



**The MLOps platform
that makes you
productive, everywhere!**

Klarna Meetup

Stockholm, 2023-05-23

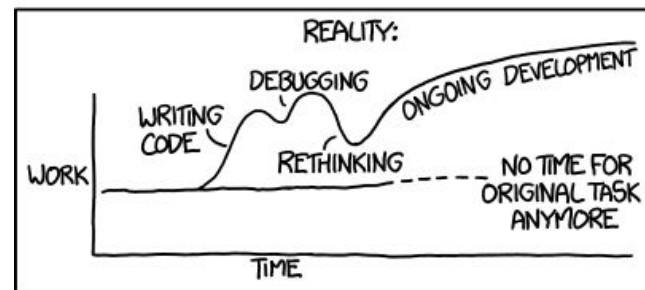
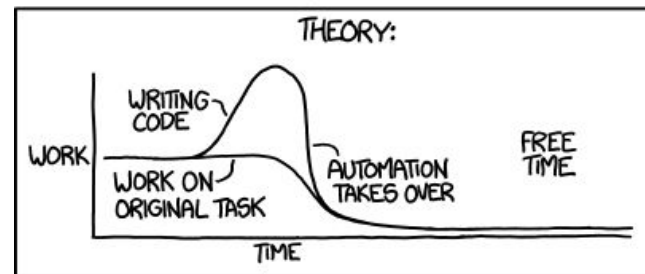


What MLOps is (not only) about ?

- Application of the DevOps principles to ML world
- Managing ML model lifecycle
- Tools and platforms
- *Automation* and processes
- Infrastructure as Code

The ultimate goal is: **PRODUCTIVITY**

"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



GID MLOps “Productivity Manifesto”

- Machine learning and data science should be ***first-class*** citizens of Data Platforms
- ***Open*** standards and cloud ***agnosticism***
- Short development ***feedback loop*** (incl. local dev)
- ***Fast*** new ML projects bootstrapping and ***standardization***
- Execution environment ***independent*** training pipelines
- Easy ***collaboration***
- ... MLOps capabilities provisioned ***in days not months***

ML projects in layers



**Data
Scientist**

Experimentation + EDA

Machine Learning frameworks

Example technologies:



ML projects in layers



Data Scientist

Experimentation + EDA

Machine Learning frameworks

?

Execution environment

Data



MLOps / ML Engineer

Example technologies:



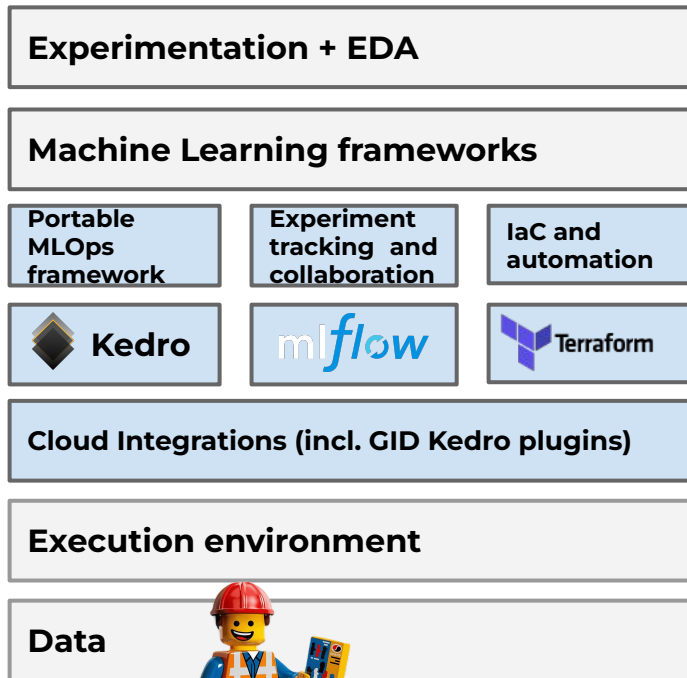
XGBoost



Building blocks of the GID MLOps



Data Scientist



MLOps / ML Engineer

Example technologies:



XGBoost

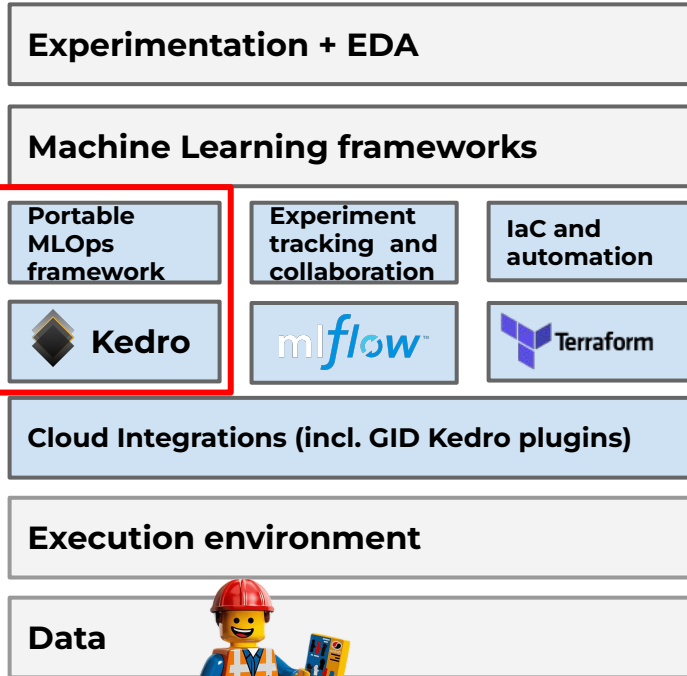
GID MLOps platform



Building blocks of the GID MLOps



Data Scientist



MLOps / ML Engineer

Example technologies:



XGBoost

GID MLOps Platform



What is Kedro?



