

DATA+AI
SUMMIT 2022

MLOps on Databricks

A how-to guide

Outline

Part I

- Introduction
 - What is MLOps
 - Why care about MLOps?
 - People and process
 - Guiding principles
- Fundamentals
 - Semantics of dev, staging, prod
 - ML deployment patterns
- Reference architectures

Part II

- The use case
- Overview of pipelines
- Tooling
 - DevOps
 - DataOps
 - ModelOps
- Project structure
- Demo
- Getting started

Part I

Introduction

- What is MLOps?
- Why care about MLOps
- People and process
- Guiding principles

Hardest Part of ML isn't ML, it's Data

"Hidden Technical Debt in Machine Learning Systems," Google NIPS 2015

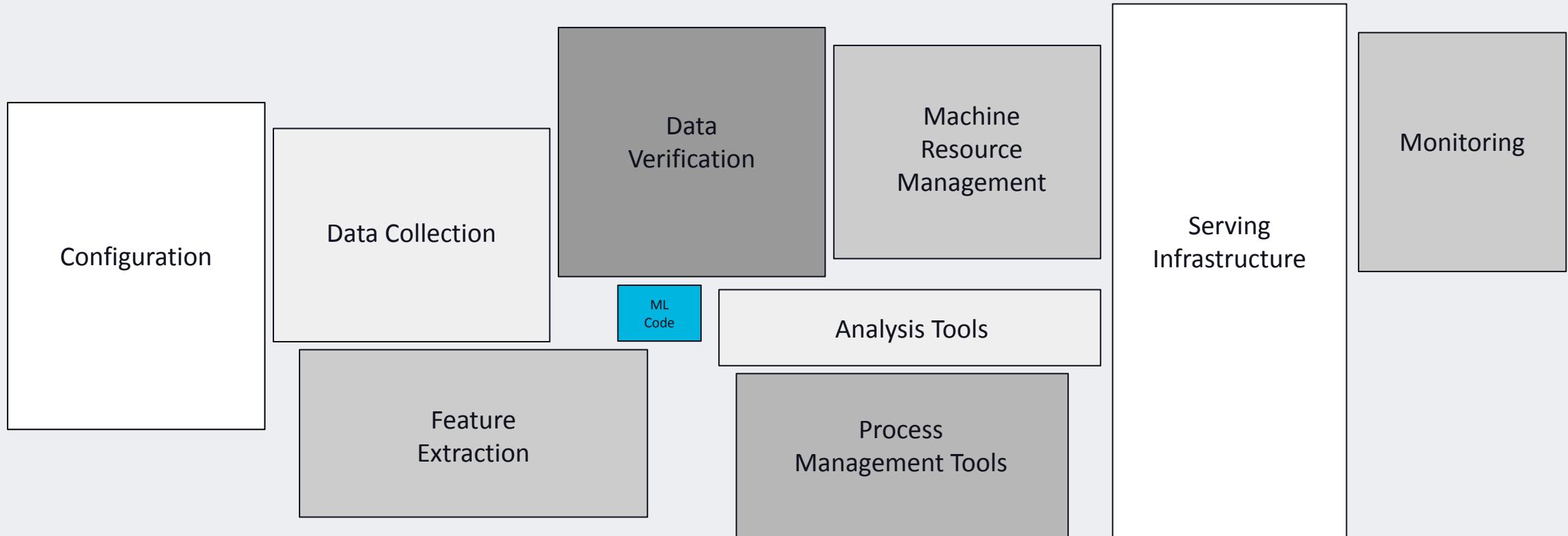


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small green box in the middle. The required surrounding infrastructure is vast and complex.

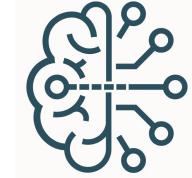
What is MLOps?

MLOps is a set of **processes and automation**

for **managing code, data, and models**

to **improve performance, stability and long-term efficiency** of ML systems

MLOps = DevOps + DataOps + ModelOps



Why should I care about MLOps?

MLOps helps you reduce risk

- Technical risk – poorly performing models, fragile infrastructure
- Compliance risk – violating regulatory or corporate policies

MLOps improves long-term efficiency through automation

- Catch errors before they hit production
- Avoid slow, manual processes

People and process



Business Stakeholder

→ Responsible for business value of the ML solution



Data Engineer

→ Builds data pipelines



Data Scientist

→ Translates business problem; trains, tunes model



ML Engineer

→ Deploys ML model to production

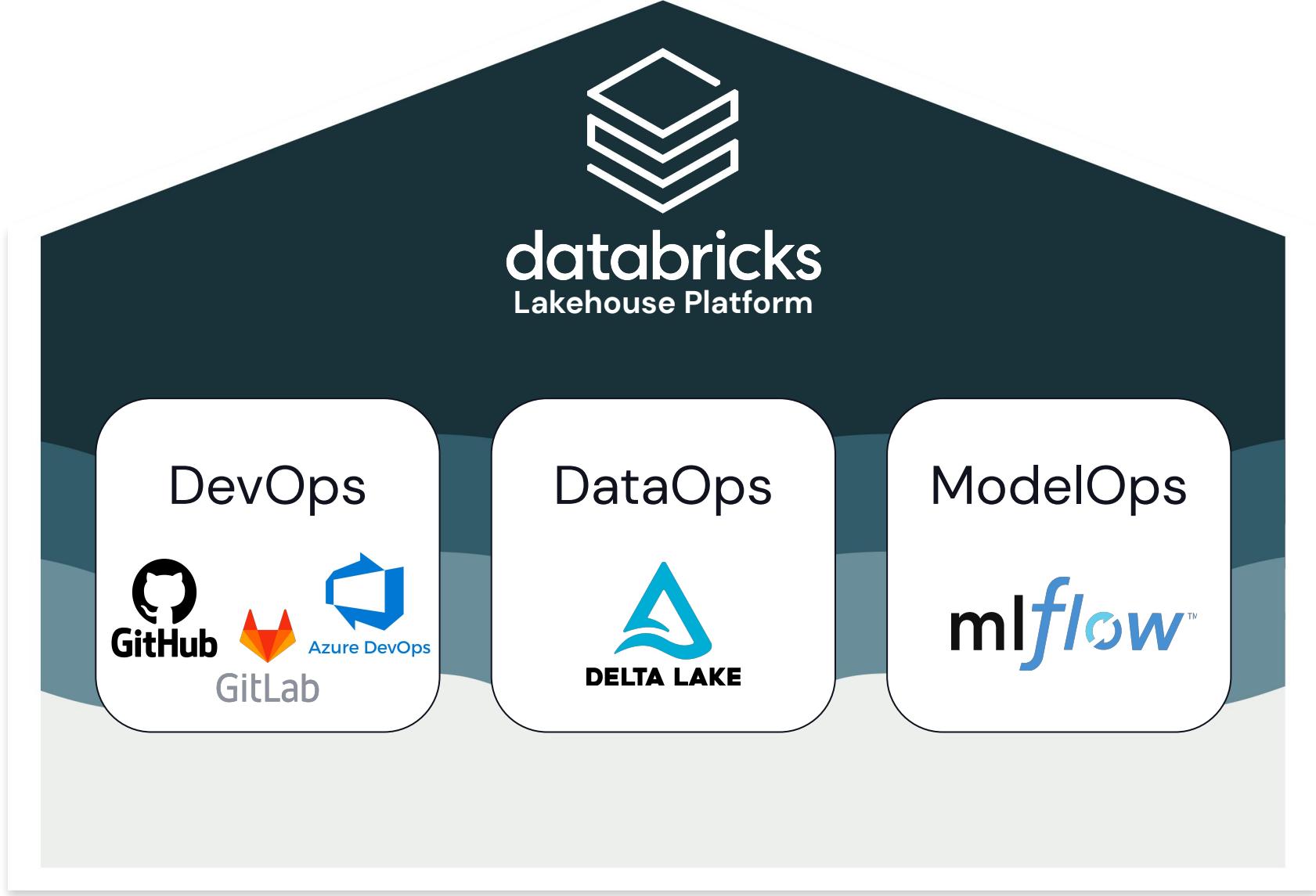


Data Governance Officer

→ Responsible for data governance and compliance

People and process





Data-centric ML platform

Guiding principles



Always keep your business goals in mind



Take a data-centric approach to machine learning



Implement MLOps in a modular fashion



Process should guide automation

Fundamentals of MLOps

- Semantics of dev, staging and prod
- Deployment patterns

Semantics of dev, staging, and prod

ML Workflow Assets:

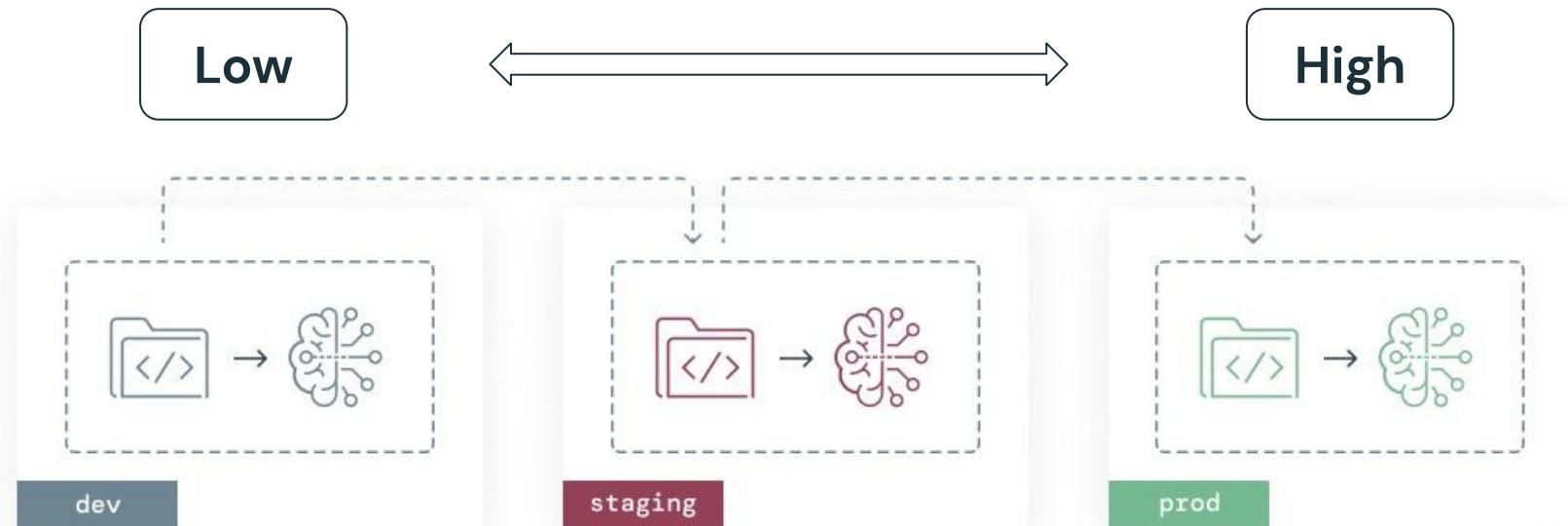


Assets need to be:



Dev vs. staging vs. prod

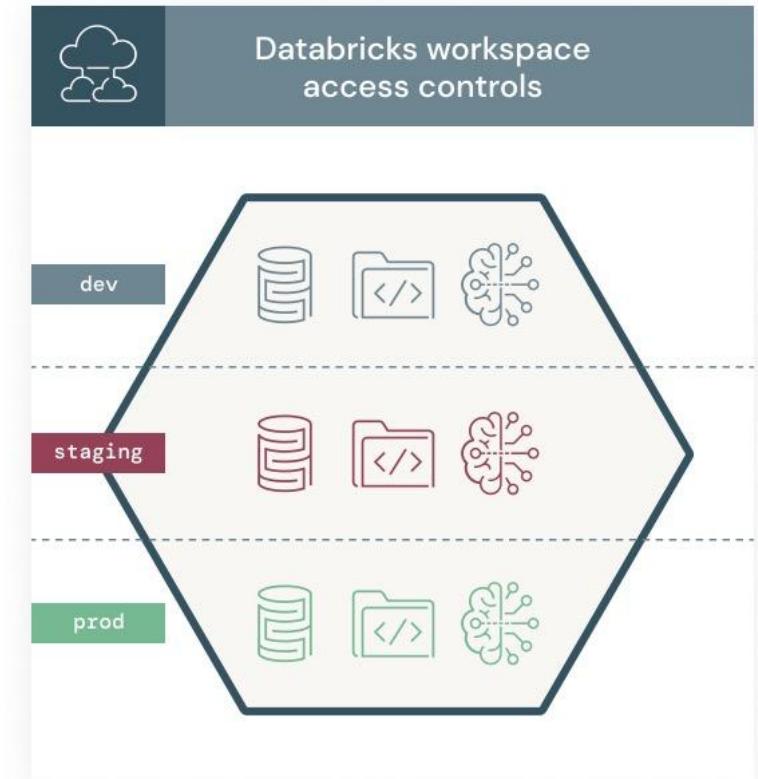
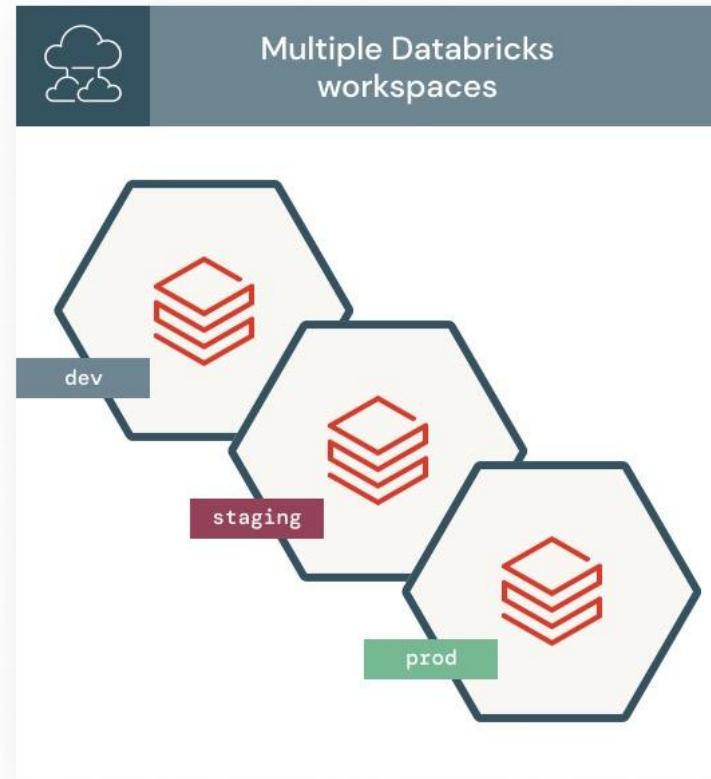
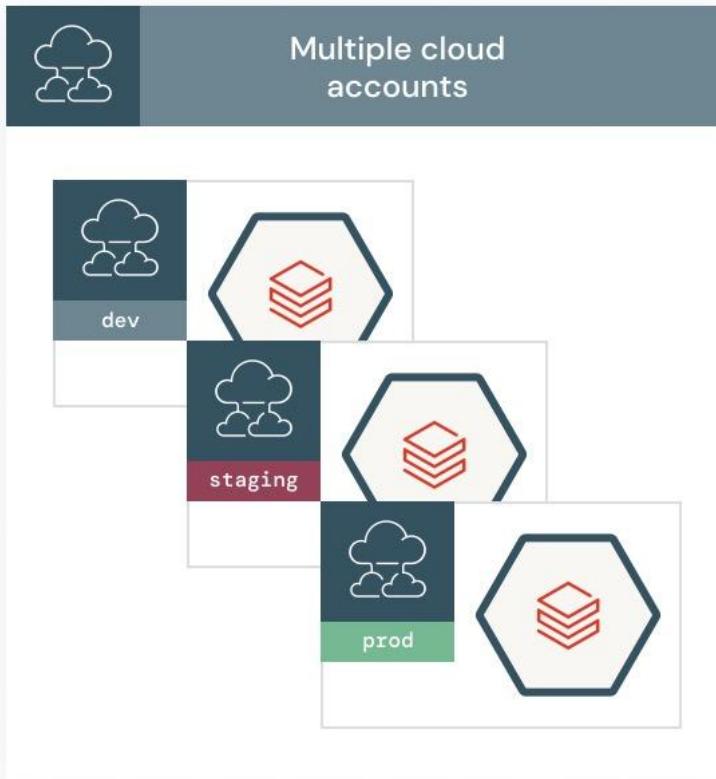
Level of trust, quality
and testing:



