

Supercharging the Machine Learning Lifecycle

Anniina Sallinen
Data Engineer @ Oura
She / her

CS major in University of Helsinki 2011-2018

- Algorithms and machine learning (+ software systems + distributed systems)

Worked in IT since 2014

- data domain since 2018
- Oura since 2021

Web development:

- Java

Data related development:

- Clojure
- Python
- Typescript
- AWS
- SQL
- IaC

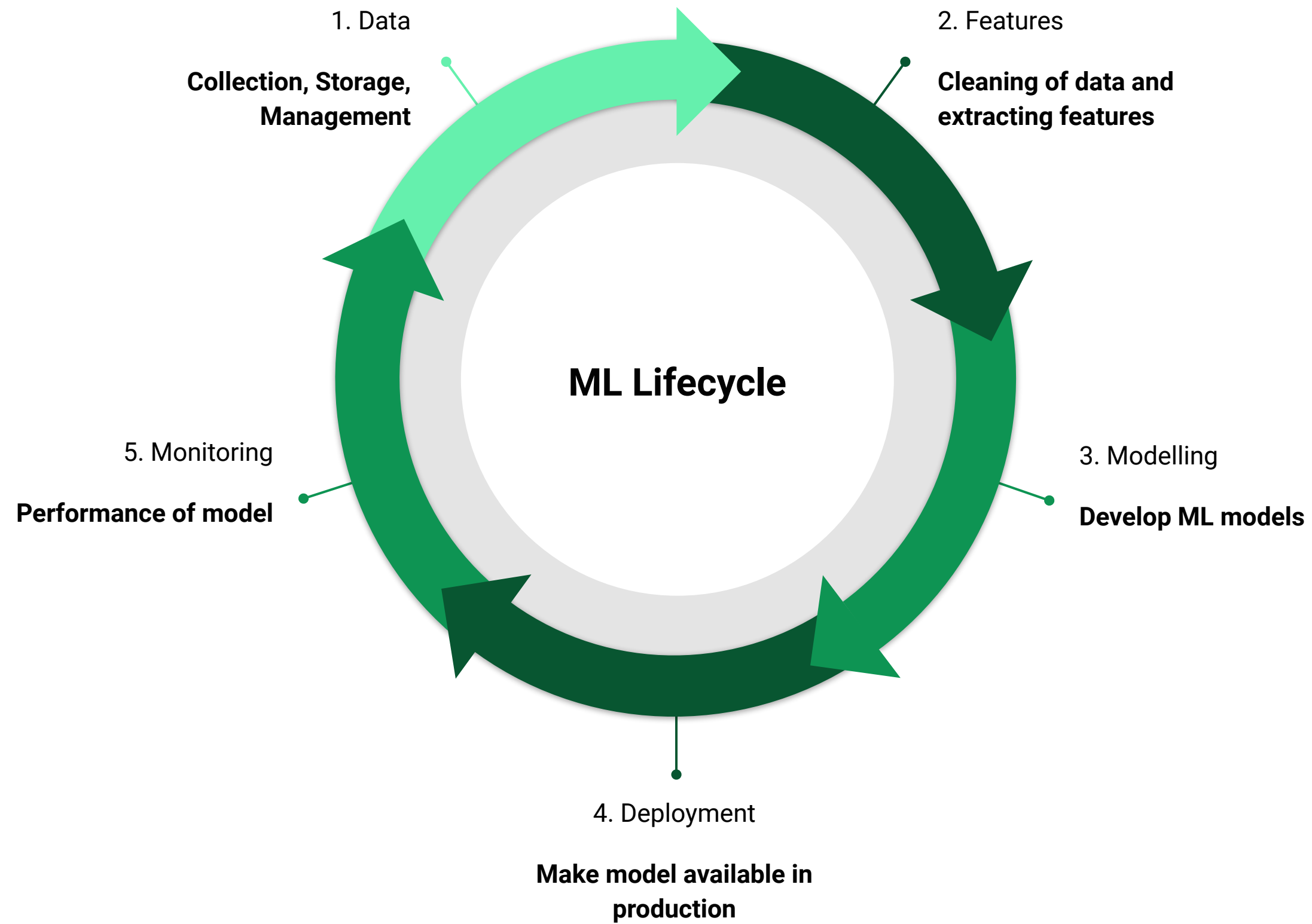


Anniina Sallinen

Agenda

- Introduction to machine learning life cycle
- Introduction to MLOps
- MLOps at Oura
- QA