Kedro

=

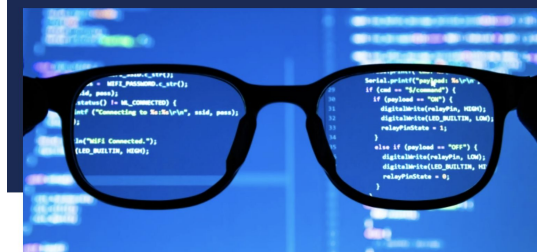
Software
Engineering
Principles

+

Data Science

Kedro is an open-source Python framework for creating reproducible, maintainable and modular data science code.

McKinsey donates machine learning pipeline tool Kedro to the Linux Foundation



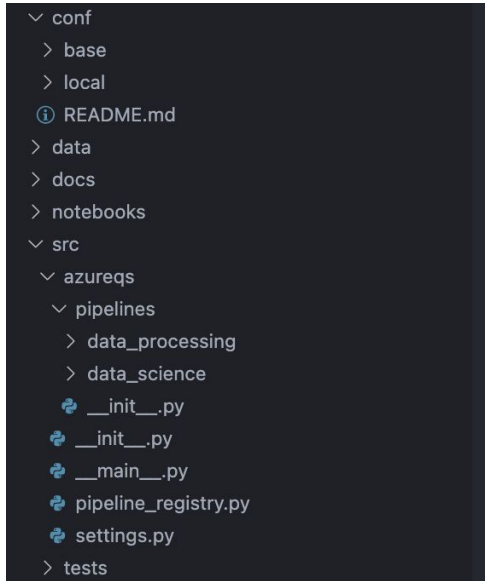
What features does Kedro have? (Part 1)

```

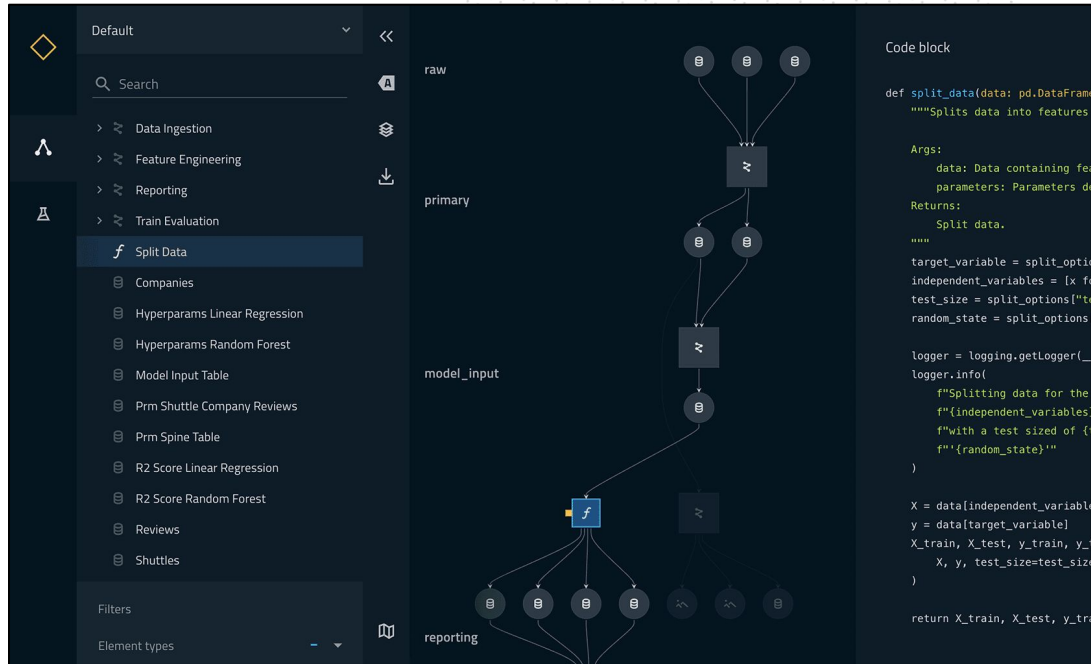
└─ conf
  ├── base
  ├── local
  └─ README.md
└─ data
└─ docs
└─ notebooks
└─ src
  └─ azureqs
    └─ pipelines
      ├── data_processing
      ├── data_science
      ├── __init__.py
      ├── __init__.py
      ├── __main__.py
      ├── pipeline_registry.py
      └── settings.py
  └─ tests
```

**Well defined
project structure**

What features does Kedro have? (Part 1)



Well defined
project structure
+ project starters



Nodes & pipelines
abstractions

Kedro pipeline - data engineering

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=preprocess_companies,
11                inputs="companies",
12                outputs="preprocessed_companies",
13                name="preprocess_companies_node",
14            ),
15            node(
16                func=preprocess_reviews,
17                inputs="reviews",
18                outputs="preprocessed_reviews",
19                name="preprocess_reviews_node",
20            ),
21            node(
22                func=create_model_input_table,
23                inputs=["preprocessed_reviews", "preprocessed_companies", "ratings"],
24                outputs="model_input_table",
25                name="create_model_input_table_node",
26            ),
27        ]
28    )
```

Kedro pipeline - data science

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=split_data,
11                inputs=["model_input_table", "params:model_options"],
12                outputs=["X_train", "X_test", "y_train", "y_test"],
13                name="split_data_node",
14            ),
15            node(
16                func=train_model,
17                inputs=["X_train", "y_train"],
18                outputs="regressor",
19                name="train_model_node",
20            ),
21            node(
22                func=evaluate_model,
23                inputs=["regressor", "X_test", "y_test"],
24                outputs=None,
25                name="evaluate_model_node",
26            ),
27        ]
28    )
29
```

