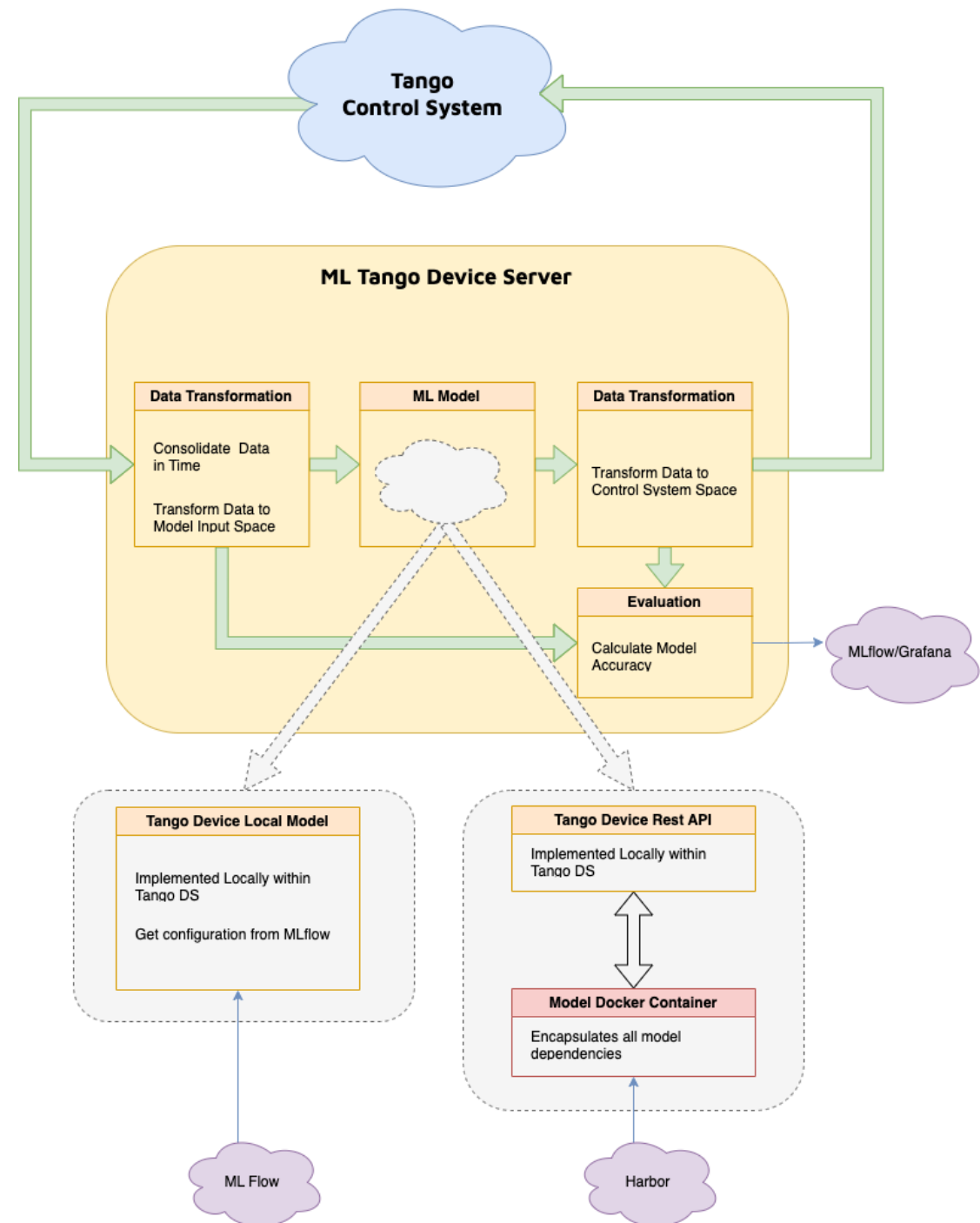
Run Name	Created	Dataset	Duration	Source	Models
bedecked-elik-947	7 months ago	-	30.5d	ipykerne...	tracking-q.../1

1 matching run



# REAL TIME DATA PROCESSING

- For use cases on control system applications (e.g. trajectory or alignment adjustment).
- Tango device deployed on the control system hosts.
- Should have access to the trained model.
- Should consolidate data collected on real time.
- Should apply model solutions on real time.
- Should have quality evaluation over time.
- Model deployment should be done via Ansible.



## Current State

- There is currently no standardized solution for data handling.
- Most challenges are being addressed as they arise.
- A machine learning project focused on Insertion Devices Compensation is currently underway, and we are using this opportunity to learn how MAX IV should approach ML-based development.