Machine learning lifecycle



Data and features

Querying data from data lake / data warehouse

Sources of data for example: user data from software, internal systems (CRM), external systems (external APIs)

Data scientists usually need to join the data, calculate features based on existing data

Real world data is messy, so it often needs cleaning

Exploring and visualizing data to understand it

Garbage in,
garbage out

Modelling

Selecting suitable machine learning algorithm

Hyperparameter tuning

Usually repeated multiple times with different algorithms and hyperparameters, possibly different data

Slow and compute heavy part

Needs careful data selection and handling

Essential that training and validation of the model is done with different sets of data

Deployment

CI/CD pipelines to deploy models and ML systems to production

Continuous and automated re-training vs manual re-training and slower updates

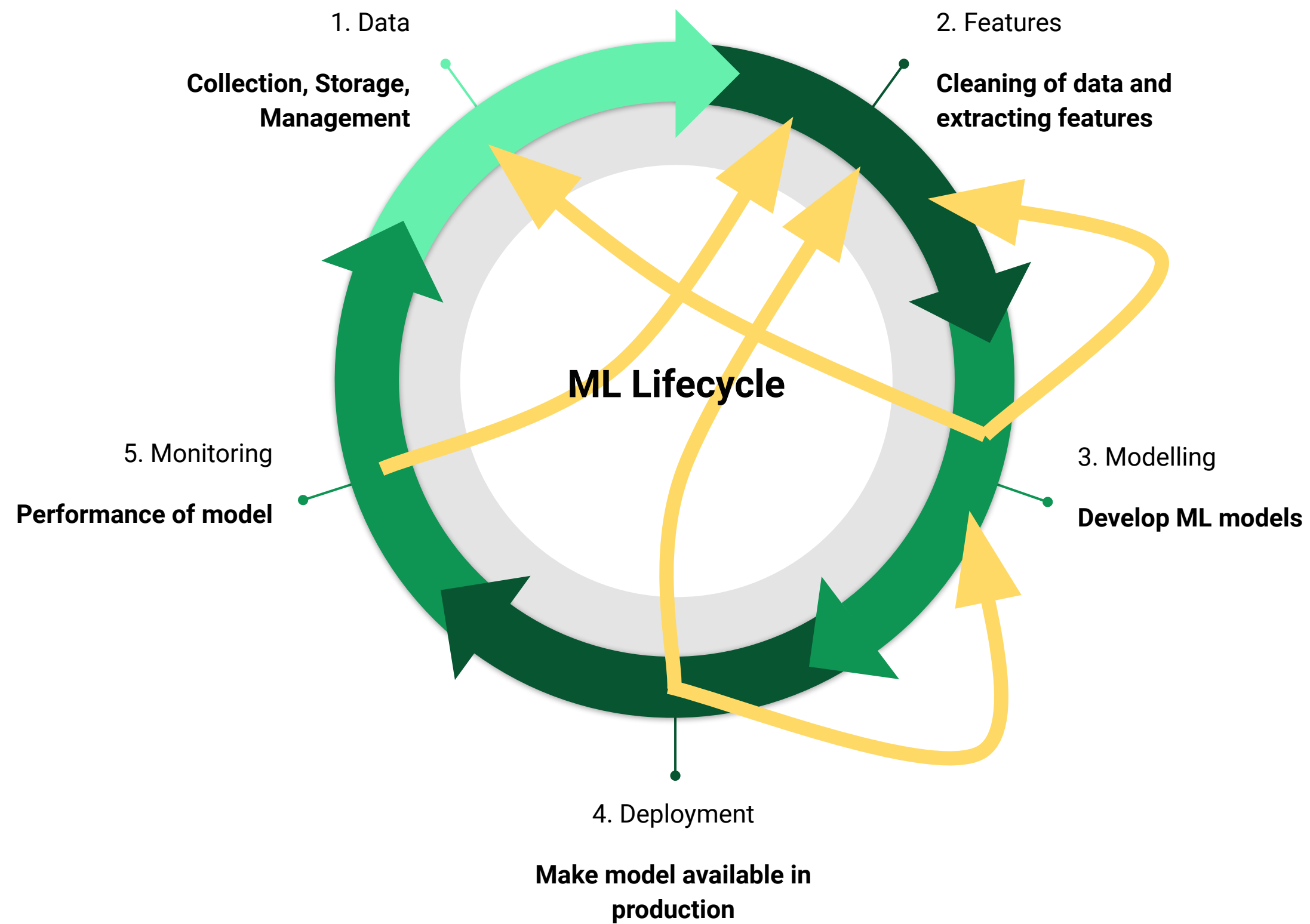
Testing that the model performs better than the previous one: A/B testing, shadow models, canary deployments