Kedro node

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=preprocess_companies,
11                inputs="companies",
12                outputs="preprocessed_companies",
13                name="preprocess_companies_node",
14            ),
15            node(
16                func=preprocess_reviews,
17                inputs="reviews",
18                outputs="preprocessed_reviews",
19                name="preprocess_reviews_node",
20            ),
21            node(
22                func=create_model_input_table,
23                inputs=["preprocessed_reviews", "
24                outputs="model_input_table",
25                name="create_model_input_table_node",
26            ),
27        ]
28    )
```

```
49 def create_model_input_table(
50     reviews: pd.DataFrame, companies: pd.DataFrame, ratings: pd.DataFrame
51 ) -> pd.DataFrame:
52     """Combines all data to create a model input table.
53
54     Args:
55         reviews: Preprocessed data for reviews.
56         companies: Preprocessed data for companies.
57         ratings: Raw data for ratings.
58
59     Returns:
60         Model input table.
61
62     """
63     reviews_with_ratings = reviews.merge(ratings, left_on="id", right_on="rating_id")
64     model_input_table = reviews_with_ratings.merge(
65         companies, left_on="company_id", right_on="id"
66     )
67     model_input_table = model_input_table.dropna()
68     return model_input_table
```

What about parameters?

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=split_data,
11                inputs=["model_input_table", "params:model_options"],
12                outputs=["X_train", "X_test", "y_train", "y_test"],
13                name="split_data_node",
14            ),
15            node(
16                func=train_model,
17                inputs=["X_train", "y_train"],
18                outputs="regressor",
19                name="train_model_node",
20            ),
21            node(
22                func=evaluate_model,
23                inputs=["regressor", "X_test", "y_test"],
24                outputs=None,
25                name="evaluate_model_node",
26            ),
27        ]
28    )
29
```

What about parameters?

```
6 def create_pipeline(**kwargs) -> Pipeline:
```

```
7     return pipeline(  
8         [  
9             node(  
10                func=split_data,  
11                inputs=["model_input_table", "y_train"],  
12                outputs=["X_train", "X_test", "y_test"],  
13                name="split_data_node",  
14            ),  
15            node(  
16                func=train_model,  
17                inputs=["X_train", "y_train"],  
18                outputs="regressor",  
19                name="train_model_node",  
20            ),  
21            node(  
22                func=evaluate_model,  
23                inputs=["regressor", "X_test", "y_test"],  
24                outputs=None,  
25                name="evaluate_model_node",  
26            ),  
27        ]  
28    )  
29
```

```
  ✓ conf
```

```
  ✓ base
```

```
  ✓ parameters
```

```
    ! data_processing.yml
```

```
    ! data_science.yml
```

```
    ! azureml.yml
```

```
    ! catalog.yml
```

```
    ! logging.yml
```

```
    ! parameters.yml
```

```
  > local
```

```
  > data
```

```
  > docs
```

```
2
```

```
3
```

```
4
```

```
5
```

```
6
```

```
7
```

```
8
```

```
9
```

```
10
```

```
11
```

```
12
```

```
13
```

```
14
```

```
15
```

```
model_options:
```

```
  test_size: 0.2
```

```
  random_state: 3
```

```
features:
```

```
  - engines
```

```
  - passenger_capacity
```

```
  - crew
```

```
  - d_check_complete
```

```
  - moon_clearance_complete
```

```
  - iata_approved
```

```
  - company_rating
```

```
  - review_scores_rating
```

ML model?

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=split_data,
11                inputs=["model_input_table", "params:model_options"],
12                outputs=["X_train", "X_test", "y_train", "y_test"],
13                name="split_data_node",
14            ),
15            node(
16                func=train_model,
17                inputs=["X_train", "y_train"],
18                outputs="regressor",
19                name="train_model_node",
20            ),
21            node(
22                func=evaluate_model,
23                inputs=["regressor", "X_test", "y_test"],
24                outputs=None,
25                name="evaluate_model_node",
26            ),
27        ]
28    )
29
```



What about data?

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9             node(
10                func=preprocess_companies,
11                inputs="companies",
12                outputs="preprocessed_companies",
13                name="preprocess_companies_node",
14            ),
15            node(
16                func=preprocess_reviews,
17                inputs="reviews",
18                outputs="preprocessed_reviews",
19                name="preprocess_reviews_node",
20            ),
21            node(
22                func=create_model_input_table,
23                inputs=["preprocessed_reviews", "preprocessed_companies", "ratings"],
24                outputs="model_input_table",
25                name="create_model_input_table_node",
26            ),
27        ]
28    )
```