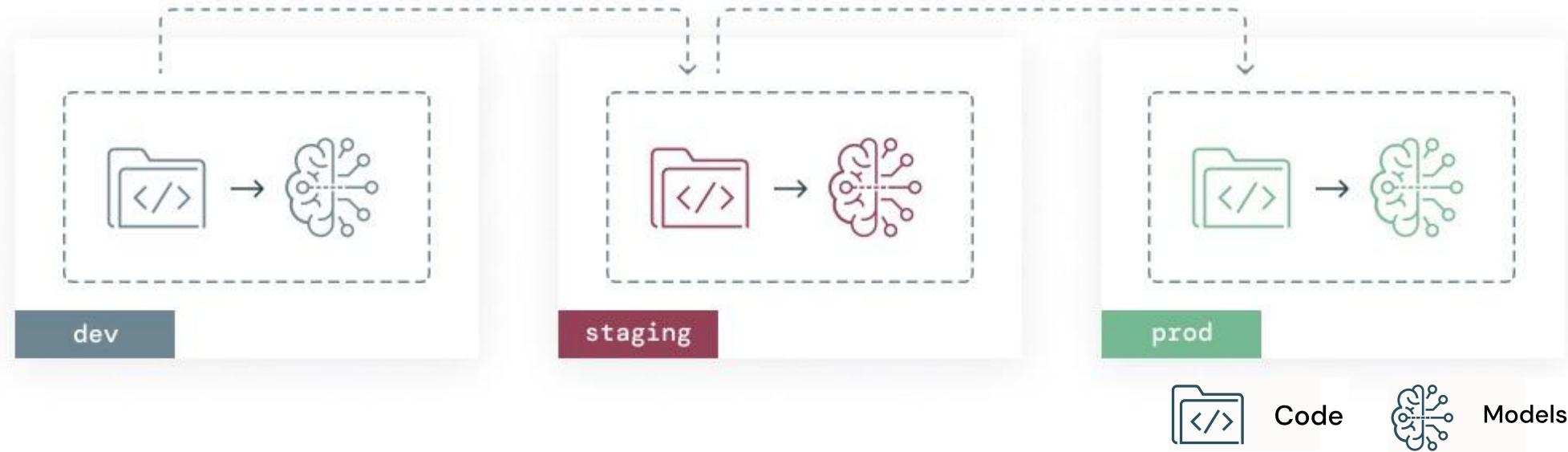
Openness of access:



Environments: dev vs. staging vs. prod



Model vs. code lifecycles



Model and code lifecycles often operate asynchronously

- Weekly fraud detection model
 - Model update; no code change
- Computer vision model, large language model
 - No model update; code change

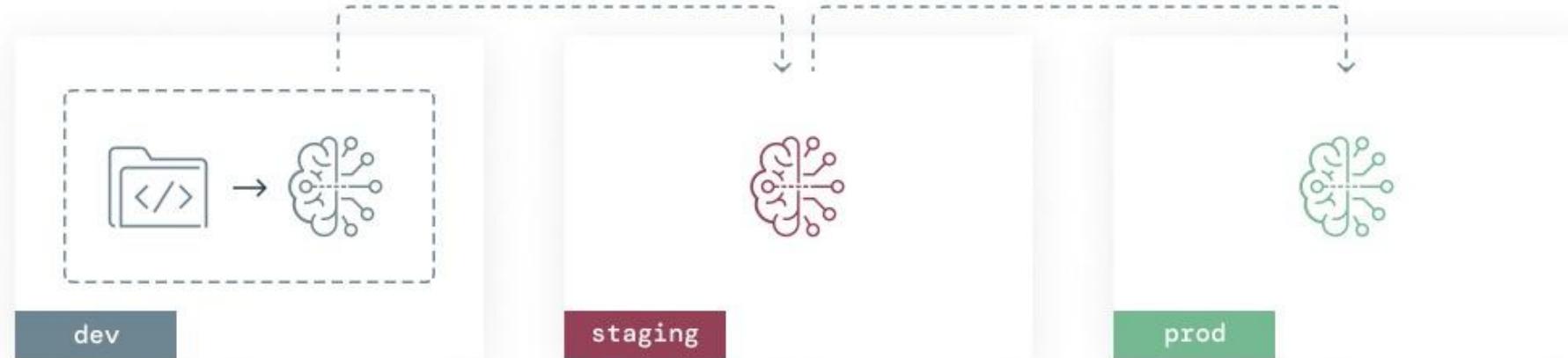
→ mlflow™

Dev vs. staging vs. prod: managing assets

| ASSET | SEMANTICS | SEPARATED BY |
|------------------------|--|--|
| Execution environments | Labeled according to where development, testing and connections with production systems happen | Cloud provider and Databricks Workspace access controls |
| Models | Labeled according to model lifecycle phase | MLflow access controls or cloud storage permissions |
| Data | Labeled according to its origin in dev, staging or prod execution environments | Table access controls or cloud storage permissions |
| Code | Labeled according to software development lifecycle phase | Git repository branches |

ML deployment patterns

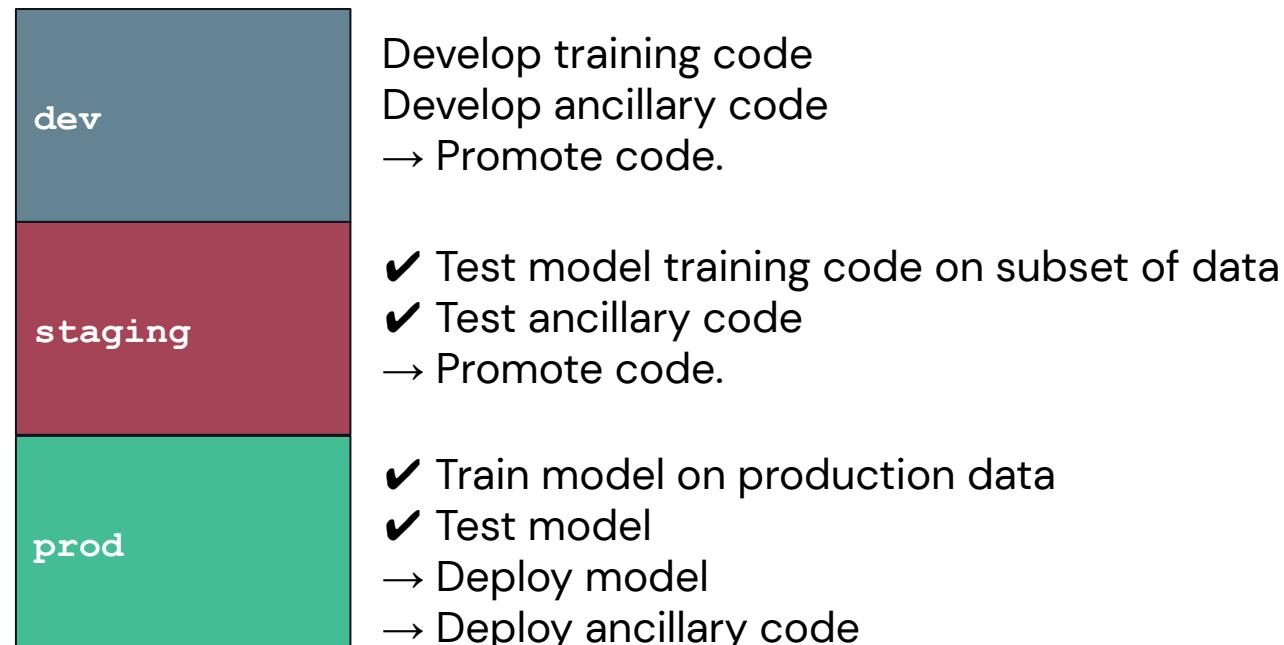
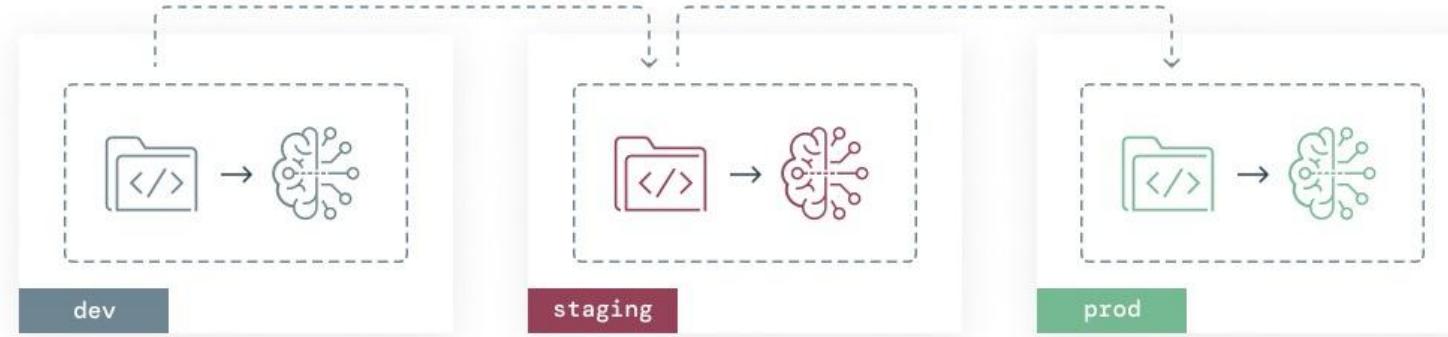
Deploy Models



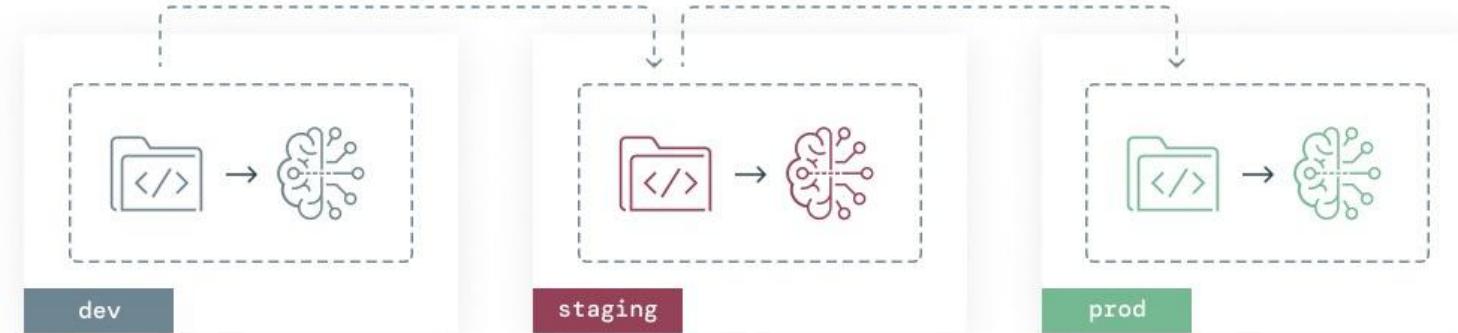
Deploy Code



Deploy code process



Benefits of deploy code

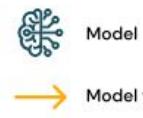
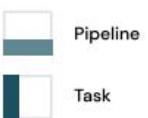
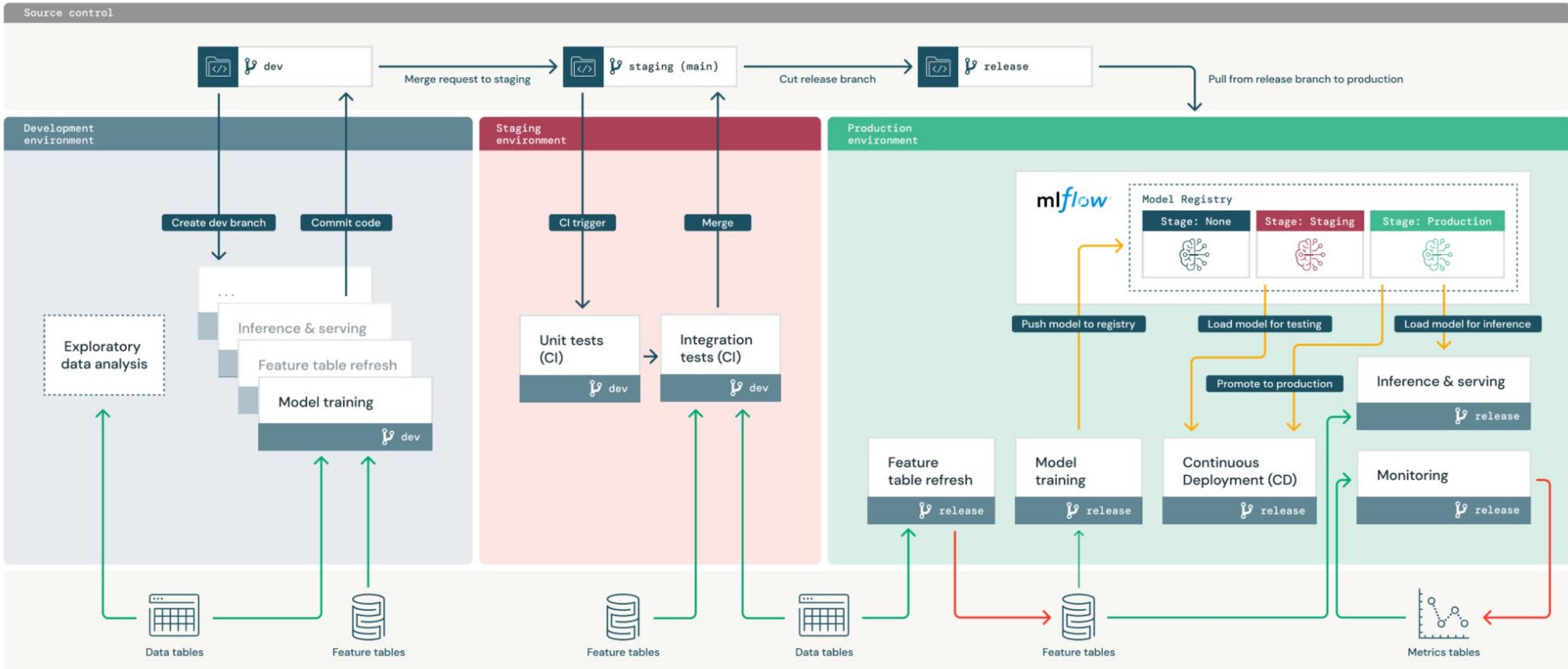


| | |
|------------------------------------|--|
| Automation | ↑ Supports automated retraining in locked-down env. |
| Data access control | ↑ Only prod env needs read access to prod training data. |
| Reproducible models | ↑ Eng control over training env, which helps to simplify reproducibility. |
| Support for large projects | ↑ This pattern forces the DS team to use modular code and iterative testing, which helps with coordination and development in larger projects. |
| Data science familiarity | ↓ DS team must learn to write & hand off modular code to Eng. |
| Eng setup & maintenance | ↓ Requires CI/CD infra for unit and integration tests, even for one-off models. |