Monitoring

Monitoring system metrics such as latency, CPU load, memory usage

Monitoring of the model and data:
performance of the model, data drift

Model monitoring often not automated,
based on analysis performed by a data
scientist



Common challenges

Lack of

- Reproducibility
- Visibility
- Reliability
- Collaboration and sharing

“MLOps (machine learning operations) is a practice that aims to make developing and maintaining production machine learning seamless and efficient.”

-Valohai

“MLOps is an ML engineering culture and practice that aims at unifying ML system development (Dev) and ML system operation (Ops). Practicing MLOps means that you advocate for automation and monitoring at all steps of ML system construction, including integration, testing, releasing, deployment and infrastructure management.”

-Google

MLOps concepts

Feature store

Model training metadata store

Model registry

Feature store

Storage for data that is used in training machine learning models and for inference.

Features can be calculated ahead of time, before training or making predictions

Standardizes data and calculations. Metadata is also stored in the feature store

Can be used for monitoring data to detect changes in it over time

Model training metadata store

Storage for metadata about training machine learning models

In the metadata store we can store:

- Algorithm and hyperparameters used
- Compute environment
- Test / validation results
- Other metrics about training (execution time etc)

Can be combined with information about data used in training and model artifact to debug and selection of a model

Model registry

Registry for models

Storage for current and old versions of models

Can also contain information about when the model was trained and references to metadata about training

MLOps at Oura

Oura

Founded in 2013 in Oulu

Smart ring to improve well-being of individuals by providing insights of sleep, readiness and activity

Ring generations launched:

- Gen1 (2015)
- Gen 2 (2017)
- Gen 3 (2021)