Kedro Data Catalog

```
6 def create_pipeline(**kwargs) -> Pipeline:
7     return pipeline(
8         [
9         node(
10            func=preprocess_companies,
11            inputs="companies",
12            outputs="preprocessed_companies",
13            name="preprocess_companies_node",
14        ),
15        node(
16            func=preprocess_reviews,
17            inputs="reviews",
18            outputs="preprocessed_reviews",
19            name="preprocess_reviews_node",
20        ),
21        node(
22            func=create_model_input_table,
23            inputs=["preprocessed_reviews", "preprocessed_companies", "ratings"],
24            outputs="model_input_table",
25            name="create_model_input_table_node",
26        ),
27    ]
28 )
```

```

  ✓ conf
  ✓ base
    > parameters
    ! azureml.yml
  ! catalog.yml
  ! logging.yml
  ! parameters.yml
  > local
  > data
  > docs
  > notebooks
  ✓ src
```

```
42 companies:
43     type: pandas.CSVDataSet
44     filepath: data/01_raw/companies.csv
45
46 reviews:
47     type: pandas.ParquetDataSet
48     filepath: data/01_raw/reviews.parquet
49
50 pictures:
51     type: pillow.ImageDataSet
52     filepath: data/01_raw/images/*.jpg
53
```

What features does Kedro have? (Part 2)

```
companies: Local catalog.yml
  type: pandas.CSVDataSet
  filepath: data/01_raw/companies.csv

reviews:
  type: pandas.ParquetDataSet
  filepath: data/01_raw/reviews.parquet

pictures:
  type: pillow.ImageDataSet
  filepath: data/01_raw/images/*.jpg
```

```
companies: Cloud catalog.yml
  type: pandas.CSVDataSet
  filepath: abfs://my_blob_container/data/01_raw/companies.csv

reviews:
  type: pandas.SQLQueryDataSet
  sql: "select * from reviews;"
  credentials: db_credentials

pictures:
  type: kedro_azureml.AzureMLFileDataSet
  dataset: my_dataset_from_azureml
  filepath: data/01_raw/images/*.jpg
```

Data Catalog

What features does Kedro have? (Part 2)

```
companies: Local catalog.yml
  type: pandas.CSVDataSet
  filepath: data/01_raw/companies.csv

reviews:
  type: pandas.ParquetDataSet
  filepath: data/01_raw/reviews.parquet

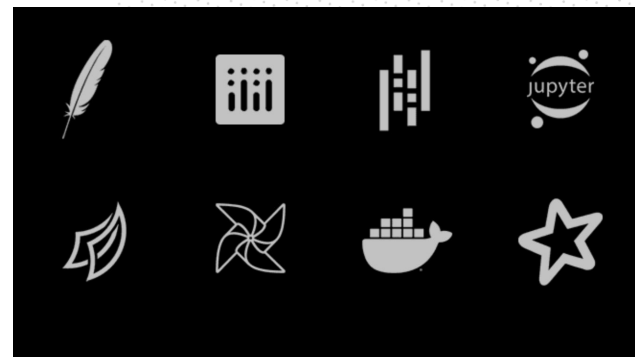
pictures:
  type: pillow.ImageDataSet
  filepath: data/01_raw/images/*.jpg
```

```
companies: Cloud catalog.yml
  type: pandas.CSVDataSet
  filepath: abfs://my_blob_container/data/01_raw/companies.csv

reviews:
  type: pandas.SQLQueryDataSet
  sql: "select * from reviews;"
  credentials: db_credentials

pictures:
  type: kedro_azureml.AzureMLFileDataSet
  dataset: my_dataset_from_azureml
  filepath: data/01_raw/images/*.jpg
```

Data Catalog

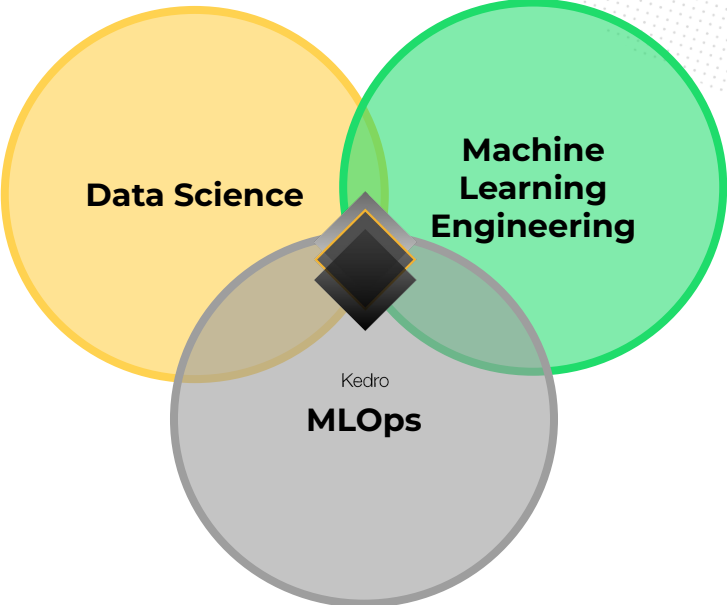


Extensibility & Integrations

Kedro can be integrated with multiple industry leading solutions, including: Apache Spark, Pandas, Dask, Matplotlib, Plotly, fsspec, Apache Airflow, Jupyter Notebook and Docker.