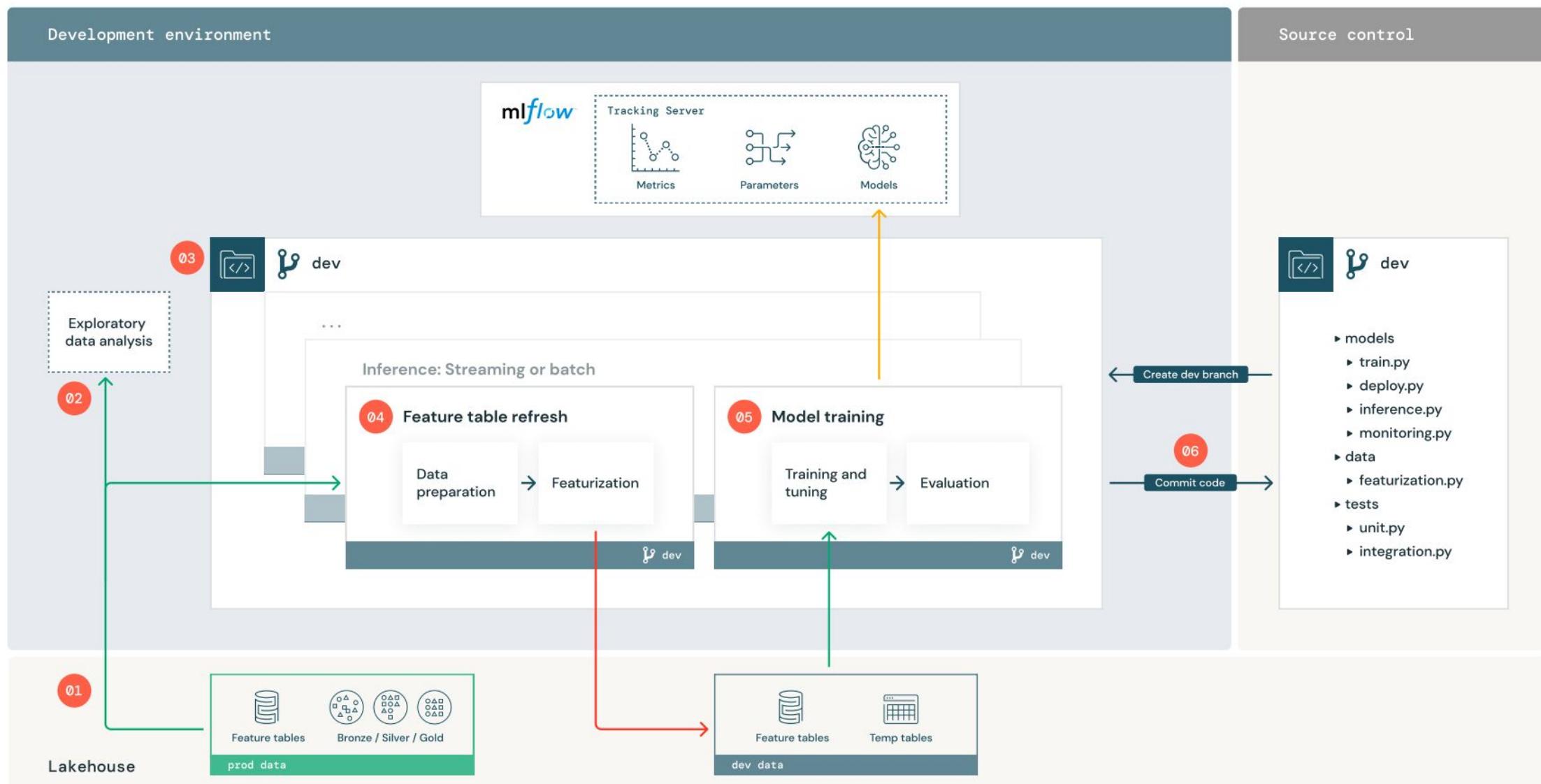
Reference Architectures

- Multi-environment view
- Dev
- Staging
- Production

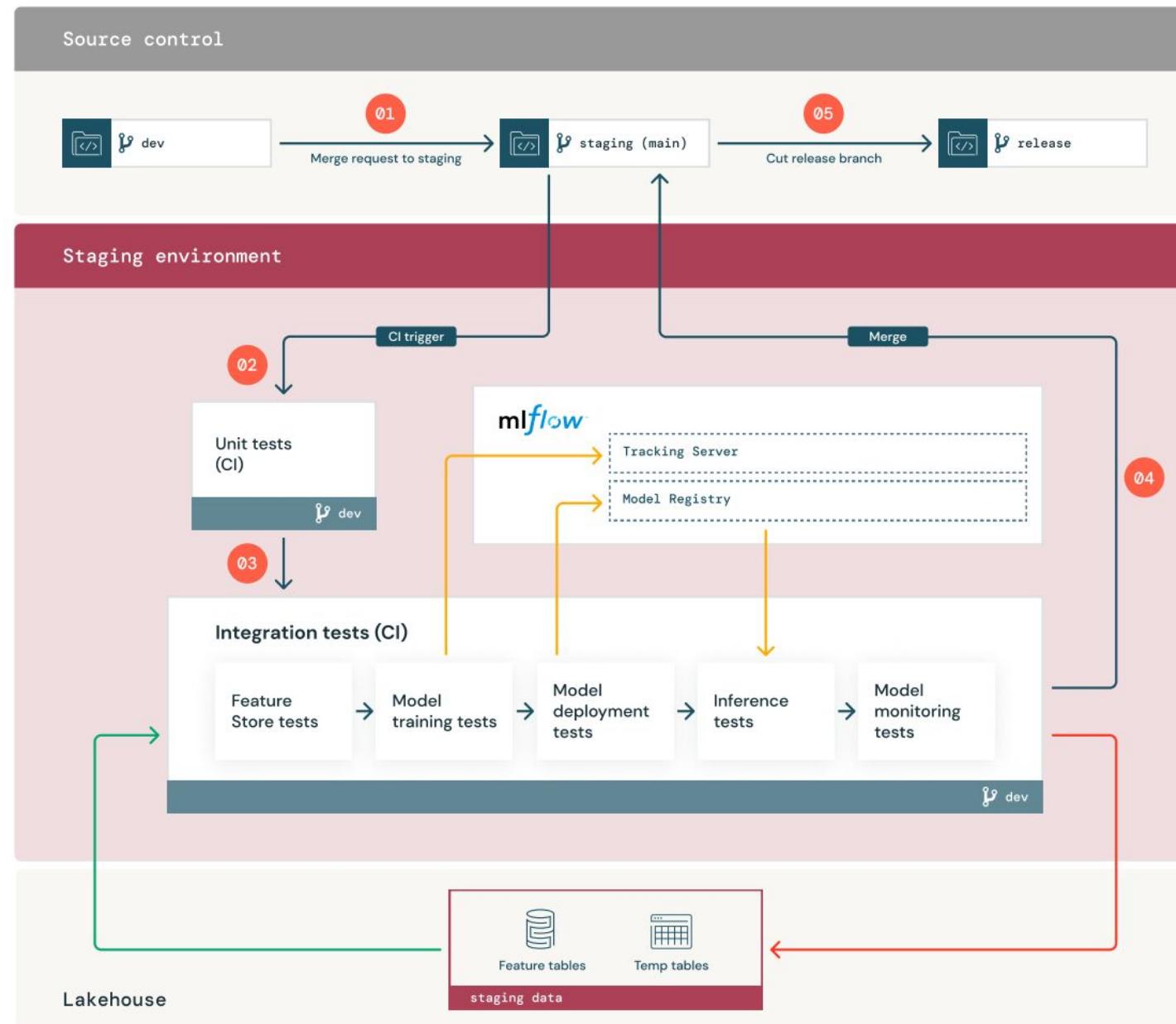
Overview

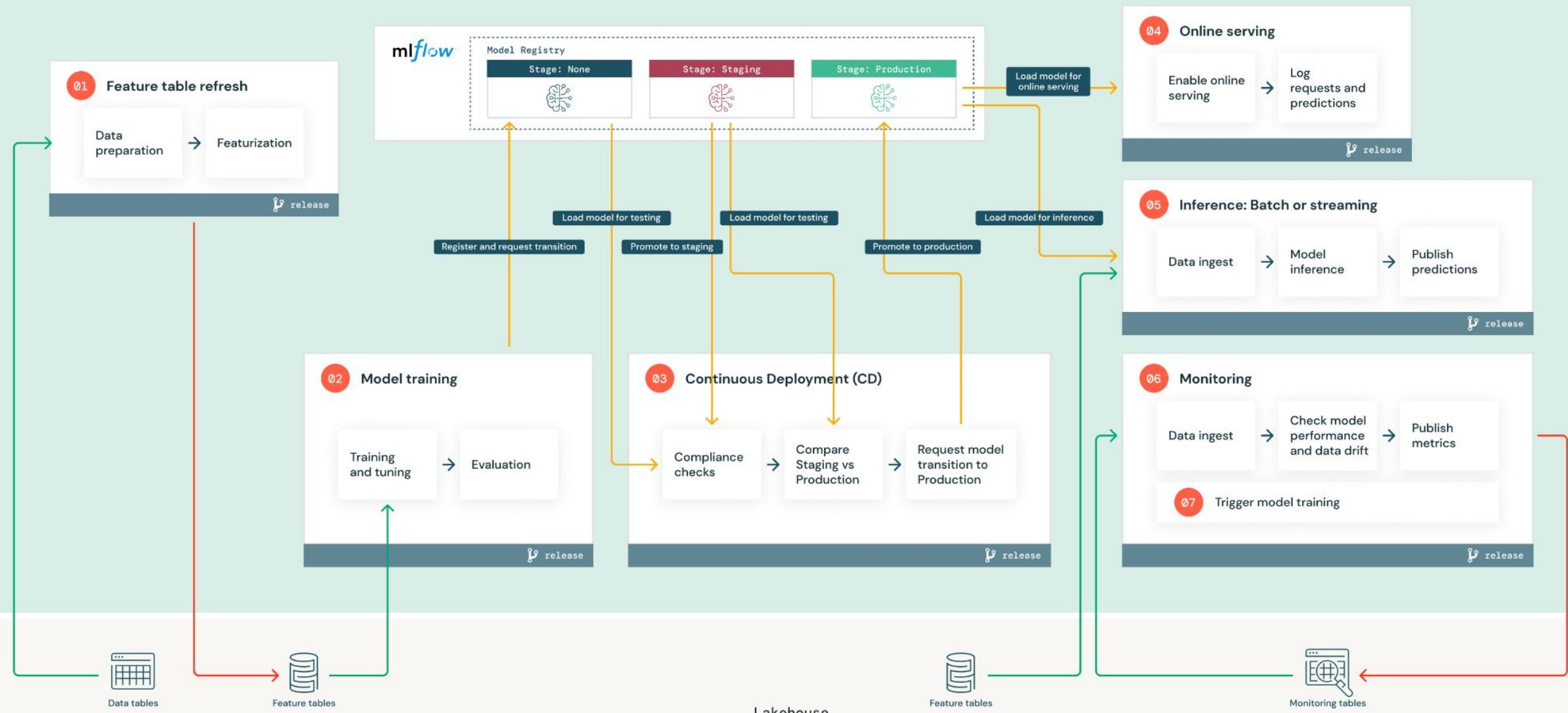


Dev



Staging





Part II

Outline

Part I

- Introduction
 - People and process
 - Goals of MLOps
 - Guiding principles
- Fundamentals
 - Semantics of dev, staging, prod
 - Moving to production
- Reference Architectures

Part II

- The use case
- Overview of pipelines
- Tooling
 - DevOps
 - DataOps
 - ModelOps
- Project Structure
- Demo
- Getting started

The use case

Churn prediction problem

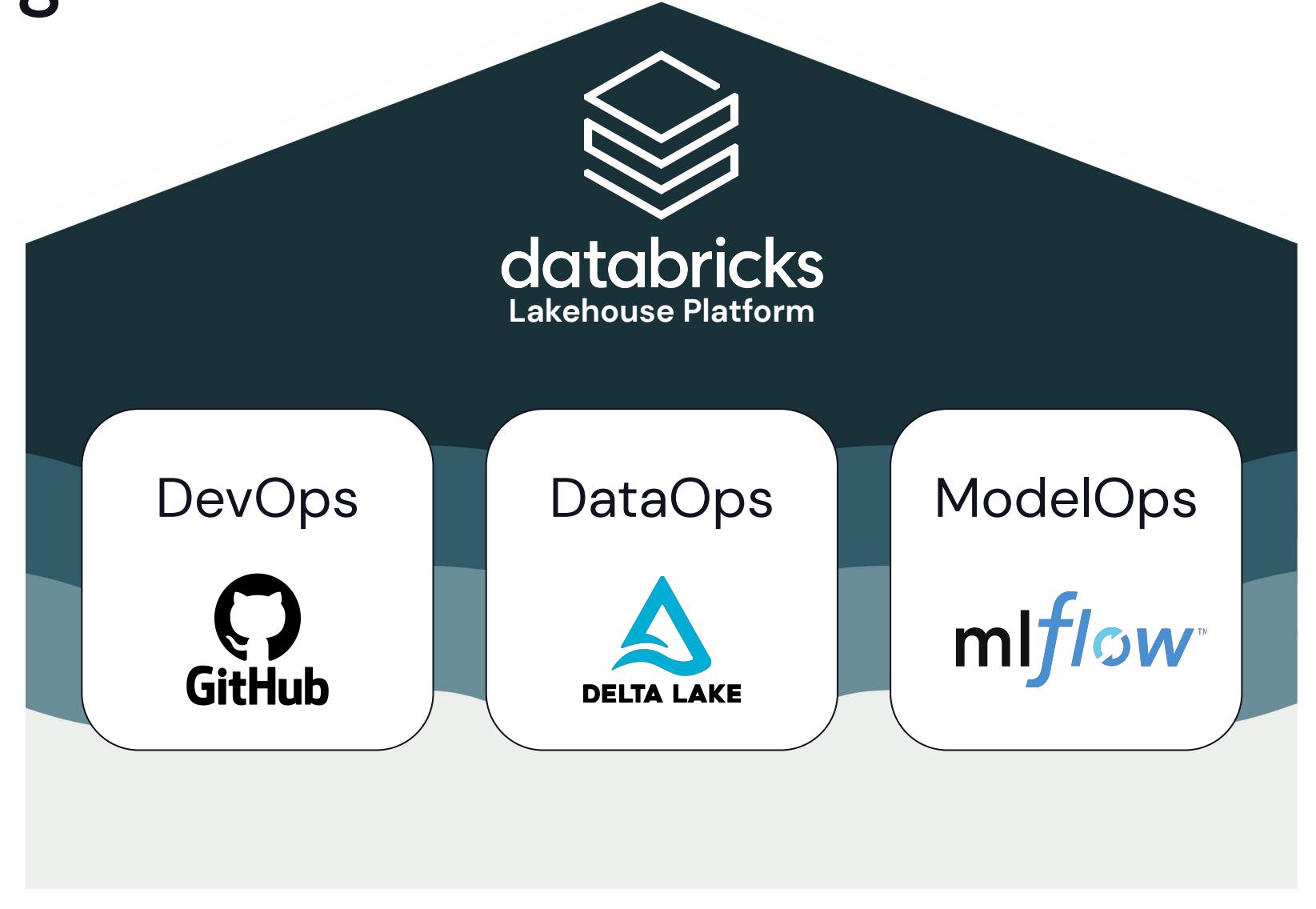
IBM telco customer churn dataset

- Fictional telco company providing phone and internet services
- 7,043 customers
- Churn label column indicates whether customer has left within the last month
- Demographic, and contract-based information

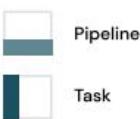
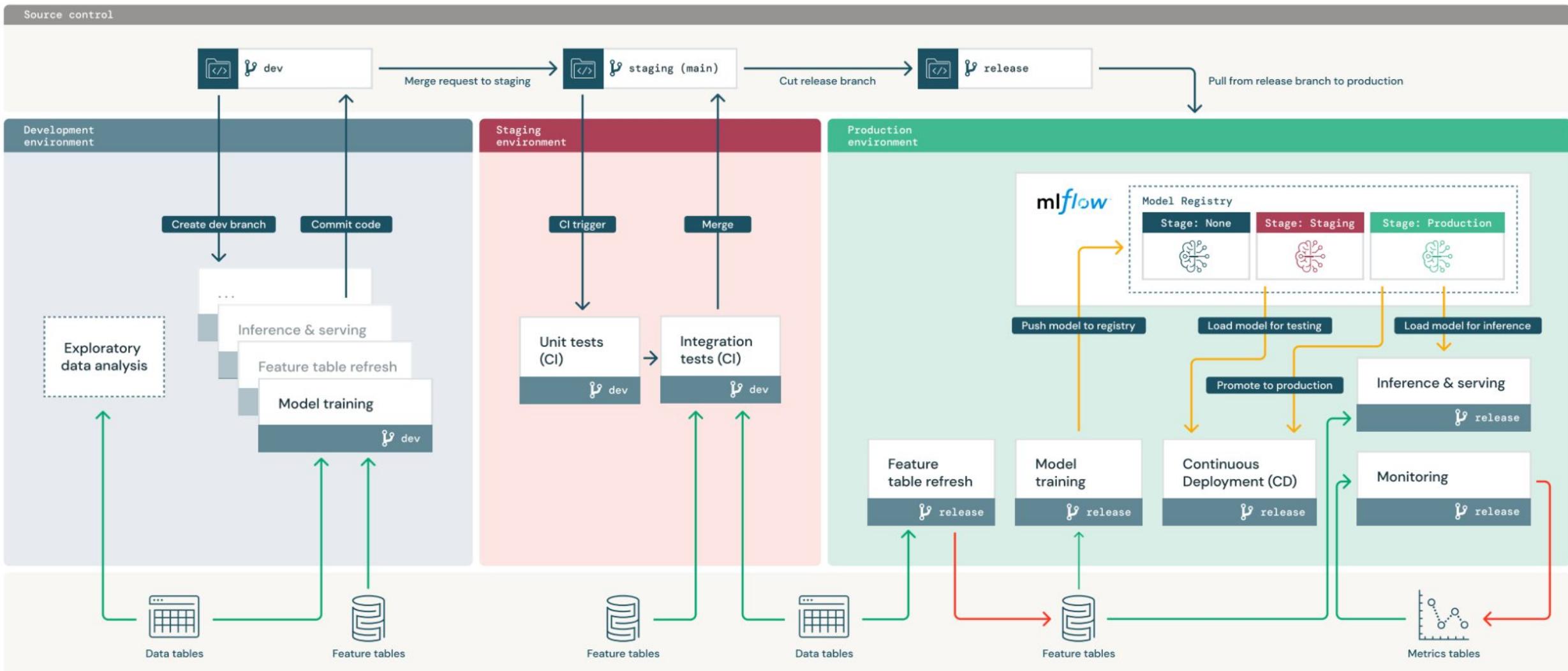
Goal

Create MLOps workflow to continuously deliver
a churn prediction model to production