Almost 500 employees globally, offices in Helsinki, Oulu and Tampere in Finland

Engineers working on hardware, mobile applications and cloud

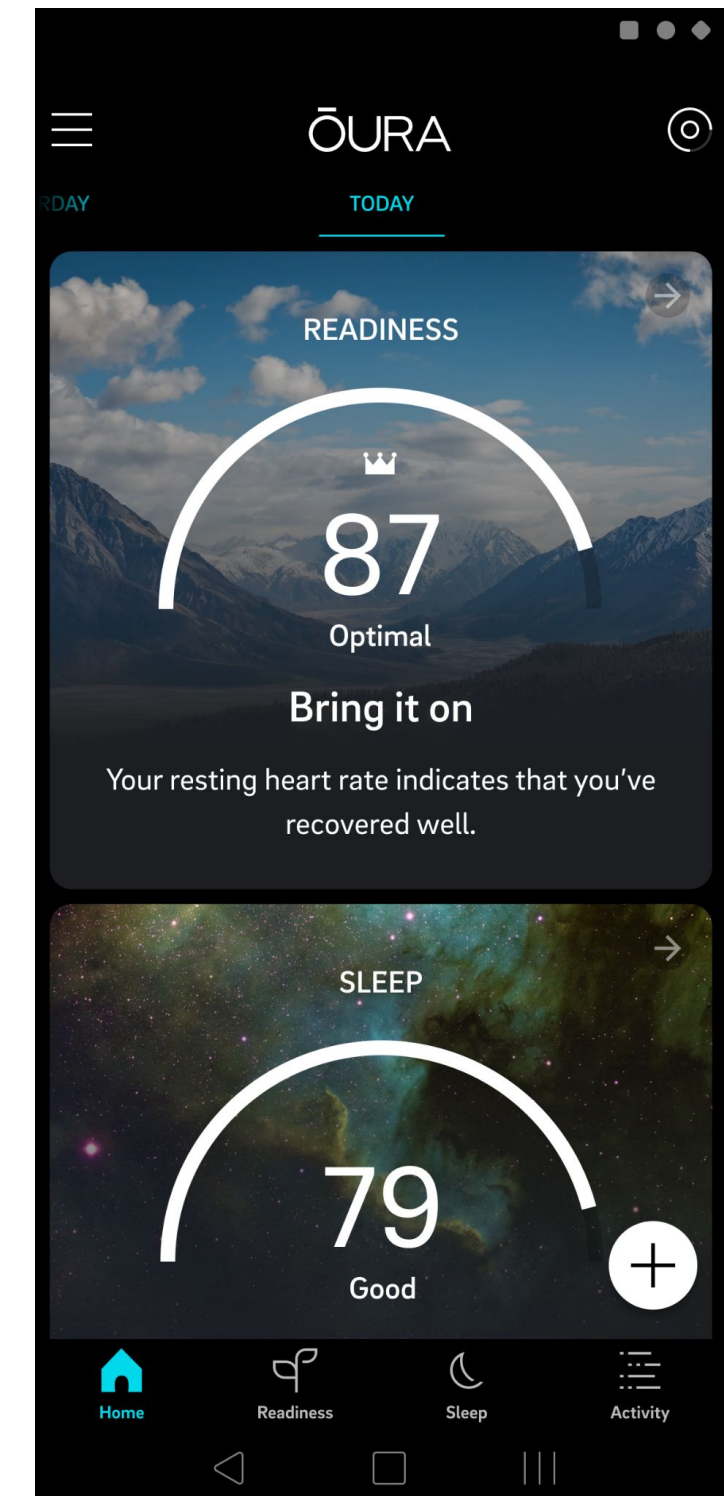
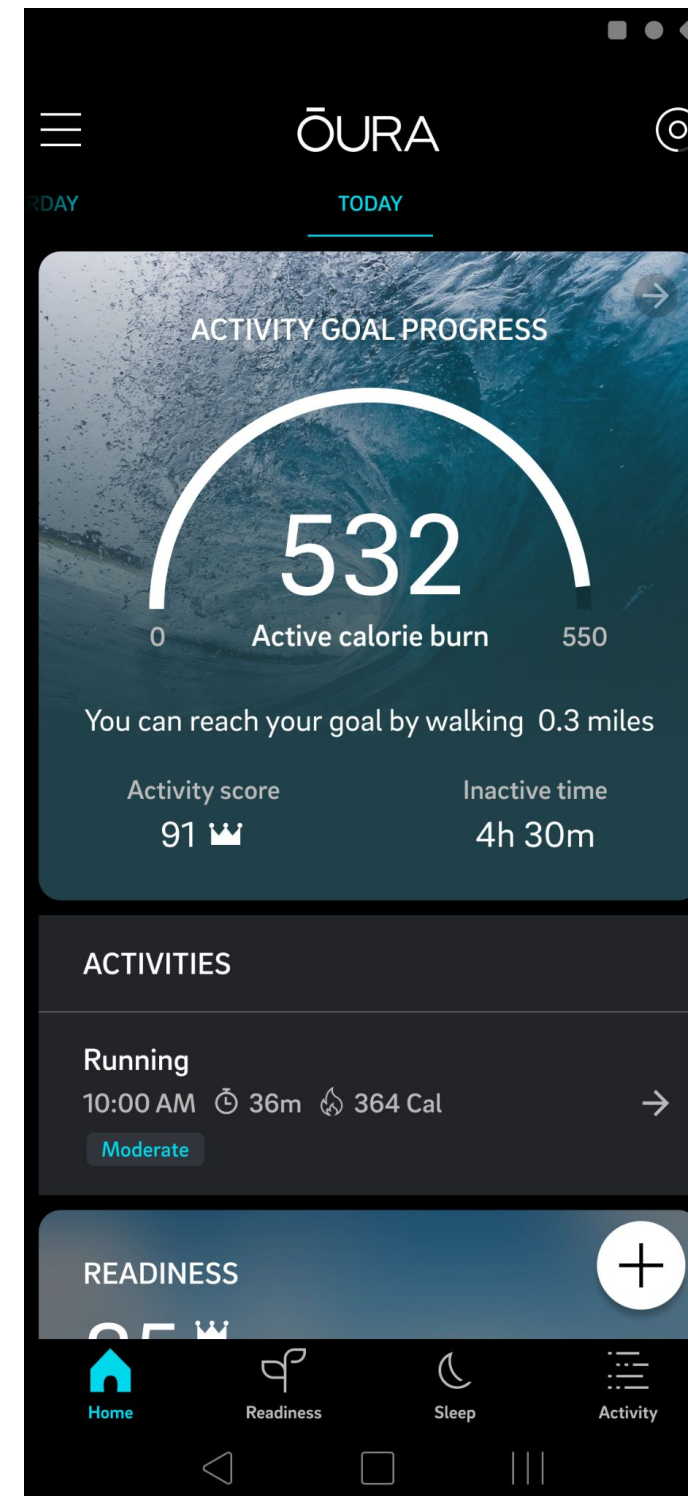