Cloud agnostic with Kedro

Machine Learning Frameworks



Model serving frameworks



Infrastructure tools





MLOps Platform by GetinData

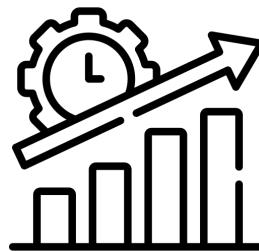
What we're achieving with MLOps Platform



Reliable experiment tracking



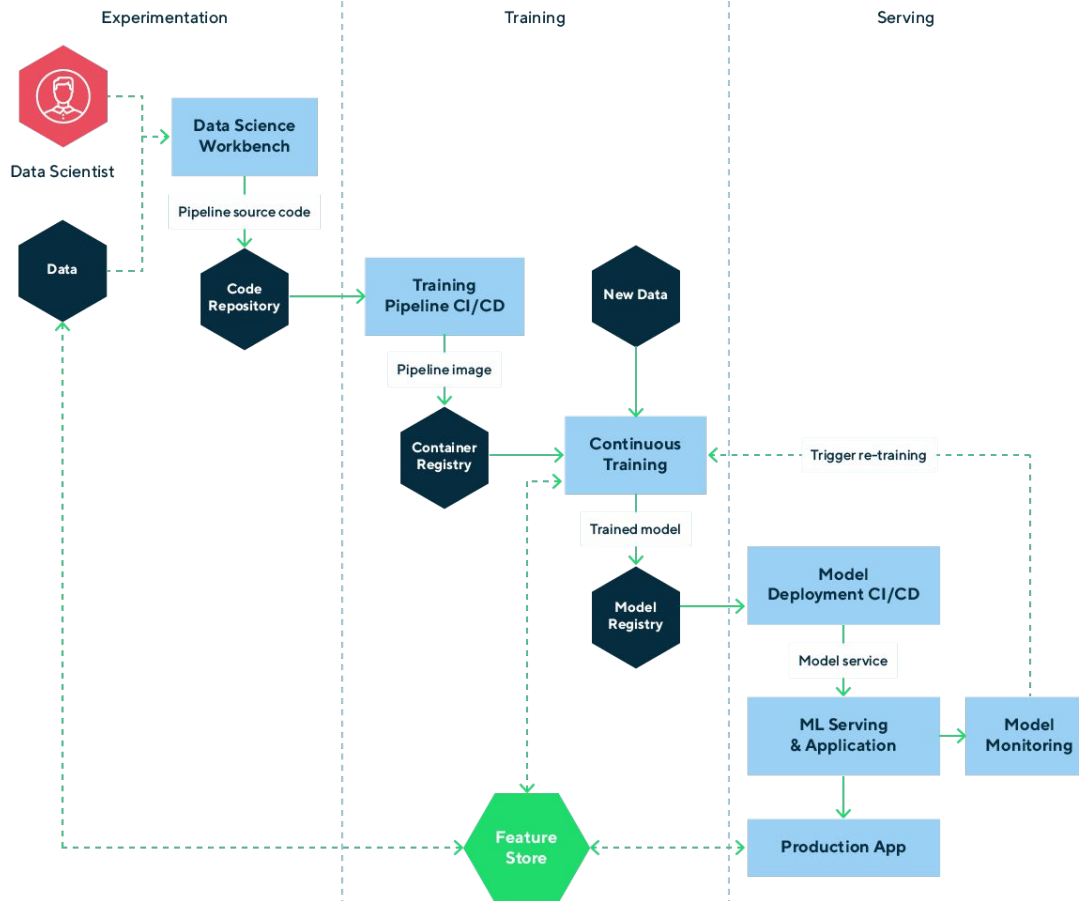
Easy path from local to cloud



Faster time to market



Model management



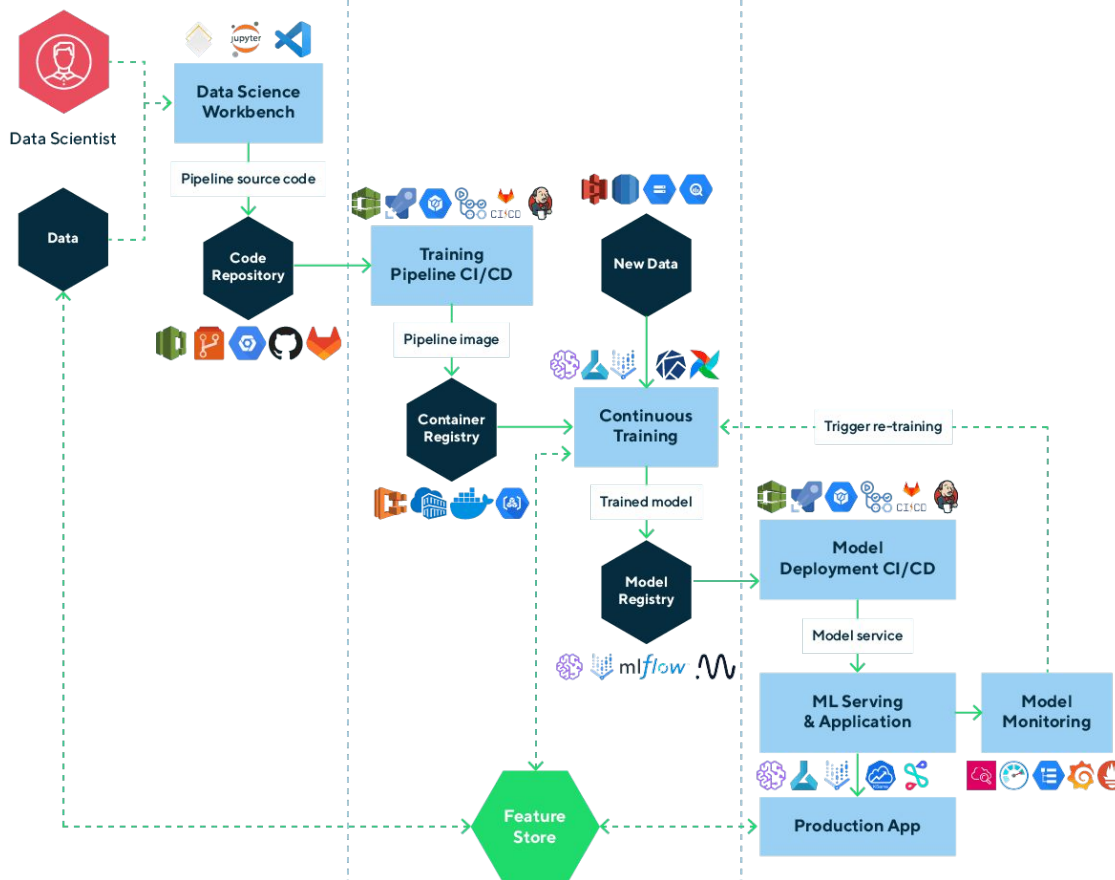
MLOps Platform

We built a best-of-breed solution instead of all-in-one

Experimentation

Training

Serving



MLOps Platform

A framework of best MLOps practices to make the process of ML experimentation, model training, and model serving efficient, secure and reliable. We built a best-of-bread solution instead of all-in-one

Supported cloud platforms



Future integrations



Write once - run (almost) everywhere



Kedro



Kedro VertexAI (GCP)

github.com/getindata/kedro-vertexai



Kedro AzureML (Azure)

github.com/getindata/kedro-azureml



Kedro Sagemaker (AWS)

github.com/getindata/kedro-sagemaker



Kedro Airflow (Kubernetes)

github.com/getindata/kedro-airflow-k8s



Kedro Kubeflow (Kubernetes)

github.com/getindata/kedro-kubeflow



In progress:

Kedro Snowflake & Databricks

How plugins work?

Pipeline definition



Cloud Native SDKs



Serialized pipeline deployment

```
def train_model(X_train: pd.DataFrame, y_train: pd.Series)
-> LinearRegression:
    regressor = LinearRegression()
    regressor.fit(X_train, y_train)
    return regressor
```

```
@dsl.component(kfp_package_path=_KFP_PACKAGE_PATH)
def train_model(X_train: pd.DataFrame, y_train: pd.Series) -> str:
    ...
    return regressor
```

```
"exec-data-science-active-modelling-pipeline-train-model-node": {
  "container": {
    "args": [
      "kedro vertexai -e local initialize-job --params='{\"data_science\": {\"active_modelling_pipeline\": {\"model_options\": {\"test_size\": 0.2, \"random_state\": 3, \"features\": [\"engines\", \"passenger_capacity\", \"crew\", \"d_check_complete\", \"moon_clearance_complete\", \"iata_approved\", \"company_rating\", \"review_scores_rating\"]}, \"candidate_modelling_pipeline\": {\"model_options\": {\"test_size\": 0.2, \"random_state\": 8, \"features\": [\"engines\", \"passenger_capacity\", \"crew\", \"review_scores_rating\"]}}}' && KEDRO_VERTEXAI_DISABLE_CONFIG_HOOK=false