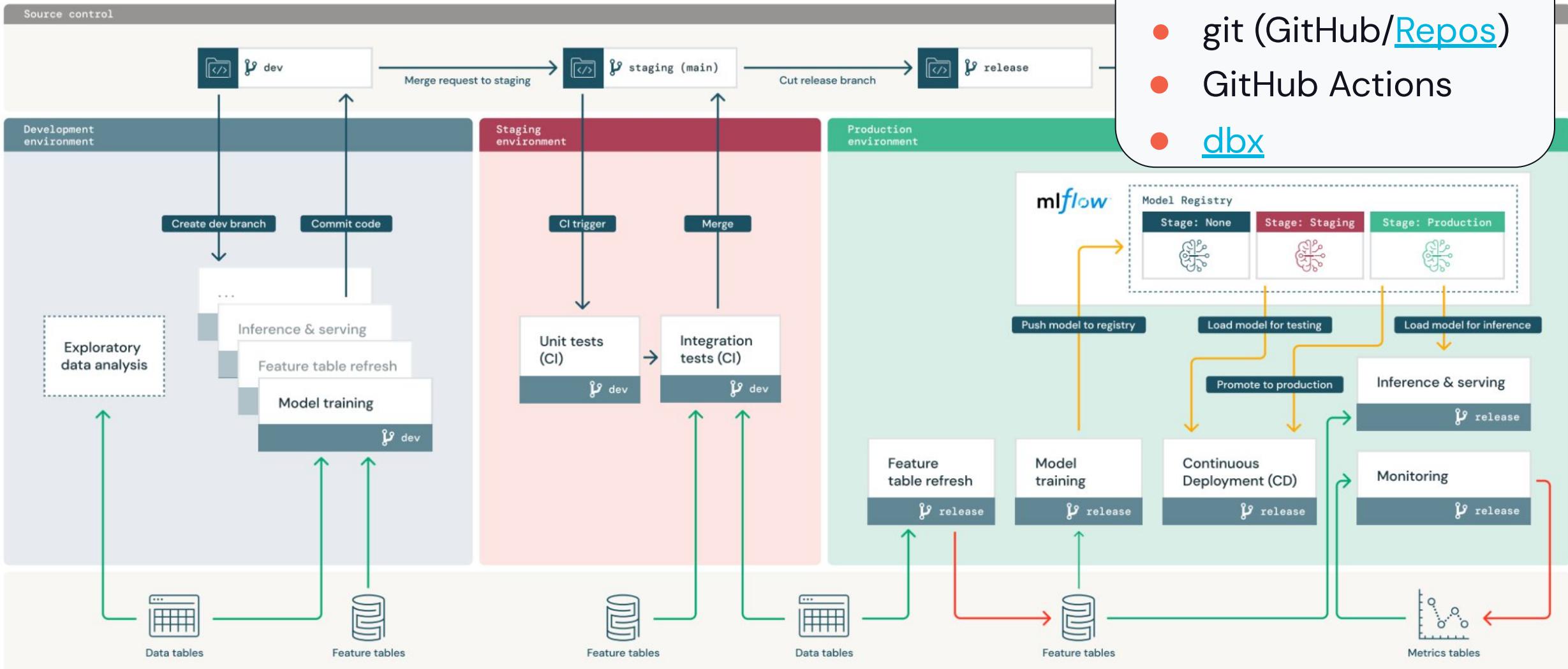
Tooling



Tooling: Architecture Perspective



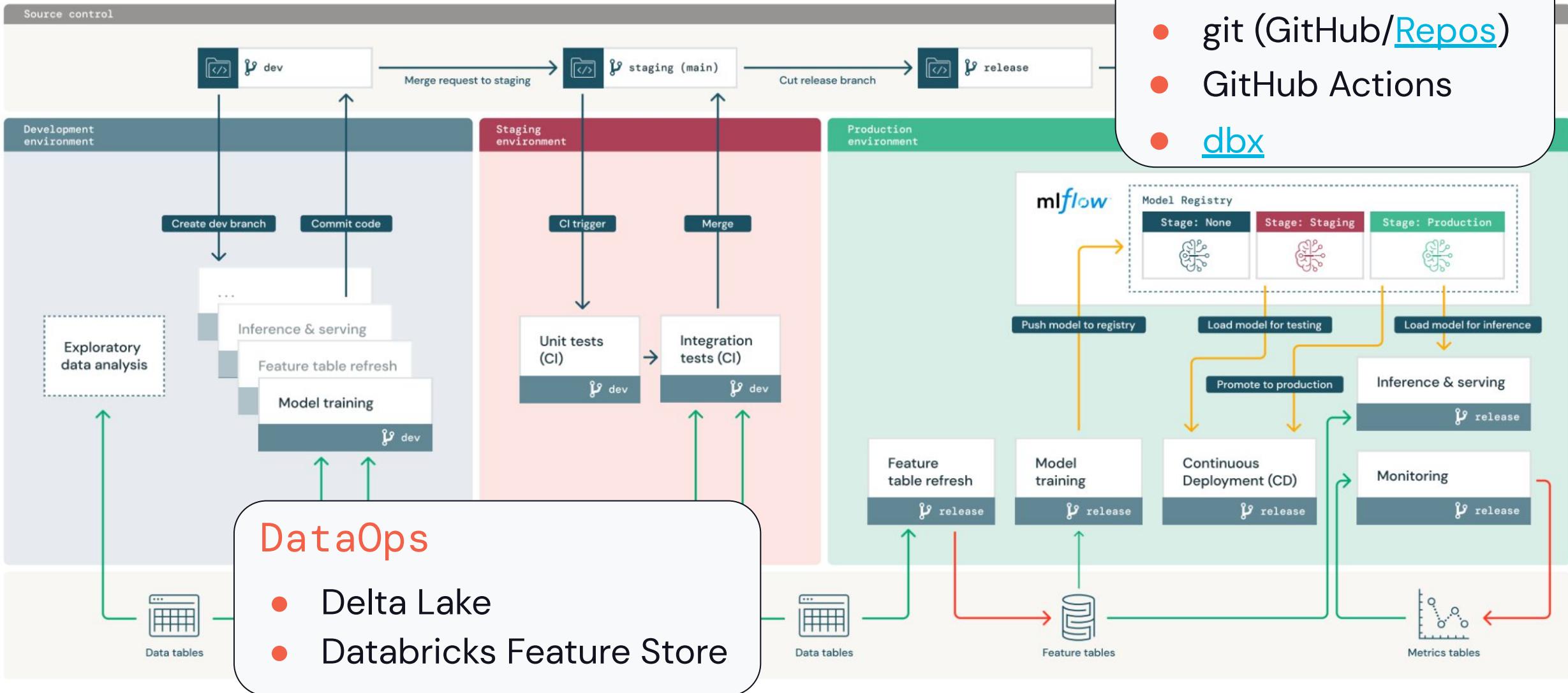
Tooling: Architecture Perspective



DevOps

- git (GitHub/[Repos](#))
- GitHub Actions
- [dbx](#)

Tooling: Architecture Perspective

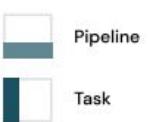


DevOps

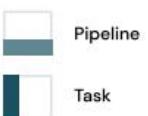
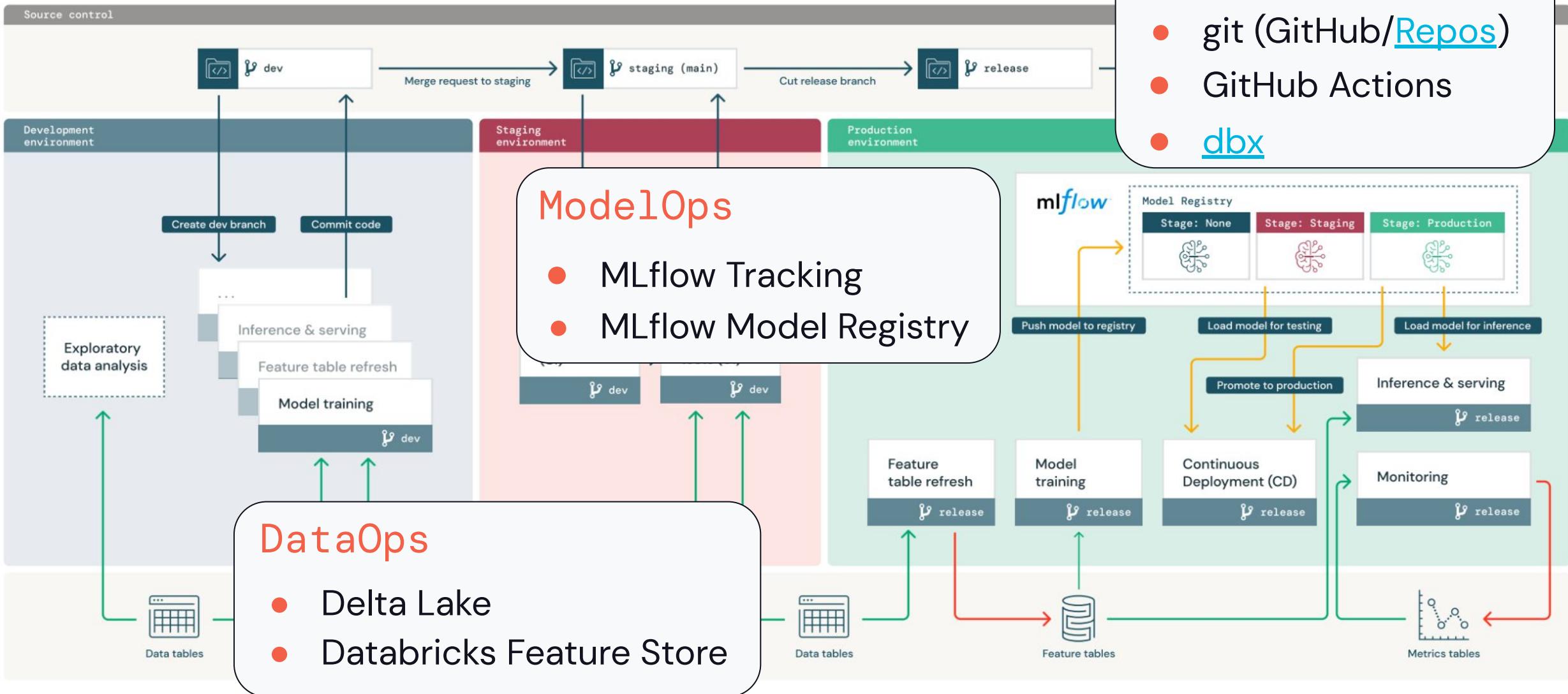
- git (GitHub/[Repos](#))
- GitHub Actions
- [dbx](#)

DataOps

- Delta Lake
- Databricks Feature Store



Tooling: Architecture Perspective



Tooling

DevOps

- dbx
 - Simplifies IDE & Databricks integration
 - Job deployment
 - Python Template
- GitHub Actions
 - Automate CI/CD workflows

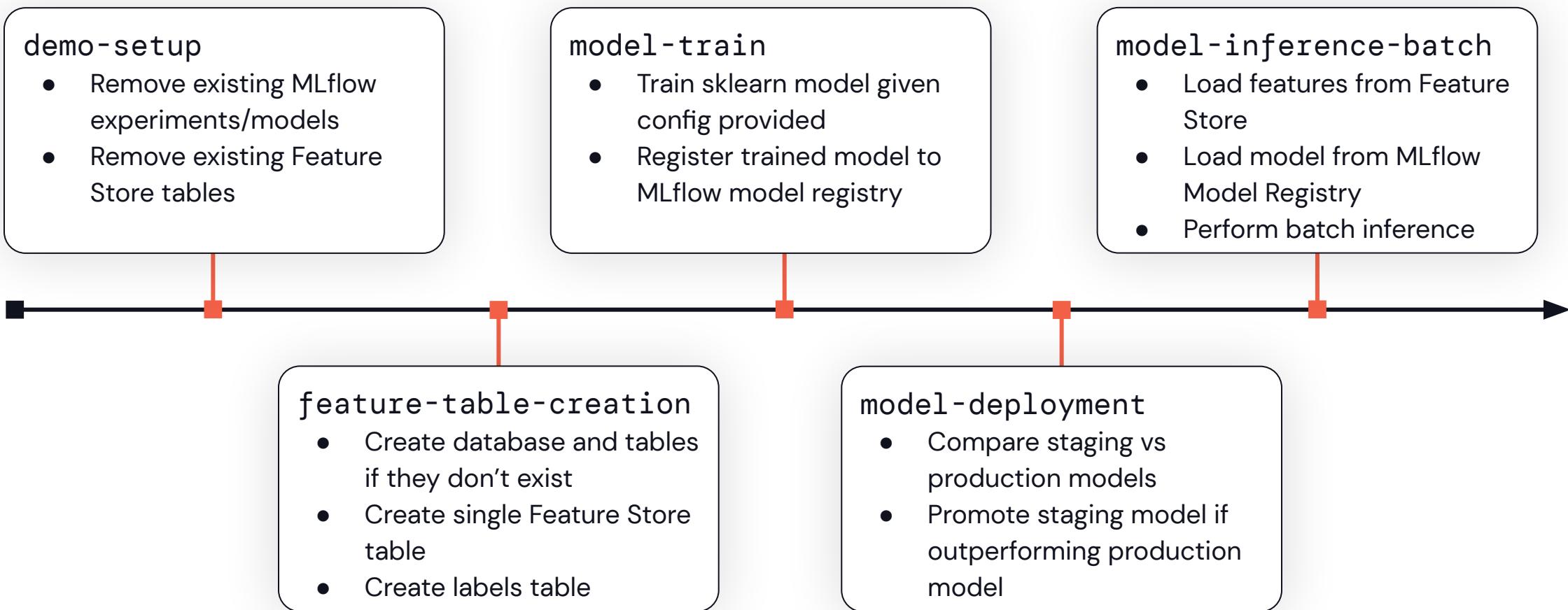
DataOps

- Delta Lake
 - Used as raw data input
- Databricks Feature Store
 - Centralized repository of features
 - Lineage tracking

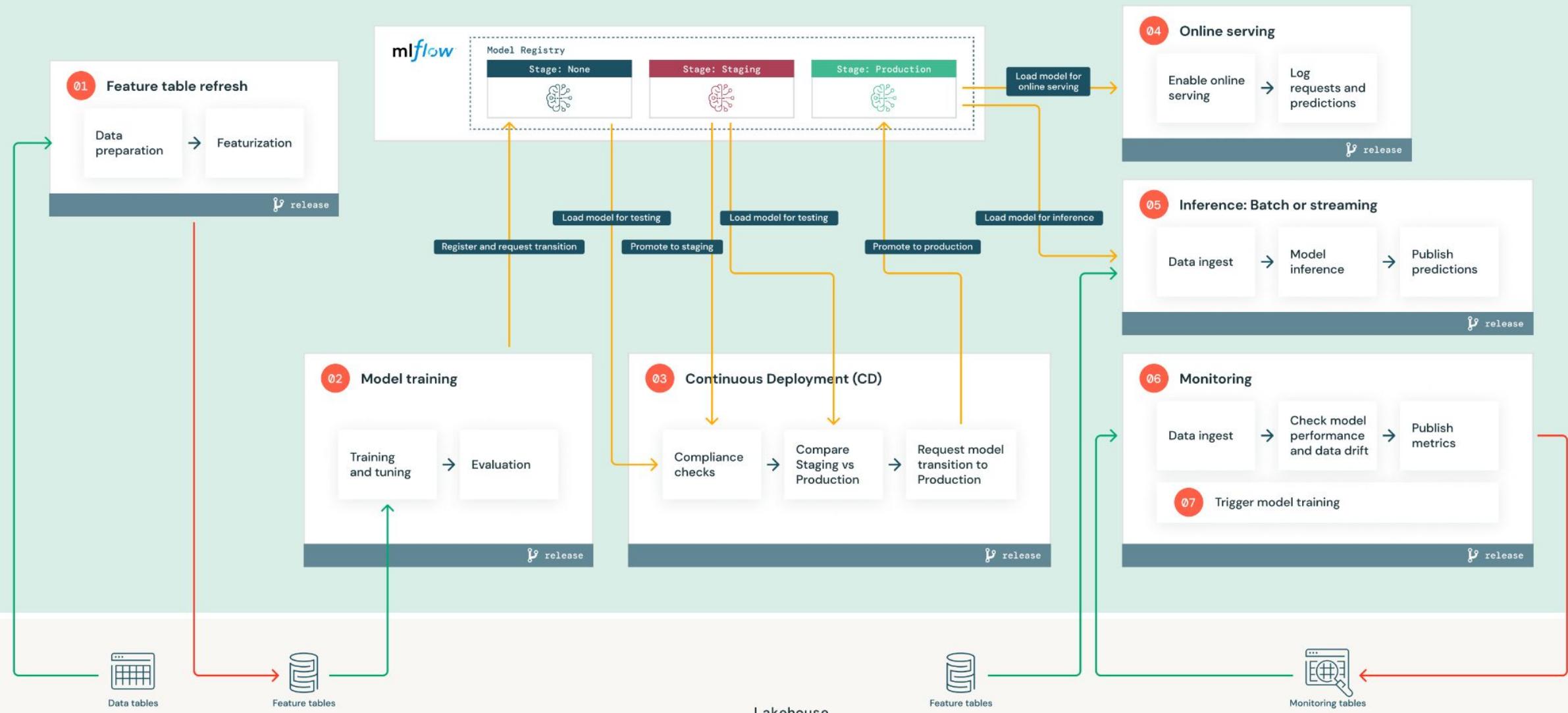
ModelOps

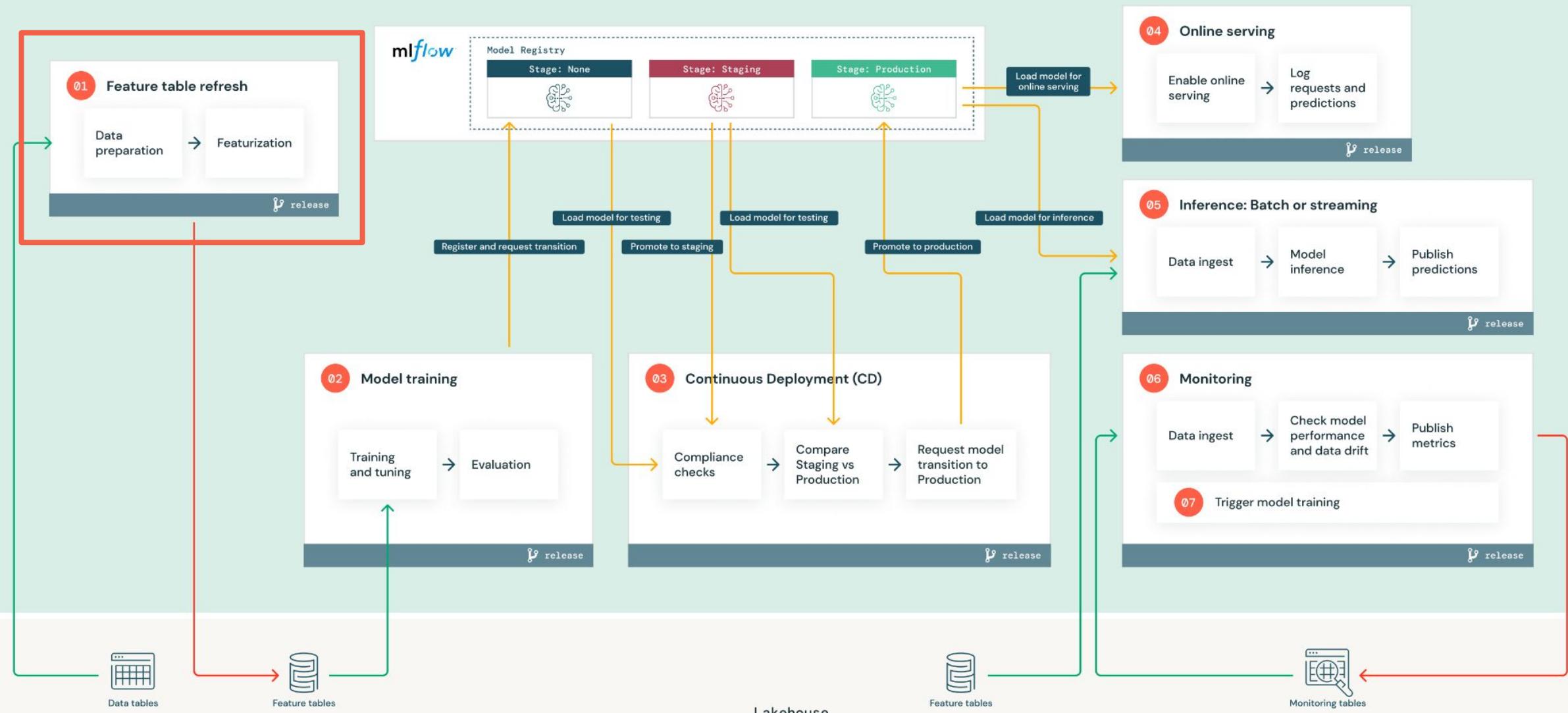
- MLflow Tracking
 - Track params, metrics and model artifacts
- MLflow Model Registry
 - Manage model lifecycle