Oura's data

Oura's data is mostly IoT data from Oura ring, in total multiple TB and increasing quickly

When the data arrives to the cloud, it is already aggregated

The data arriving cloud might be already several days old, because ring and the mobile app can be used without internet connection



Machine learning models at Oura

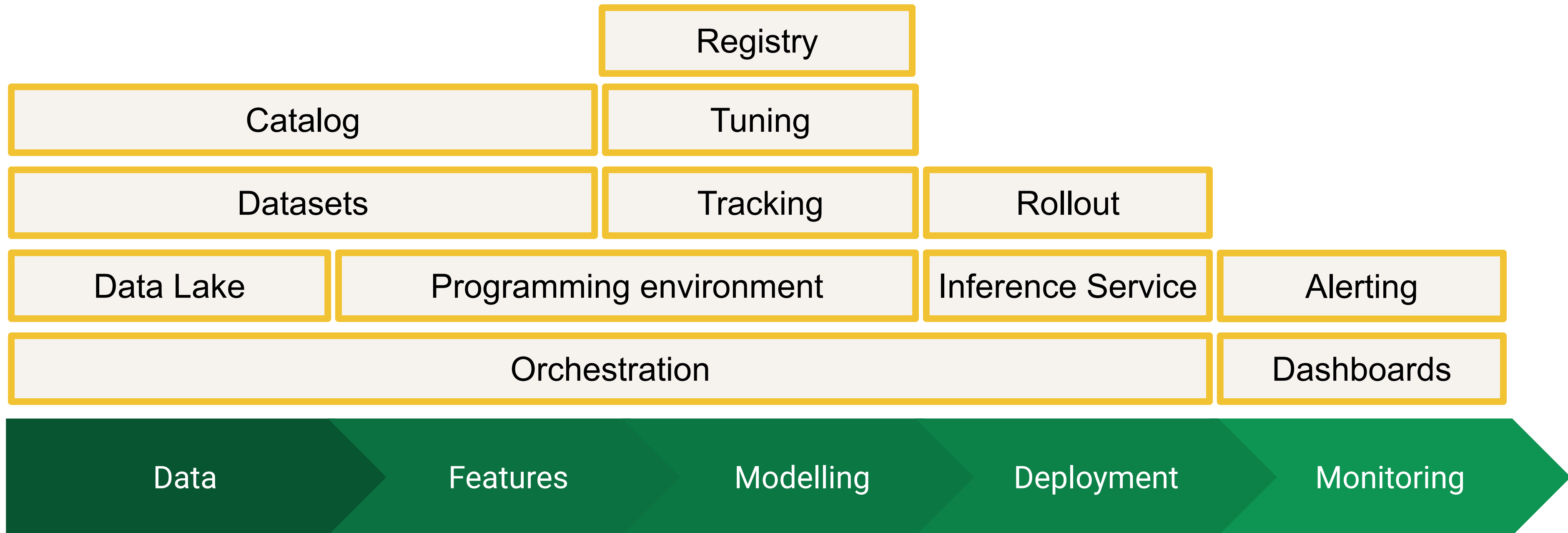
Due to the nature of our data, machine learning models are quite stable per ring release

Models running in ring, on mobile app and on cloud

When deciding whether to run the models on a device or in the cloud, we need to consider pros and cons of both approaches:

- Online vs offline usage
- Latency requirements
- Cost
- Compute power / hardware requirements