      KEDRO_CONFIG_RUN_ID={{$.pipeline_job_uuid}} KEDRO_CONFIG_JOB_NAME={{$.pipeline_job_name}}
      KEDRO_VERTEXAI_RUNNER_CONFIG='{\"storage_root\": \"mb-temp/mlops-webinar-demo\"}' kedro run -e local --pipeline
      __default__ --node \"data_science.active_modelling_pipeline.train_model_node\" --runner
      kedro_vertexai.vertex_ai.runner.VertexAIPipelinesRunner --config config.yaml"
    ], ...
  }
}
```

Our **Plugins** translate the ML Pipeline from Kedro to selected execution environment.

Vertex AI

- Dashboard
- Datasets
- Feature Store
- Labeling tasks
- Workbench
- Pipelines**
- Training
- Experiments
- Model Registry
- Endpoints
- Batch predictions
- Metadata
- Matching Engine

Marketplace

spaceflights-20230202112306

CLONE

STOP

DELETE

Runtime Graph

9/9 steps completed

Expand Artifacts

100%



Amazon SageMaker Studio

File Edit View Run Kernel Git

default-1669297364894 / Personal

Launcher Pipelines ml-ops-housing-demo execution-167240312159

Less than 20 seconds ago

execution-1672403121591

Status 30/12/2022, 13:25 23m10s Elapsed time Stop

Graph Parameters Settings

Search for step...

```

    graph TD
      A(select_features_node) --> B(preprocess_data_node)
      B --> C(split_data_node)
      C --> D(prepare_features_node)
      C --> E(train_model_node)
      D --> F(evaluate_model_node)
      E --> F
  
```

152% zoom

Simple 0 0 0 0

execution-1672403121591

Microsoft Azure Machine Learning

Search within your workspace

This workspace ml-ops-sandbox

getindata.com > ml-ops-sandbox > Jobs > xbia-demo > _default_

Refresh Clone Resubmit Publish Schedule Show lineage Delete Cancel

default Completed Job overview

```

    graph TD
      A(select_features_node) --> B(preprocess_data_node)
      B --> C(split_data_node)
      C --> D(prepare_features_node)
      C --> E(train_model_node)
      D --> F(evaluate_model_node)
      E --> F
      G(feature_transformer)
  