Overview of Pipelines



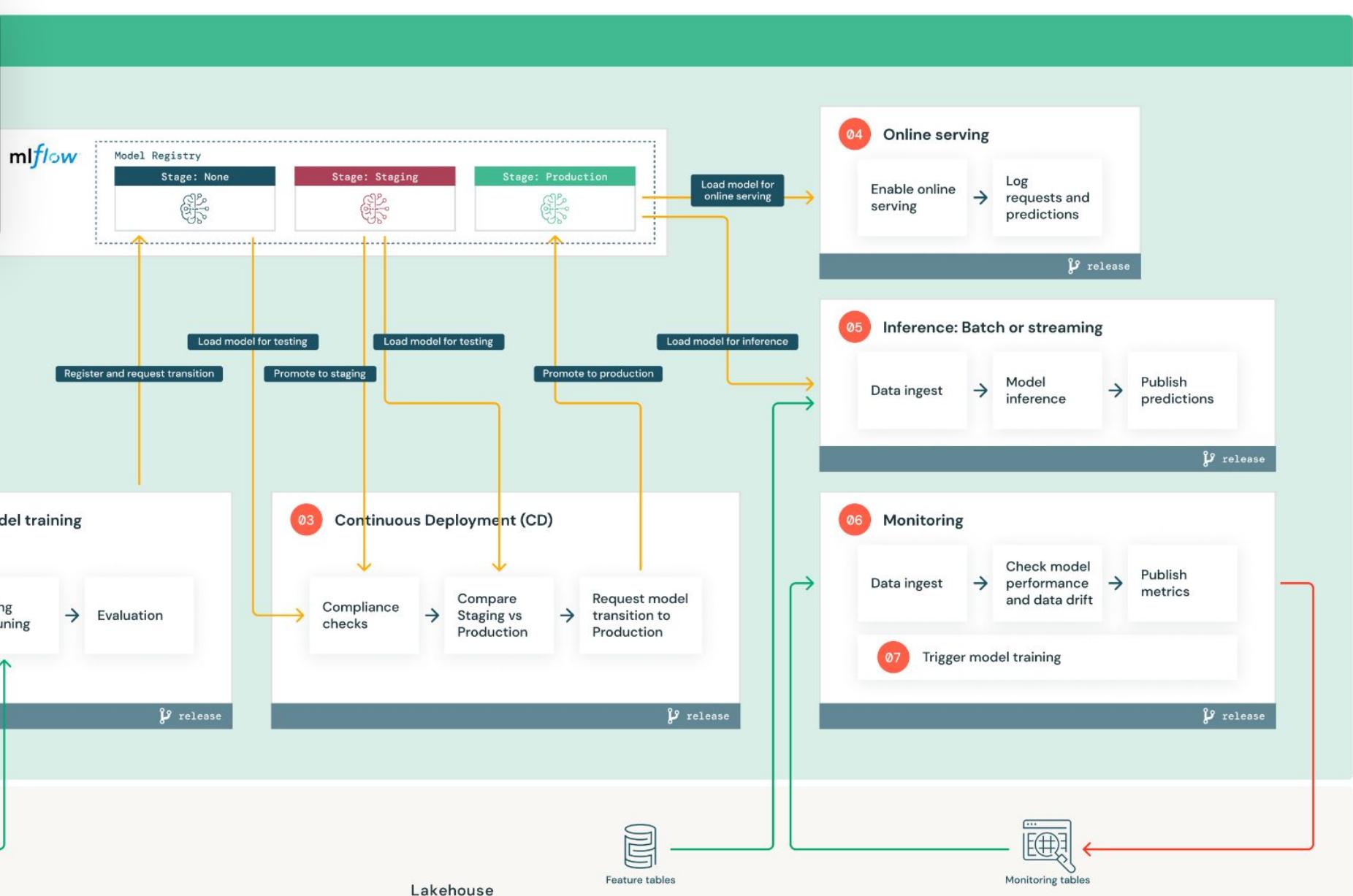
Overview of pipelines





demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

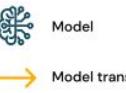
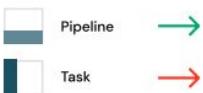
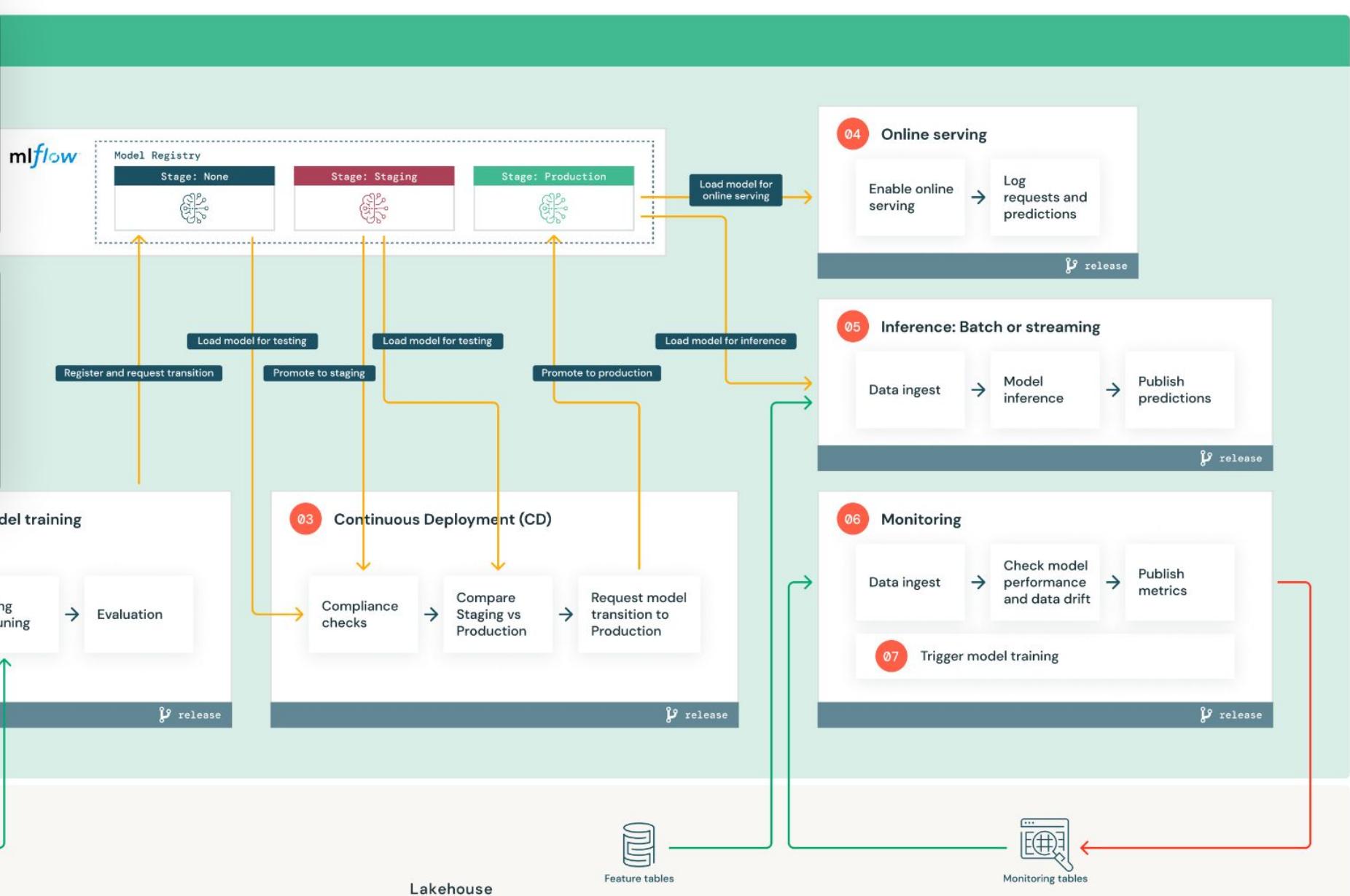


demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table

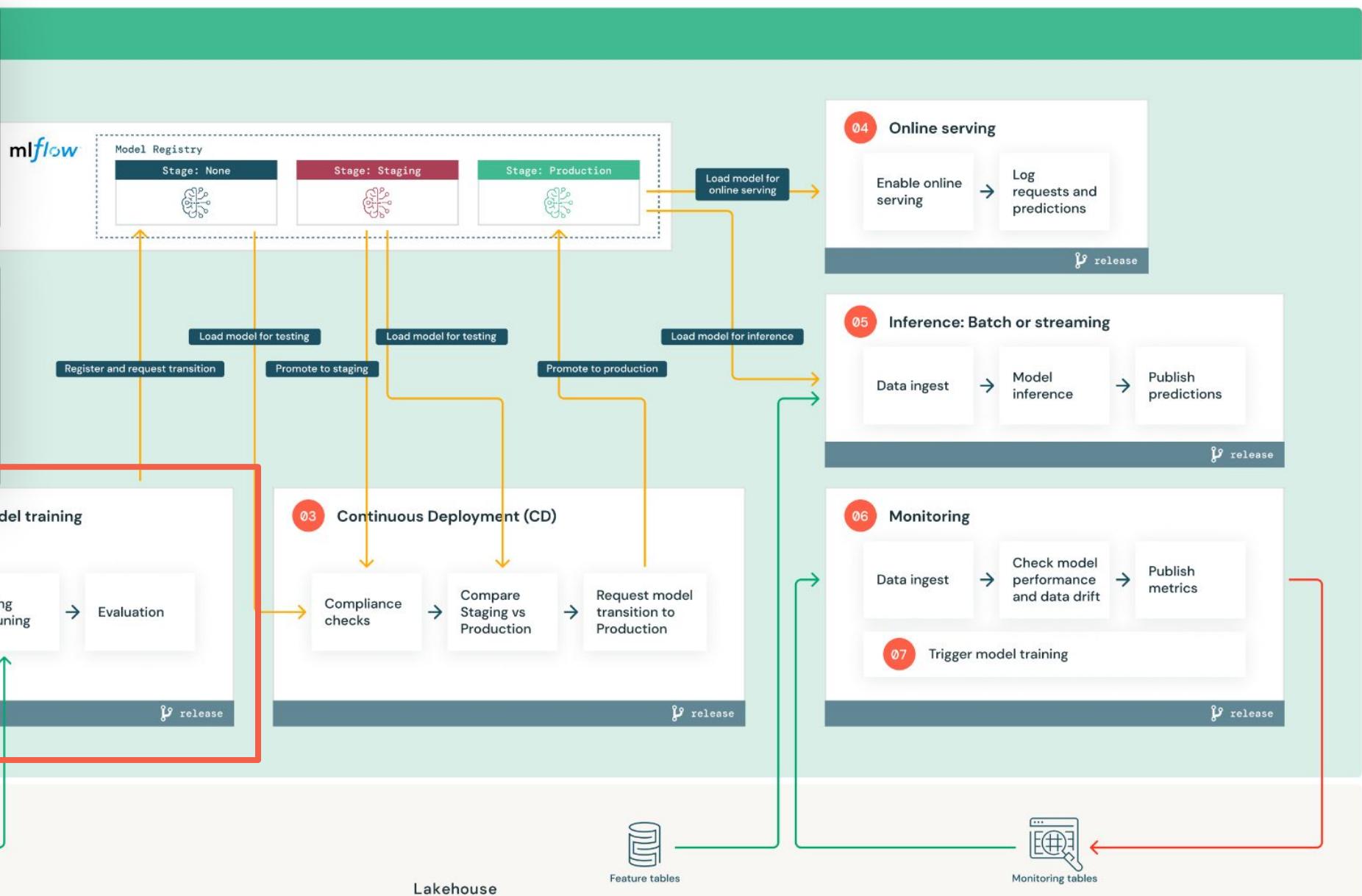


demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table



demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table

model-train

- Train sklearn model given config provided
- Register trained model to MLflow model registry

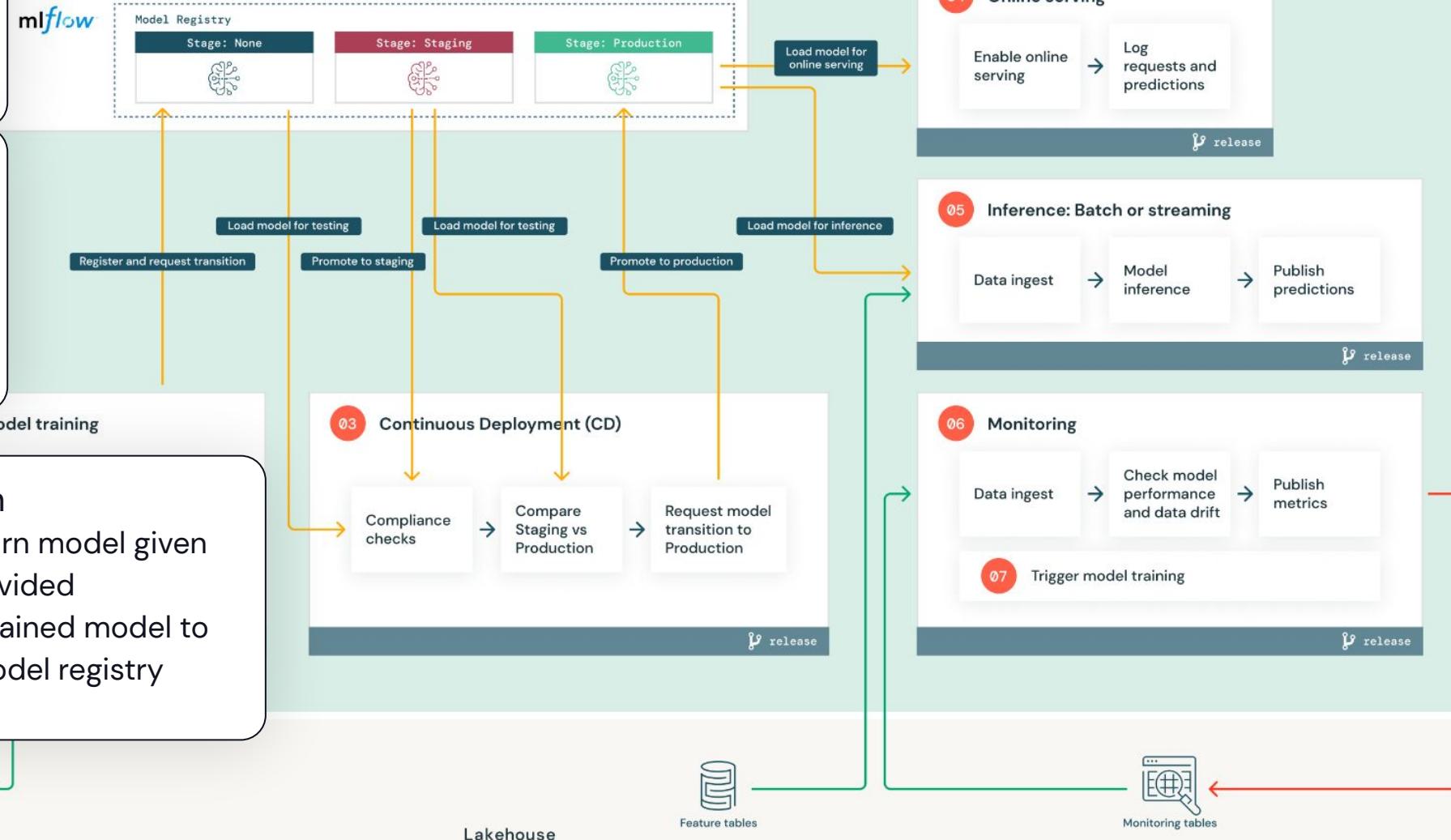


Data tables



Feature tables

mlflow



Pipeline



Reads



Model



Repo



Task



Writes



Model transition



Branch

demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table

model-train

- Train sklearn model given config provided
- Register trained model to MLflow model registry