ML Stack



Data: DBT

Data from multiple sources ingested to data lake

To prepare data sets, we use dbt transformation tooling

DBT together with Elementary enables us to monitor our data and get alerts automatically when something is up with the data

DBT docs contains documentation for the data, including data schemas and dependencies between data

Data Lake

Alerting

Orchestration

Datasets

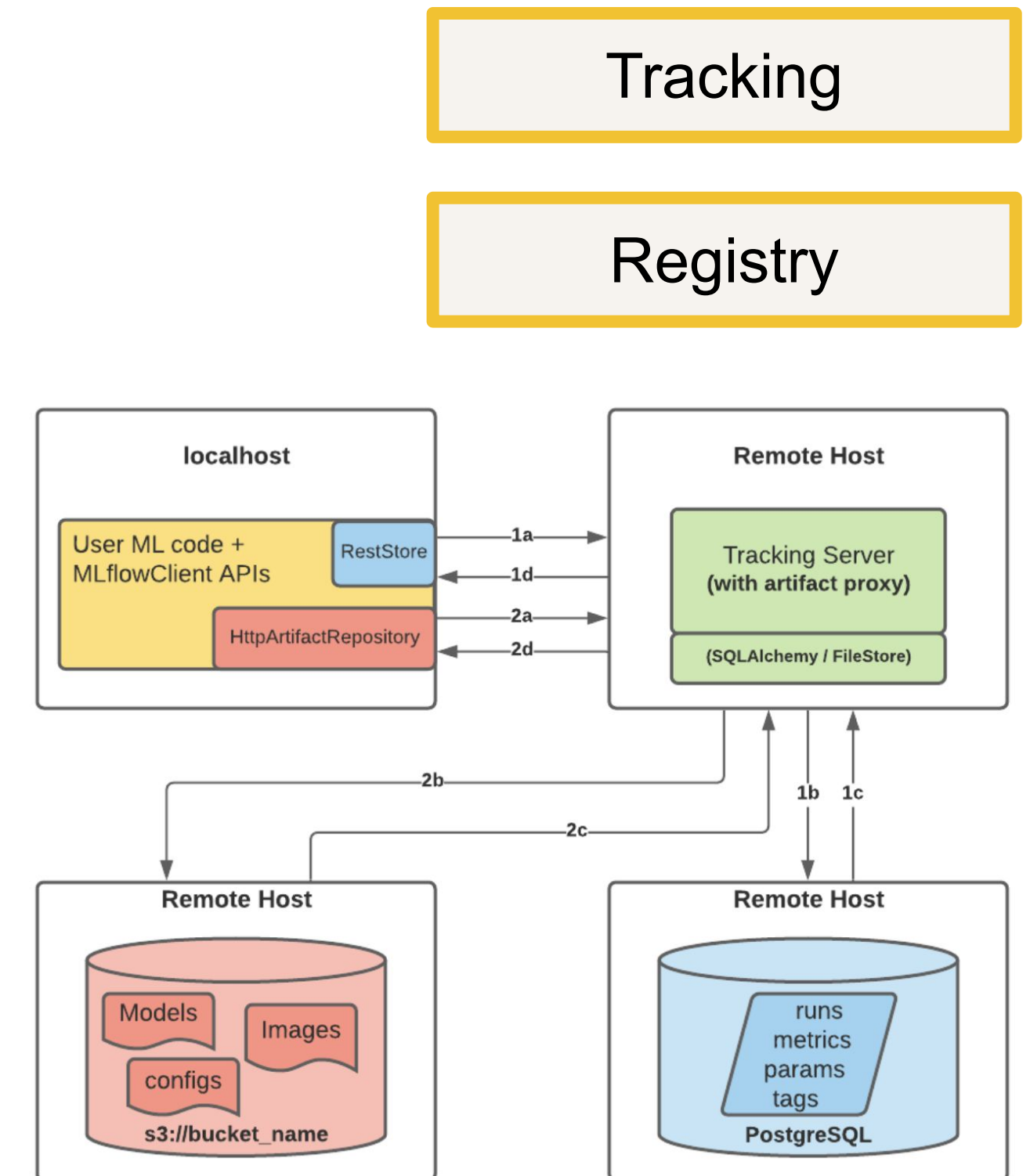
Tracking: MLFlow

Open source platform for tracking machine learning training runs / experiments

Stores metadata about the runs, metrics, artifacts

Can be also used as model registry

API + UI



mlflow 2.8.0 Experiments Models GitHub Docs

MLflow Quickstart >
kindly-stork-713

Run ID: e4bafb1bf4954d869dc72d8374c4efb9 Date: 2023-11-07 11:28:00 Source: ipykernel_launcher.py User: benjamin.wilson
Duration: 3.7s Status: FINISHED Lifecycle Stage: active