```

feature.transformer

Google Cloud

gid-ml-ops-sandbox

mlops-cloud-agnostic-demo-20221230123318

Clone Stop Delete

Runtime Graph 7/7 steps completed Expand Artifacts 100%

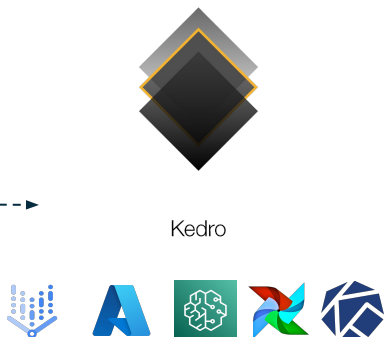
```

    graph TD
      A(mlflow-start-run) --> B(select-features-node)
      A --> C(preprocess-data-node)
      A --> D(split-data-node)
      A --> E(prepare-features-node)
      A --> F(train-model-node)
      A --> G(evaluate-model-node)
      B --> C
      C --> D
      D --> E
      D --> F
      E --> G
      F --> G
  
```

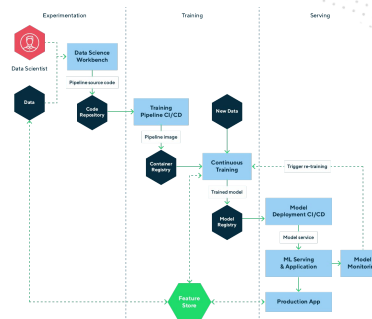
Logs

Our approach: cloud agnostic MLOps Platform

Client 1
Client 2
Client 3



MLOps Framework



MLOps Platform



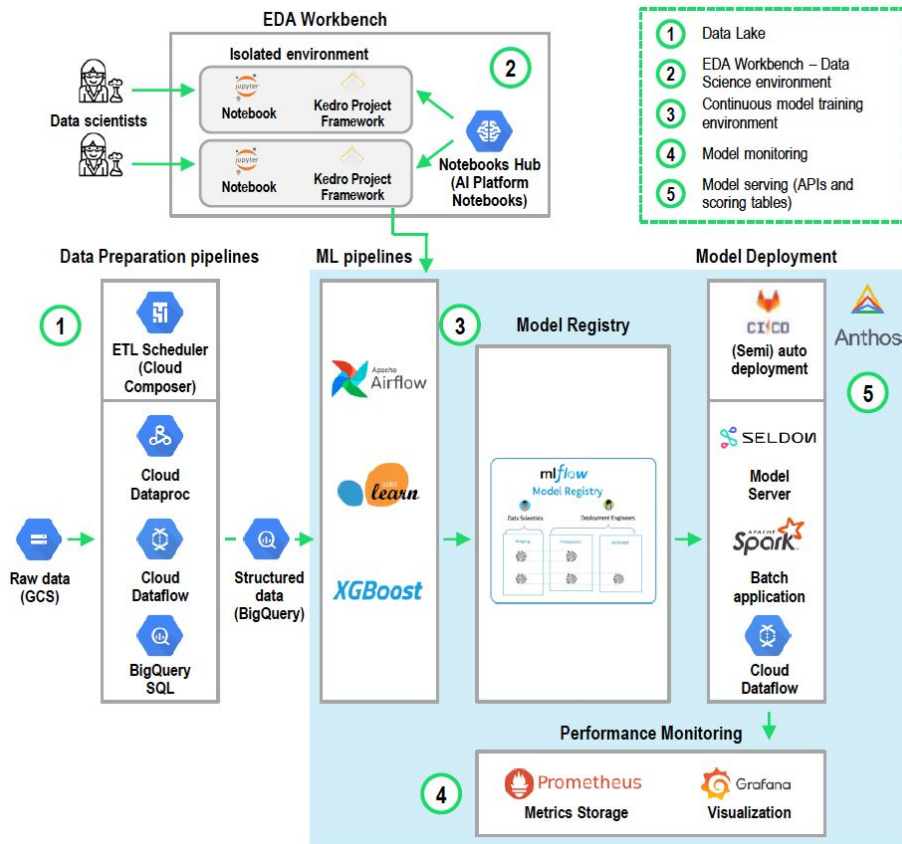
Completed project!

With **cloud agnostic** approach, we can unify MLOps projects with one framework (Kedro) and deploy it into various cloud platforms with Kedro Plugins.

MLOps Platform to Run Production ML/AI Models (banking)

CONFIDENTIAL

Leading Bank in CEE



Key technologies

A collection of logos for the key technologies used in the MLOps platform:

- Apache Airflow
- Google Cloud Platform
- Jupyter
- Apache Spark
- mlflow
- SELDON
- Grafana

Modern Data Platform, analytics and MLOps (FinTech)

We have delivered **Data Platform** that **supports general analytics** (ad-hoc querying, reporting in BI tools) and Machine Learning initiatives. Data Platform is build on top of Google Cloud services (BigQuery, Cloud Composer, Data Studio) and open-source projects (dbt, Terraform).

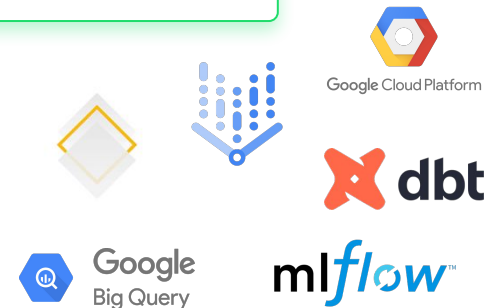
We have delivered an **MLOps Platform** that **supports a production-grade ML model lifecycle**. Our solution is using heavily Google Cloud Platform services (Vertex AI, BigQuery, Cloud Build) and open-source projects (mlflow, Kedro).

We have co-developed a number of **production AI/ML models** such as User Suspension Model, Invoice Model, Activation Model.