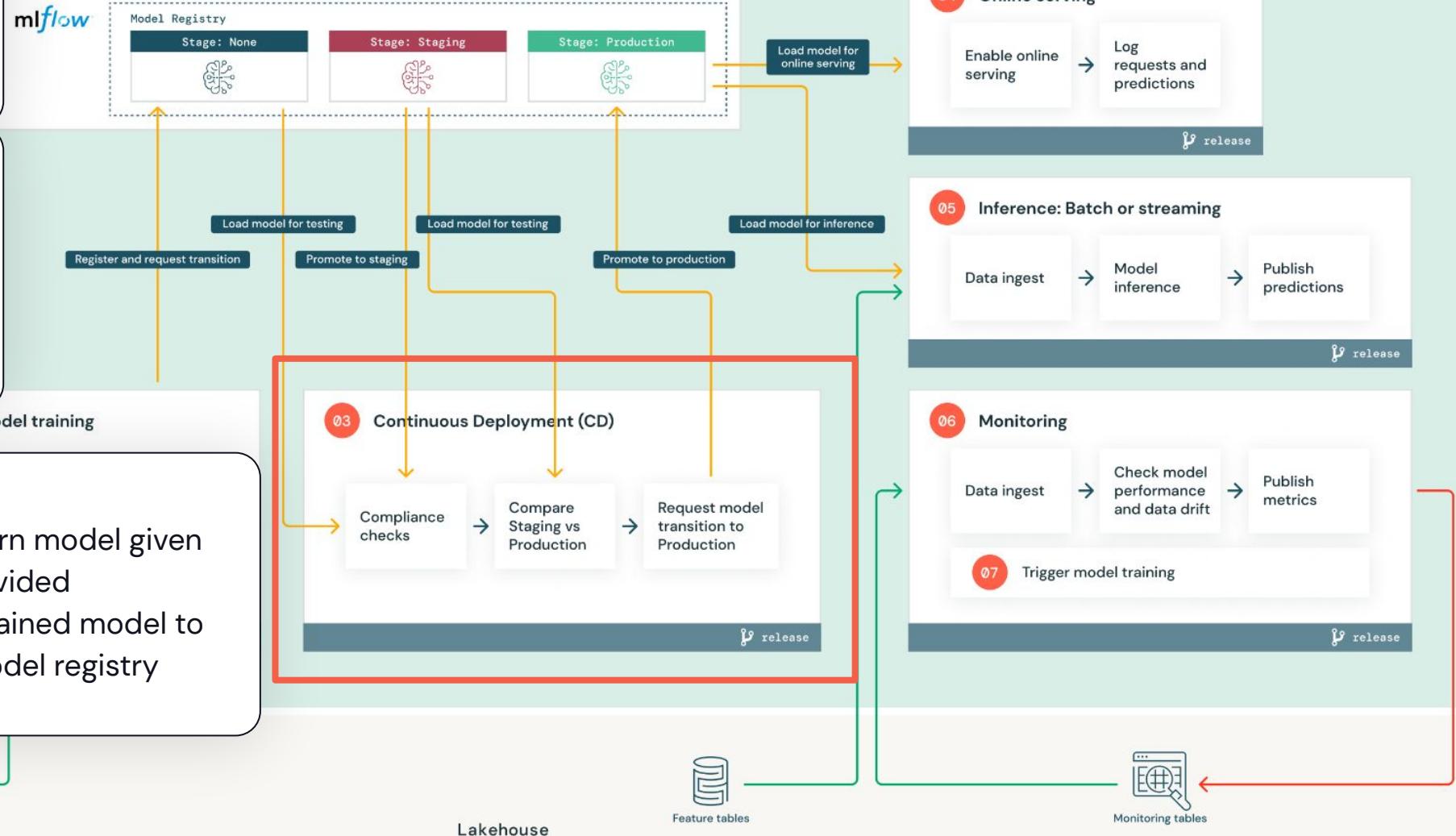


Data tables



Feature tables

mlflow



Pipeline



Reads



Model



Repo



Task



Writes



Model transition



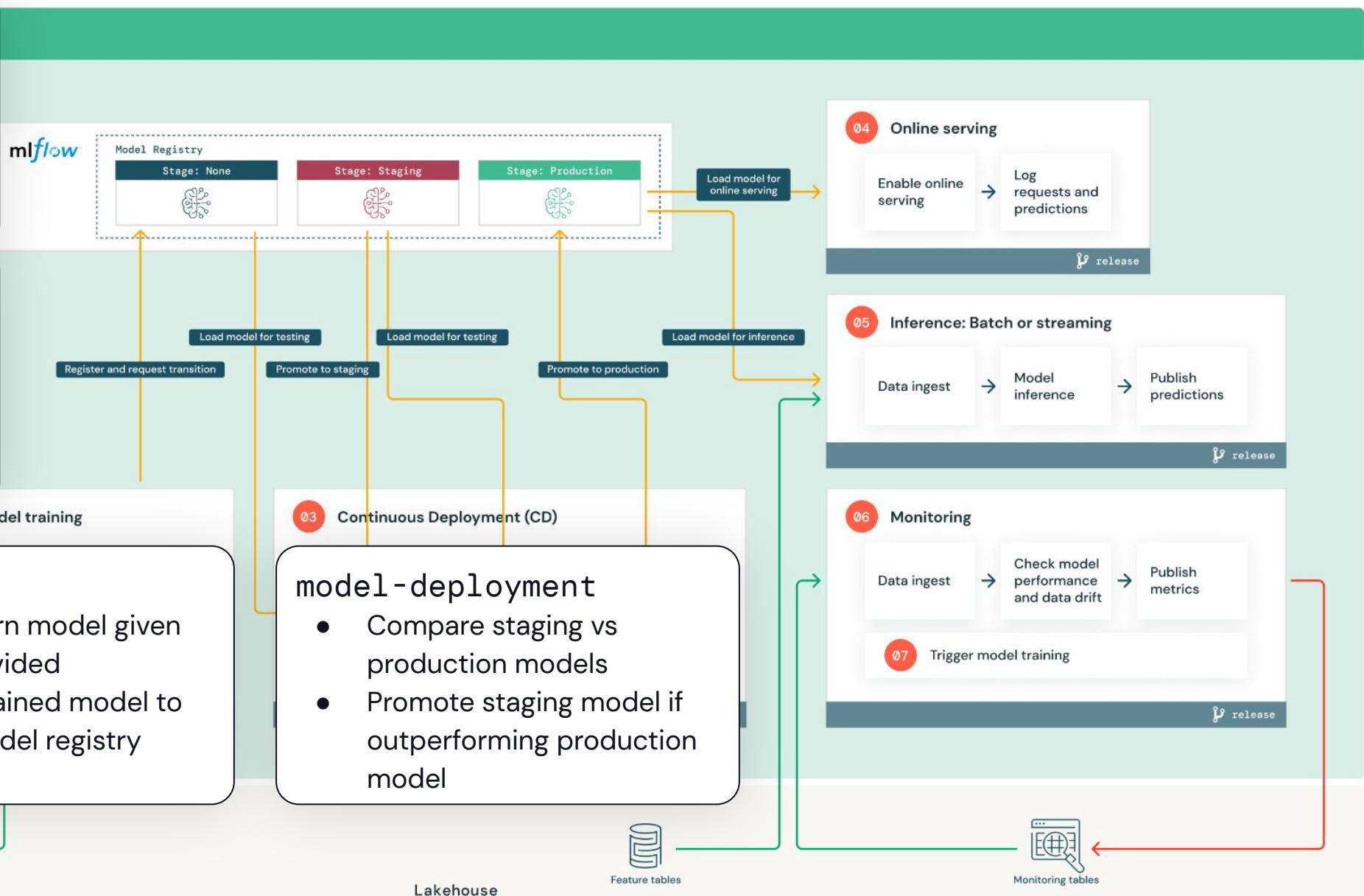
Branch

demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table

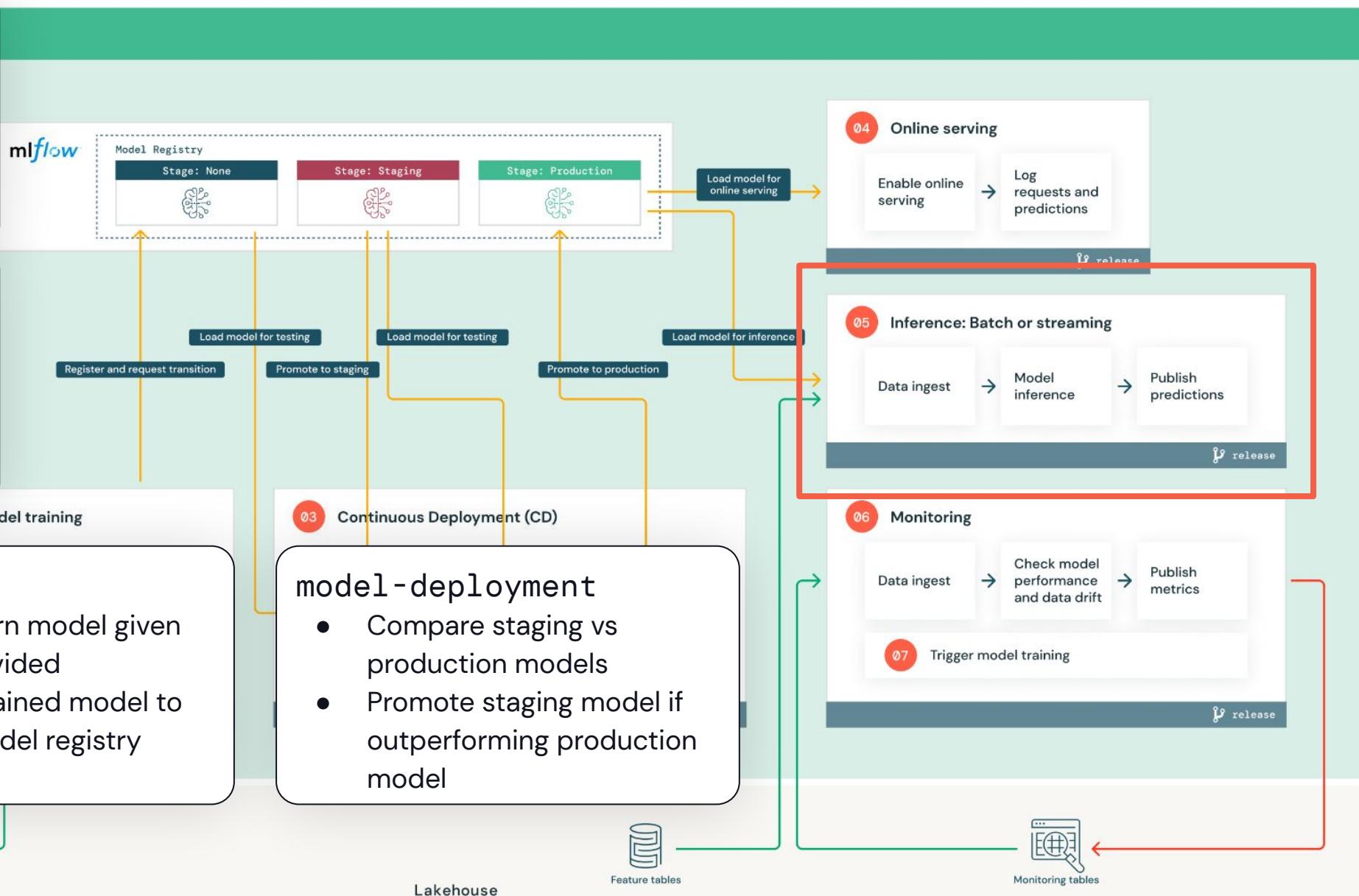


demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table

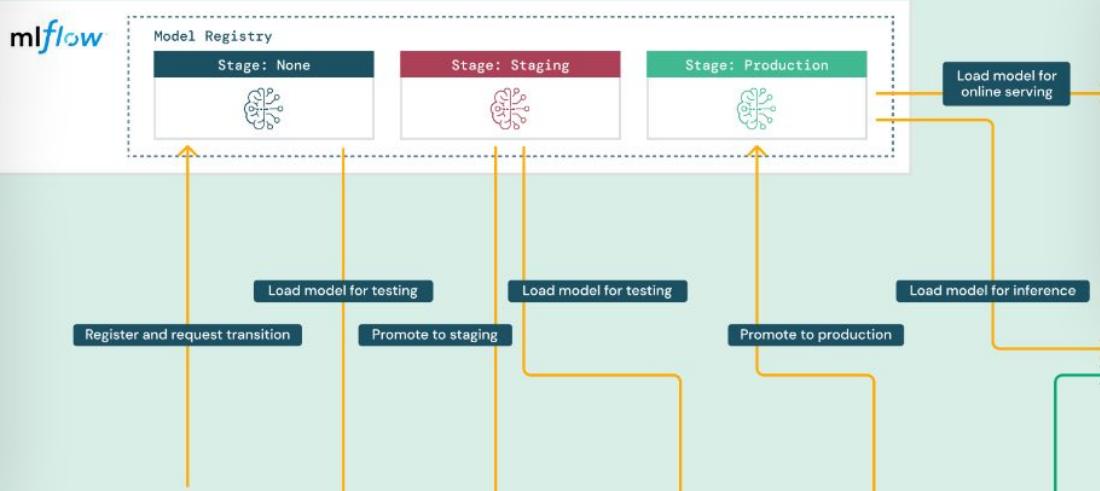


demo - setup

- Remove existing MLflow experiments/models
- Remove existing Feature Store tables

feature-table-creation

- Create database and tables if they don't exist
- Create single Feature Store table
- Create labels table



model-train

- Train sklearn model given config provided
- Register trained model to MLflow model registry



Data tables



Feature tables

model-deployment

- Compare staging vs production models
- Promote staging model if outperforming production model

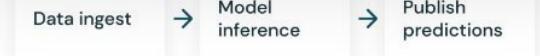
Lakehouse



Feature tables

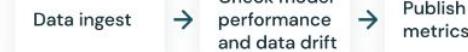
model-inference-batch

- Load features from Feature Store
- Load model from MLflow Model Registry
- Perform batch inference



release

06 Monitoring



07 Trigger model training

release



Monitoring tables

Project Structure

Project Structure

```
├── README.md  
├── .dbx           <- Directory containing environment configuration.  
└── project.json   <- Used to configure the project artifact locations and Databricks CLI profiles.
```

Project Structure

```
├── README.md  
├── .dbx           <- Directory containing environment configuration.  
├── project.json   <- Used to configure the project artifact locations and Databricks CLI profiles.  
├── .github  
└── workflows      <- Directory containing GitHub Actions Workflows for CI testing.
```

Project Structure

```
├── README.md
├── .dbx           <- Directory containing environment configuration.
├── project.json   <- Used to configure the project artifact locations and Databricks CLI profiles.
├── .github
│   ├── workflows    <- Directory containing GitHub Actions Workflows for CI testing.
├── requirements.txt <- Requirements file for reproducing the package.
├── unit-requirements.txt <- Requirements file for unit tests.
└── setup.py        <- Used to build the package .whl
```

Project Structure

```
├── README.md
├── .dbx                                <- Directory containing environment configuration.
├── project.json                         <- Used to configure the project artifact locations and Databricks CLI profiles.
├── .github
│   ├── workflows                         <- Directory containing GitHub Actions Workflows for CI testing.
│   └── requirements.txt                  <- Requirements file for reproducing the package.
├── requirements.txt                      <- Requirements file for unit tests.
├── unit-requirements.txt                <- Requirements file for unit tests.
└── setup.py                            <- Used to build the package .whl

|
├── conf
│   ├── job_configs                       <- Directory containing .yml configuration files for pipeline definitions.
│   └── deployment.yml                   <- File to configure job deployment properties (cluster configs, job specs).
```

Project Structure

```
├── README.md
├── .dbx                                     <- Directory containing environment configuration.
├── project.json                            <- Used to configure the project artifact locations and Databricks CLI profiles.
├── .github
│   ├── workflows                            <- Directory containing GitHub Actions Workflows for CI testing.
│   └── requirements.txt                     <- Requirements file for reproducing the package.
├── requirements.txt                         <- Requirements file for unit tests.
├── unit-requirements.txt                   <- Requirements file for unit tests.
└── setup.py                                <- Used to build the package .whl

|
├── conf                                    <- Directory containing configuration files and environment variables files.
│   ├── job_configs                         <- Directory containing .yml configuration files for pipeline definitions.
│   └── deployment.yml                     <- File to configure job deployment properties (cluster configs, job specs).
|
└── telco_churn                            <- Python source code for the project pipelines.
```

Project Structure

```
├── README.md
├── .dbx                                     <- Directory containing environment configuration.
├── project.json                             <- Used to configure the project artifact locations and Databricks CLI profiles.
├── .github
│   ├── workflows                            <- Directory containing GitHub Actions Workflows for CI testing.
│   └── requirements.txt                     <- Requirements file for reproducing the package.
├── requirements.txt                         <- Requirements file for unit tests.
├── unit-requirements.txt                   <- Requirements file for unit tests.
└── setup.py                                 <- Used to build the package .whl

|
├── conf                                    <- Directory containing configuration files and environment variables files.
│   ├── job_configs                         <- Directory containing .yml configuration files for pipeline definitions.
│   └── deployment.yml                     <- File to configure job deployment properties (cluster configs, job specs).
|
└── telco_churn                            <- Python source code for the project pipelines.

|
├── tests                                   <- Unit and integration test code for the various pipelines.
│   ├── unit                                <- Directory containing unit tests.
│   └── integration                         <- Directory containing integration tests.
```

Project Structure

Who owns what?