> Description [Edit](#)

> Datasets

Parameters (4)

Name	Value
max_iter	1000
multi_class	auto
random_state	8888
solver	lbfgs

The parameters that we logged

Metrics (1)

Name	Value
accuracy 🔗	1

The loss metric (accuracy) that we logged

Tags (1)

Name	Value	Actions
Training Info	Basic LR model for iris data	✎ 🗑
<input type="text"/>	<input type="text"/>	<input type="button" value="Add"/>

Our tag that we set for future reference

Artifacts

iris_model

Full Path: mlflow-artifacts:/846578415685150448/e4bafb1bf4954d869dc72d8374c4efb9/artifacts/iris_model [🔗](#)

tracking-quickstart, v1
Registered on 2023/11/07

Our registration link

Our model and its metadata

MLflow Model

The code snippets below demonstrate how to make predictions using the logged model. This model is also registered to the [model registry](#).

Model schema

Input and output schema for your model. [Learn more](#)

Name	Type
Inputs (1)	
-	Tensor (dtype: float64, shape: [-1,4])
Outputs (1)	
-	Tensor (dtype: int64, shape: [-1])

Our model signature →

Make Predictions

Predict on a Spark DataFrame:

```
import mlflow
from pyspark.sql.functions import struct, col
logged_model = 'runs:/e4bafb1bf4954d869dc72d8374c4efb9/iris_model'

# Load model as a Spark UDF. Override result_type if the model does not return double values.
loaded_model = mlflow.pyfunc.spark_udf(spark, model_uri=logged_model, result_type='double')

# Predict on a Spark DataFrame.
df.withColumn('predictions', loaded_model(struct(*map(col, df.columns))))
```

Predict on a Pandas DataFrame:

```
import mlflow
logged_model = 'runs:/e4bafb1bf4954d869dc72d8374c4efb9/iris_model'

# Load model as a PyFuncModel.
loaded_model = mlflow.pyfunc.load_model(logged_model)
```

> Datasets

Parameters (4)

Name	Value
max_iter	1000
multi_class	auto
random_state	8888
solver	lbfgs

Metrics (1)

Name	Value
accuracy	1

Screenshot source:
<https://mlflow.org/docs/latest/getting-started/intro-quickstart/index.html>

▼ Artifacts

▼ iris_model

MLmodel

conda.yaml

input_example.json

model.pkl

python_env.yaml

requirements.txt

Model schema

Input and output schema for your model. [Learn more](#)

Name	Type
Inputs (1)	
-	Tensor (dtype: float64, shape: [-1,4])
Outputs (1)	
-	Tensor (dtype: int64, shape: [-1])

Screenshot source:
<https://mlflow.org/docs/latest/getting-started/intro-quickstart/index.html>

Orchestration: Flyte

Workflow tooling that enables easy remote execution for code

Distributed computation: speed up, scale easily

Storing artifacts

Works well with MLFlow

Plain python, organise to workflows and tasks

Orchestration

Registry

Tuning


```

86 @task
87 def train_iris_model(iris_dataset: pd.DataFrame) -> dict:
88     y = iris_dataset[target_col]
89     X = iris_dataset[features]
90     X_train, X_test, y_train, y_test = train_test_split(
91         X, y, test_size=0.5, random_state=42
92     )
93     log_reg = LogisticRegression(max_iter=1000)
94     log_reg.fit(X_train, y_train)
95
96     test_prediction = log_reg.predict(X_test)
97     return metrics.classification_report(
98         y_train, test_prediction, digits=3, output_dict=True
99     )
100
101 @workflow
102 def test_workflow() -> dict:
103     df = load_iris_dataset()
104     visualize_iris_data(iris_data=df)
105     df_chunks = create_chunks(iris_dataset=df, chunk_amount=3)
106     processed_dfs = map_task(process_iris_data)(iris_data=df_chunks)
107     full_df = combine_data(chunks=processed_dfs)
108     return train_iris_model(iris_dataset=full_df)
109