Willa is a Sweden and U.S.-based FinTech that helps professional freelancers, influencers, and social media content creators get paid immediately by brands for their freelance work and paid collaborations.

Key technologies



[Building robust ML & Analytics capability very early at a FinTech](#)



MLOps Labs



Michał Bryś

Machine Learning Engineer at Getindata

michal.brys@getindata.com



Marek Wiewiórka

Chief Data Architect at Getindata

marek@getindata.com



Marcin Zabłocki

MLOps Architect at Getindata

marcin.zablocki@getindata.com

Let's meet our (hypothetical) clients



Client 1

- A few ML models, no pipelines, manual / script based on VMs
- Data Scientists distributed across different teams
- Strict compliance requirements



Client 2

- Small data volume, but this will change quickly as the business grows
- One model, trained on the DS workstation / Jupyter Notebook
- Limited budget for the infrastructure
- Small team, overwhelmed by tasks



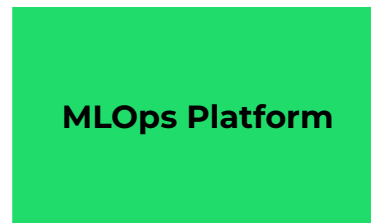
Client 3

- Current ML workloads running on-premises on Kubernetes
- Exploring the cloud services
- Lots of technical debt



Our approach: cloud agnostic MLOps Platform

Client 1
Client 2
Client 3



Completed
project!

With **cloud agnostic** approach, we can unify **MLOps** projects with one framework (Kedro) and deploy them to various cloud platforms with Kedro Plugins.

ML model training in layers



Experimentation + EDA

Machine Learning framework

MLOps Framework

Integrations (plugins)

Execution environment (local, cloud)

The data

Example technologies:



XGBoost



Kedro



Kedro abstracts the pipeline from the execution environment SDK.

Under-engineering refers to building with reduced complexity resulting in a **less robust, efficient and capable product**. It happens because of tight deadlines or because of a lack of understanding.

We **under-engineer** machine-learning prototypes and create code with a lot of **technical debt**.

Technical debt is intentional or accidental decisions that make **code difficult to understand, maintain, extend and fix errors**. Much like a loan, you pay a higher cost later, as it decreases the team's agility as the project matures.

What features does Kedro have?

```
project-template # Project folder
├── conf          # Configuration files
├── data          # Local project data
├── docs          # Documentation
├── logs          # Logs of pipeline runs
├── notebooks     # Exploratory Jupyter notebooks
├── pyproject.toml # Identifies the project root
├── setup.cfg     # Configuration options for tes
├── README.md    # README.md explaining your pro
├── setup.cfg     # Configuration options for tes
└── src          # Source code for pipelines
```

FEATURES

Project Templates

Kedro starter contains code in the form of a [Cookiecutter](#) template for a ML project. Metaphorically, a starter is similar to using a pre-defined layout when creating a presentation or document.

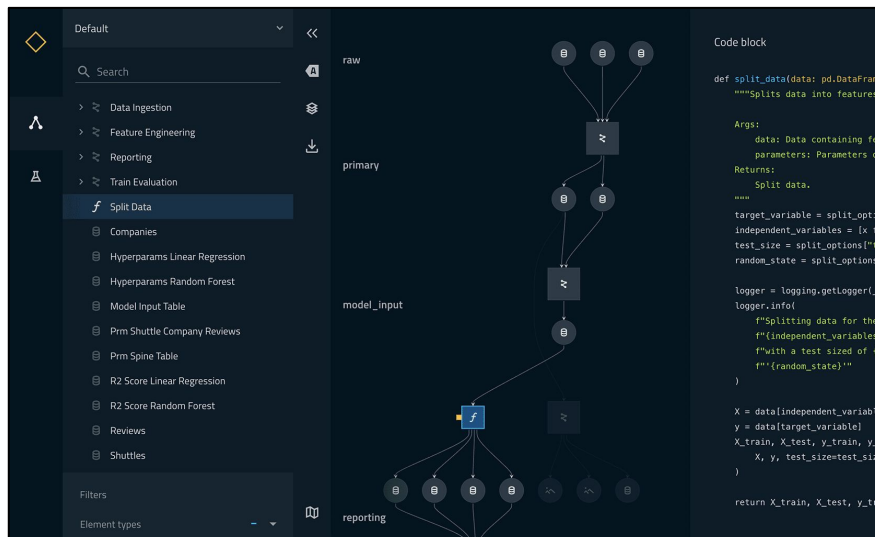


FEATURES

Extensibility & Integrations

Kedro can be integrated with multiple industry leading solutions, including Apache Spark, **Pandas**, Dask, Matplotlib, Plotly, fsspec, Apache Airflow, **Jupyter Notebook** and **Docker**.

What features does Kedro have?



FEATURES

Pipeline Visualisations

[Kedro's pipeline visualisation plugin](#) shows a blueprint of your developing data and machine-learning workflows, provides **data lineage**, keeps track of machine-learning experiments and makes it easier to collaborate with business stakeholders.

The image shows two code blocks in a dark-themed editor. The first block is a comment: '# Load a Spark DataFrame on S3'. Below it is a function definition for 'flight_patterns' that returns a Spark DataFrame with specific parameters. The second block is a comment: '# Save an image created with Matplotlib on Google Cloud Storage'. Below it is a function definition for 'results_plot' that saves a plot to GCS with specific parameters.

FEATURES

Data Catalog

A series of lightweight **data connectors** used to save and load data across many different file formats and file systems.



Pillars of the GID MLOps approach