| | |
|-----------------------|--|
| README.md | |
| .dbx | <- Directory containing environment configuration. |
| project.json | <- Used to configure the project artifact locations and Databricks CLI profiles. |
| .github | |
| workflows | <- Directory containing GitHub Actions Workflows for CI testing. |
| requirements.txt | <- Requirements file for reproducing the package. |
| unit-requirements.txt | <- Requirements file for unit tests. |
| setup.py | <- Used to build the package .whl |
| conf | |
| job_configs | <- Directory containing configuration files and environment variables files. |
| deployment.yml | <- File to configure job deployment properties (cluster configs, job specs). |
| telco_churn | <- Python source code for the project pipelines. |
| tests | |
| unit | <- Unit and integration test code for the various pipelines. |
| integration | <- Directory containing unit tests. <- Directory containing integration tests. |

Project Structure

Who owns what?

| | | |
|-----------------------|--|-------------|
| README.md | | ML Engineer |
| .dbx | <- Directory containing environment configuration. | |
| project.json | <- Used to configure the project artifact locations and Databricks CLI profiles. | |
| .github | | |
| workflows | <- Directory containing GitHub Actions Workflows for CI testing. | |
| requirements.txt | <- Requirements file for reproducing the package. | |
| unit-requirements.txt | <- Requirements file for unit tests. | |
| setup.py | <- Used to build the package .whl | |
| conf | | |
| job_configs | <- Directory containing configuration files and environment variables files. | |
| deployment.yml | <- Directory containing .yml configuration files for pipeline definitions. <- File to configure job deployment properties (cluster configs, job specs). | |
| telco_churn | <- Python source code for the project pipelines. | |
| tests | | |
| unit | <- Unit and integration test code for the various pipelines. | |
| integration | <- Directory containing unit tests. <- Directory containing integration tests. | |

Project Structure

Who owns what?

| | | |
|-----------------------|---|----------------|
| README.md | | ML Engineer |
| .dbx | <- Directory containing environment configuration. | |
| project.json | <- Used to configure the project artifact locations and Databricks CLI profiles. | |
| .github | | |
| workflows | <- Directory containing GitHub Actions Workflows for CI testing. | |
| requirements.txt | <- Requirements file for reproducing the package. | |
| unit-requirements.txt | <- Requirements file for unit tests. | |
| setup.py | <- Used to build the package .whl | |
| conf | | |
| job_configs | <- Directory containing .yml configuration files for pipeline definitions. | |
| deployment.yml | <- File to configure job deployment properties (cluster configs, job specs). | |
| telco_churn | <- Python source code for the project pipelines. | |
| tests | | |
| unit | <- Unit and integration test code for the various pipelines. | |
| integration | <- Directory containing unit tests. <- Directory containing integration tests. | Data Scientist |

Project Structure

```
├── README.md  
├── .dbx                                     <- Directory containing environment configuration.  
└── project.json                            <- Used to configure the project artifact locations and Databricks CLI profiles.  
├── .github  
│   └── workflows                            <- Directory containing GitHub Actions Workflows for CI testing.  
├── requirements.txt                         <- Requirements file for reproducing the package.  
├── unit-requirements.txt                    <- Requirements file for unit tests.  
└── setup.py                                 <- Used to build the package .whl  
  
├── conf  
│   ├── job_configs                          <- Directory containing .yml configuration files for pipeline definitions.  
│   └── deployment.yml                      <- File to configure job deployment properties (cluster configs, job specs).  
└── telco_churn                             <- Python source code for the project pipelines.  
    ├── tests  
    │   ├── unit                                <- Unit and integration test code for the various pipelines.  
    │   └── integration                         <- Directory containing unit tests.  
    └── integration                           <- Directory containing integration tests.
```

Data Scientist

Project Structure

```
|--- telco_churn
```

Project Structure

```
|── telco_churn  
|   ├── __init__.py  
|   └── common.py
```

<- Python file containing a generic class called Job to run jobs via dbx.

Project Structure

```
telco_churn
├── __init__.py
└── common.py      ← Python file containing a generic class called Job to run jobs via dbx.

├── featurize.py   ← Perform data preprocessing steps on features and label.
└── feature_table_creator.py  ← Create Databricks Feature Store table.
```

Project Structure

```
telco_churn
├── __init__.py
└── common.py
    ← Python file containing a generic class called Job to run jobs via dbx.

feature-table-creation
├── featurize.py
└── feature_table_creator.py
    ← Perform data preprocessing steps on features and label.
    ← Create Databricks Feature Store table.
```

Project Structure

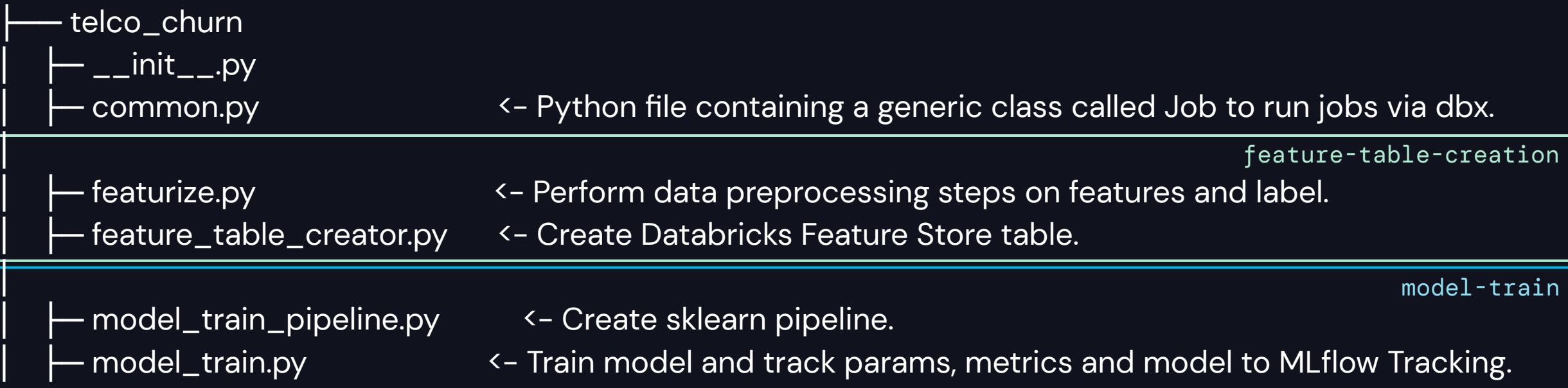
```
telco_churn
├── __init__.py
└── common.py
    ← Python file containing a generic class called Job to run jobs via dbx.

featurize.py
feature_table_creator.py
    ← Perform data preprocessing steps on features and label.
    ← Create Databricks Feature Store table.

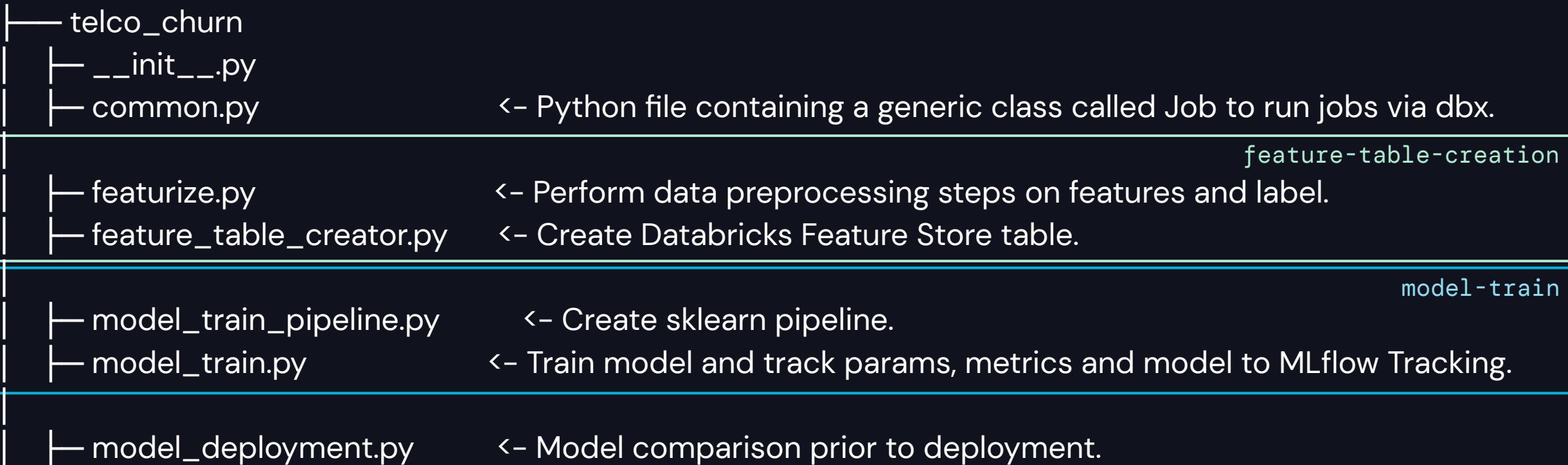
feature-table-creation

model_train_pipeline.py
model_train.py
    ← Create sklearn pipeline.
    ← Train model and track params, metrics and model to MLflow Tracking.
```