```

Domain	Version	Cluster	Time	Duration	IAM Role	Service Account	Raw Output Prefix	Parallelism	Interruptible override	Overwrite cached outputs
development	vF-1H1meDtkOFGMmQHRnDA==		8/11/2023 6:55:40 AM UTC	36s	default	default	----	0	----	false
Nodes	Graph	Timeline								

Status ▾

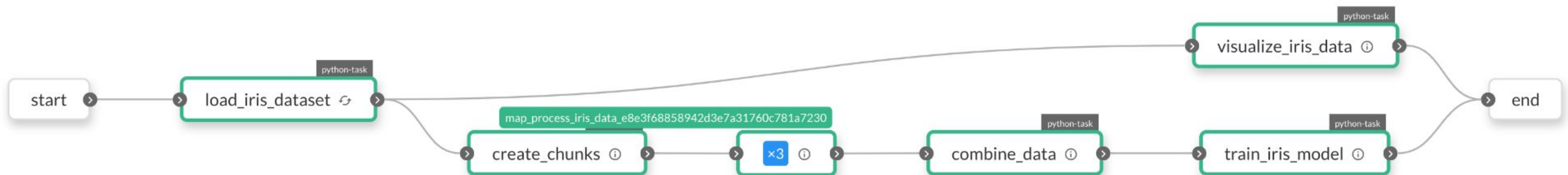
Start Time ▾

Duration ▾

<div>load_iris_dataset</div> <div>iris_example.load_iris_dataset</div>	n0	Python Task	SUCCEEDED	↺↻				<div>🖥️</div>	<div>↺</div>
<div>visualize_iris_data</div> <div>iris_example.visualize_iris_data</div>	n1	Python Task	SUCCEEDED	ⓘ	8/11/2023 6:55:41 AM UTC 8/11/2023 9:55:41 AM EEST	8s		<div>🖥️</div>	<div>↺</div>
<div>create_chunks</div> <div>iris_example.create_chunks</div>	n2	Python Task	SUCCEEDED	ⓘ	8/11/2023 6:55:41 AM UTC 8/11/2023 9:55:41 AM EEST	8s		<div>🖥️</div>	<div>↺</div>
<div>map_process_iris_data_e8e3f68858942d3e7a31760c781a7230</div> <div>iris_example.map_process_iris_data_e8e3f68858942d3e7a31760c781a7230</div>	n3	Map Task	SUCCEEDED	ⓘ	8/11/2023 6:55:51 AM UTC 8/11/2023 9:55:51 AM EEST	9s		<div>🖥️</div>	<div>↺</div>
<div>combine_data</div> <div>iris_example.combine_data</div>	n4	Python Task	SUCCEEDED	ⓘ	8/11/2023 6:56:01 AM UTC 8/11/2023 9:56:01 AM EEST	7s		<div>🖥️</div>	<div>↺</div>
<div>train_iris_model</div> <div>iris_example.train_iris_model</div>	n5	Python Task	SUCCEEDED	ⓘ	8/11/2023 6:56:09 AM UTC 8/11/2023 9:56:09 AM EEST	7s		<div>🖥️</div>	<div>↺</div>

Domain development	Version vF-1H1meDtkOFGMmQHRnDA==	Cluster	Time 8/11/2023 6:55:40 AM UTC	Duration 36s	IAM Role default	Service Account default	Raw Output Prefix ----	Parallelism 0	Interruptible override ----	Overwrite cached outputs false
------------------------------	--	---------	---	------------------------	----------------------------	-----------------------------------	----------------------------------	-------------------------	---------------------------------------	--

[Nodes](#) [Graph](#) [Timeline](#)



Challenges at Oura

Running models in ring, mobile, cloud

All of them have unique challenges:

- Updating models in ring / mobile
- Models that work both on iOS and android
- Selected technologies need to be relatively easy for the data scientists to learn
- Monitoring for models in mobile or ring
- Supporting multiple versions in cloud
- Privacy and security

Dependencies between models

Dependencies between predictions and user metrics

End to end testing with ring, mobile and cloud

Debugging of the models

Q&A