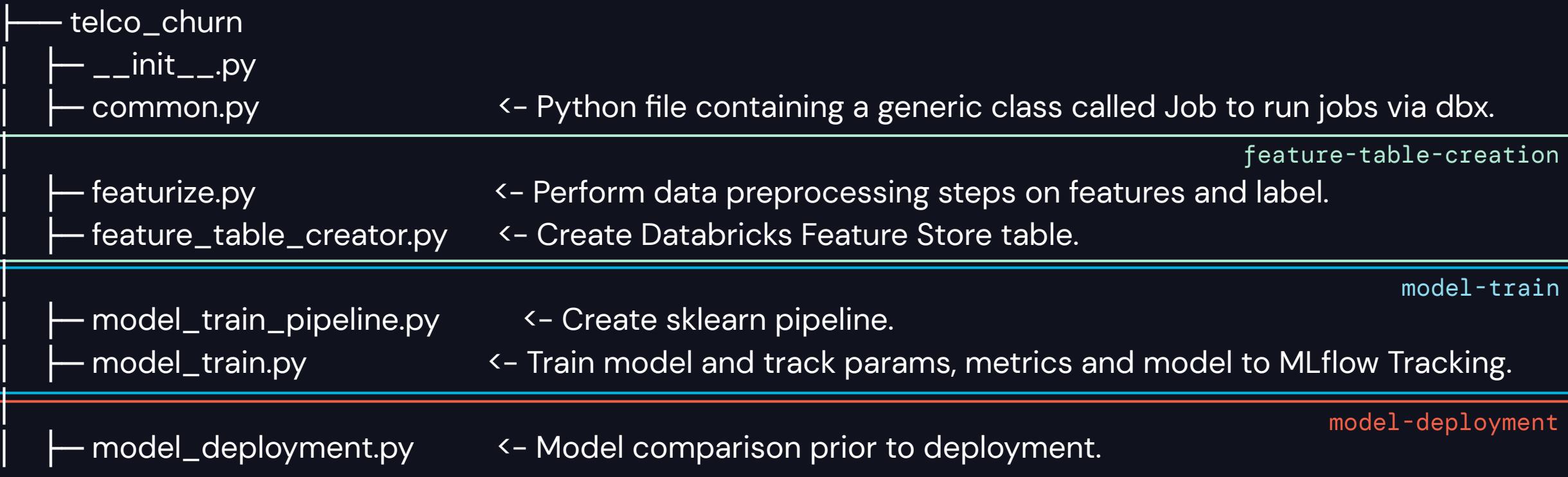
Project Structure



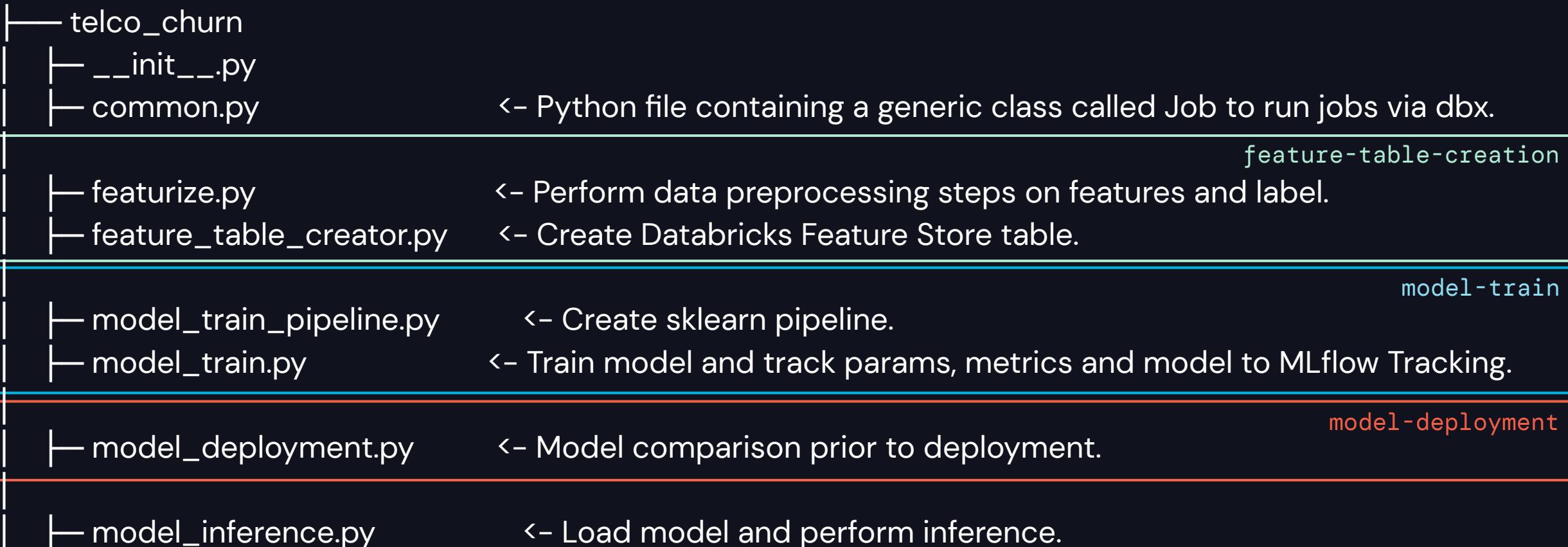
Project Structure



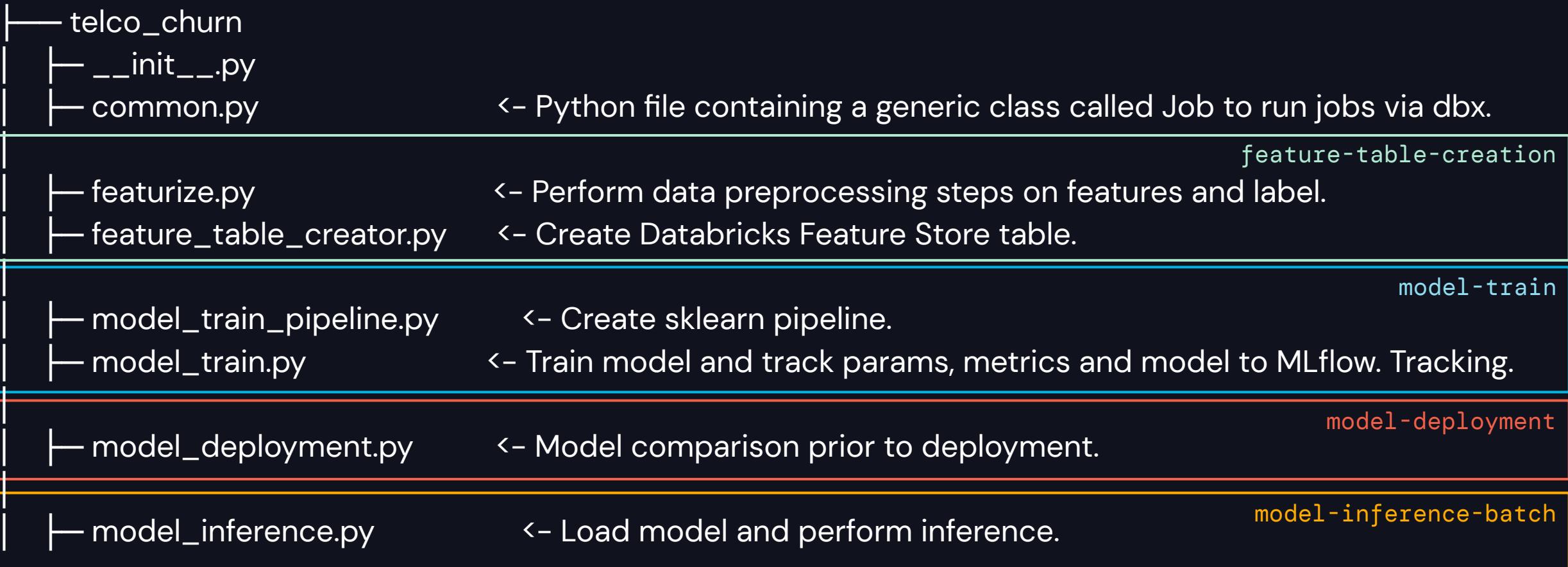
Project Structure



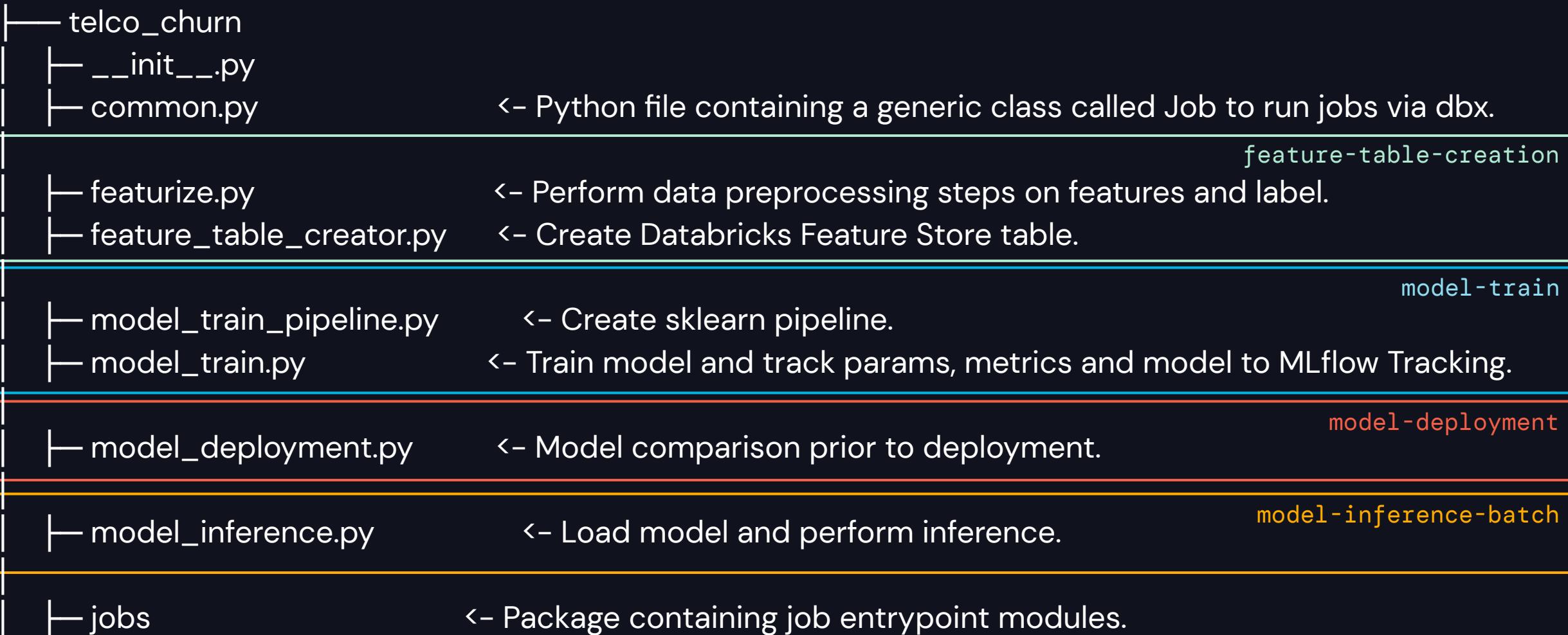
Project Structure



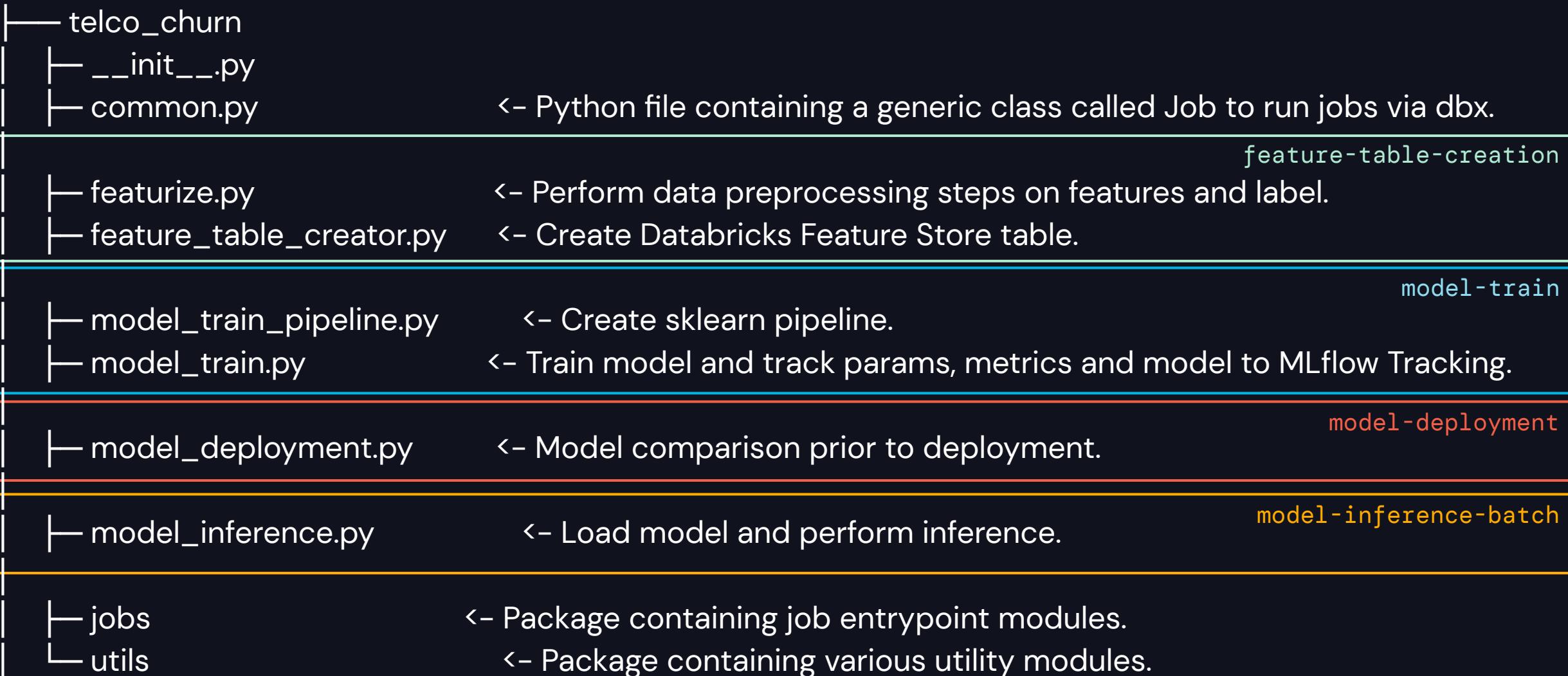
Project Structure



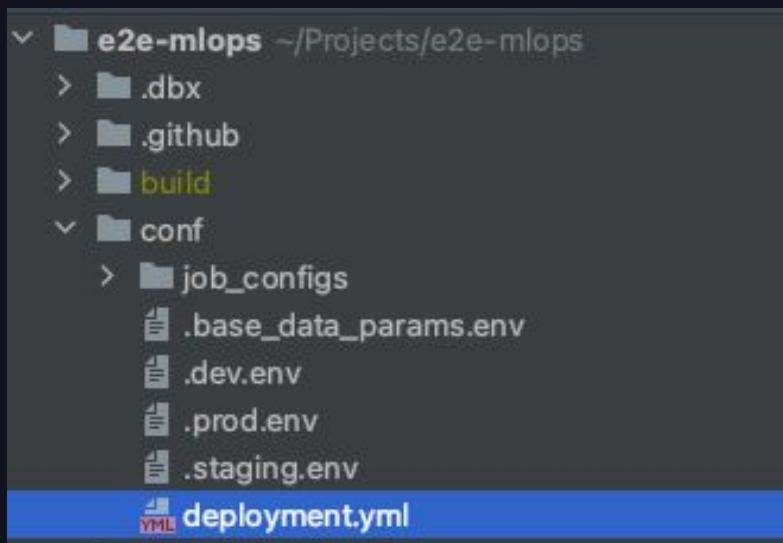
Project Structure



Project Structure

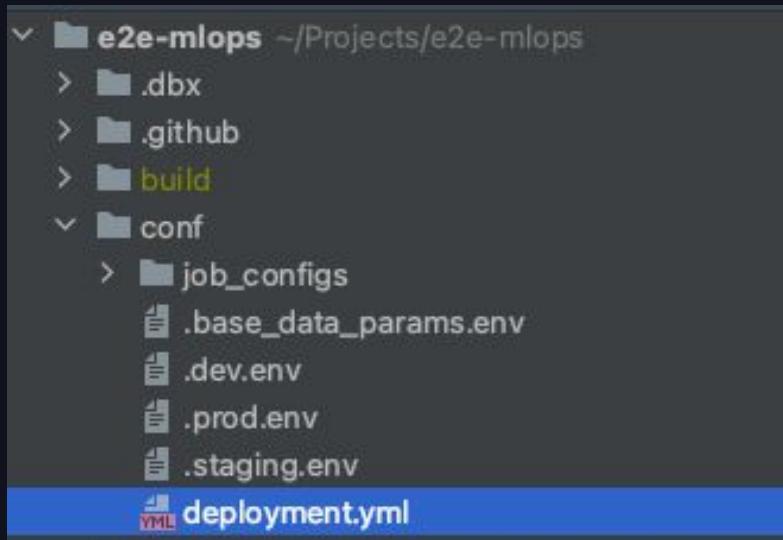


Mechanics of pipeline



deployment.yml

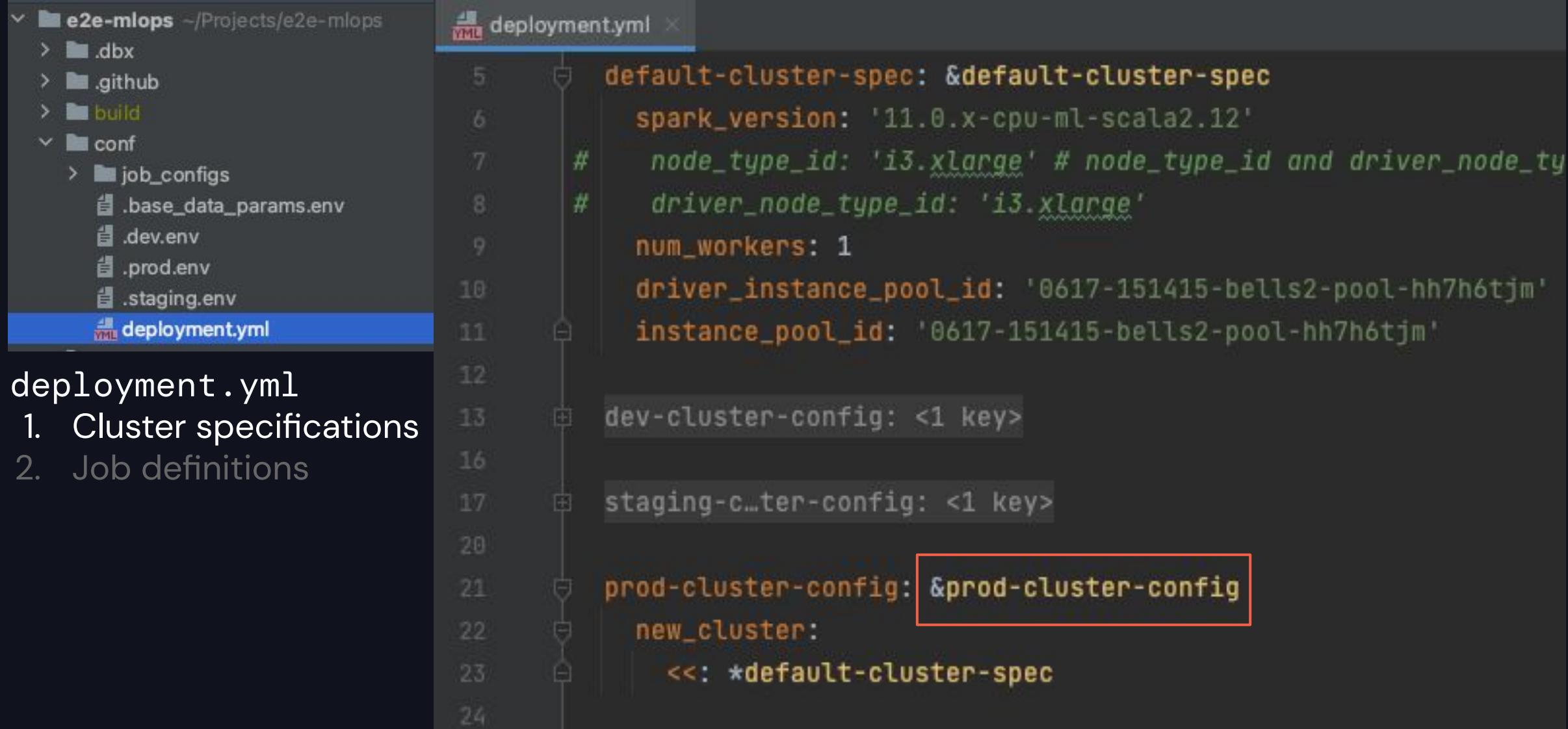
Mechanics of pipeline



`deployment.yml`

1. Cluster specifications
2. Job definitions

Mechanics of pipeline



The image shows a file explorer on the left and a code editor on the right. The file explorer lists a directory structure for 'e2e-mlops' containing '.dbx', '.github', 'build', 'conf' (which contains 'job_configs' and environment files '.base_data_params.env', '.dev.env', '.prod.env', '.staging.env'), and 'deployment.yml'. The 'deployment.yml' file is selected in the file explorer and is also the active tab in the code editor.

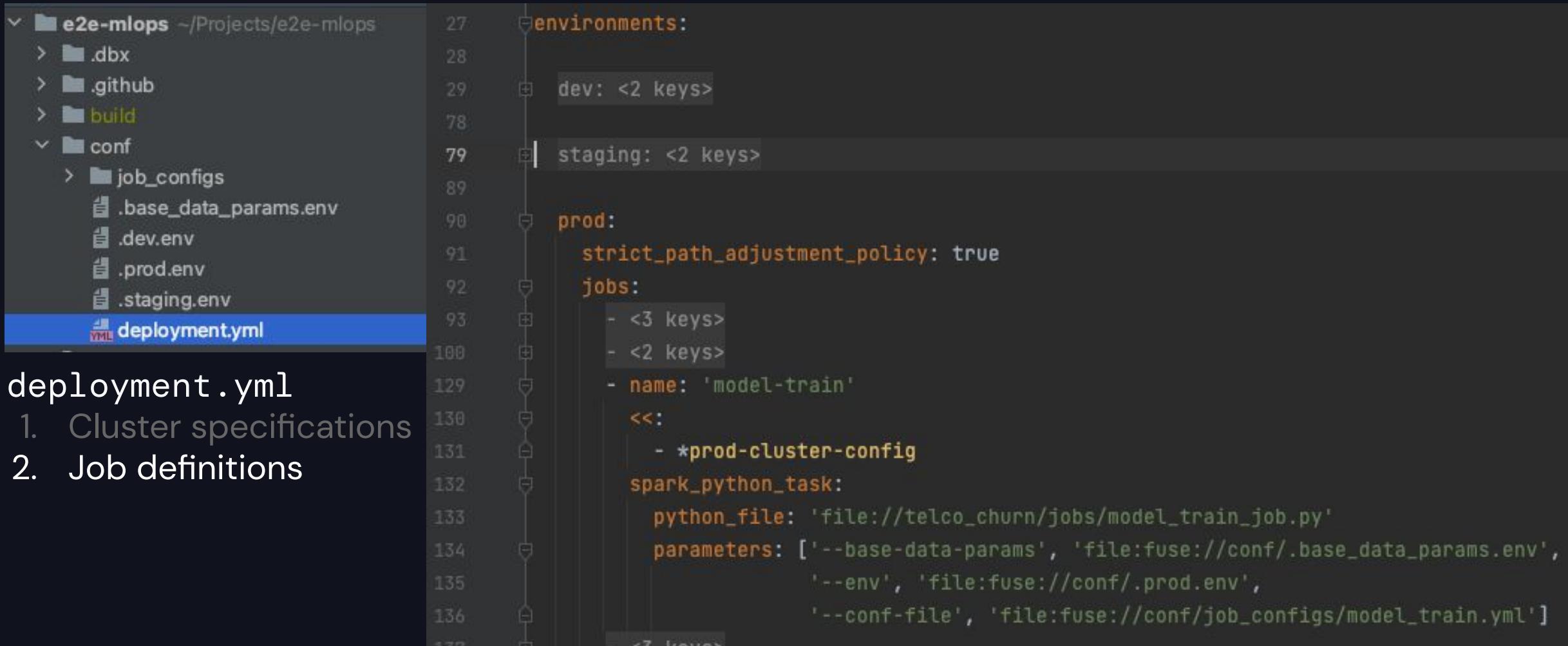
```
YML deployment.yml
5     default-cluster-spec: &default-cluster-spec
6         spark_version: '11.0.x-cpu-ml-scala2.12'
7         #     node_type_id: 'i3.xlarge' # node_type_id and driver_node_ty
8         #     driver_node_type_id: 'i3.xLarge'
9         num_workers: 1
10        driver_instance_pool_id: '0617-151415-bells2-pool-hh7h6tjm'
11        instance_pool_id: '0617-151415-bells2-pool-hh7h6tjm'
12
13        dev-cluster-config: <1 key>
14
15        staging-cluster-config: <1 key>
16
17        prod-cluster-config: &prod-cluster-config
18        new_cluster:
19            <<: *default-cluster-spec
20
21
22
23
24
```

The code editor displays the 'deployment.yml' configuration file. It includes cluster specifications for 'dev', 'staging', and 'prod' environments, along with a 'new_cluster' section that inherits from the 'default-cluster-spec'.

deployment.yml

1. Cluster specifications
2. Job definitions

Mechanics of pipeline



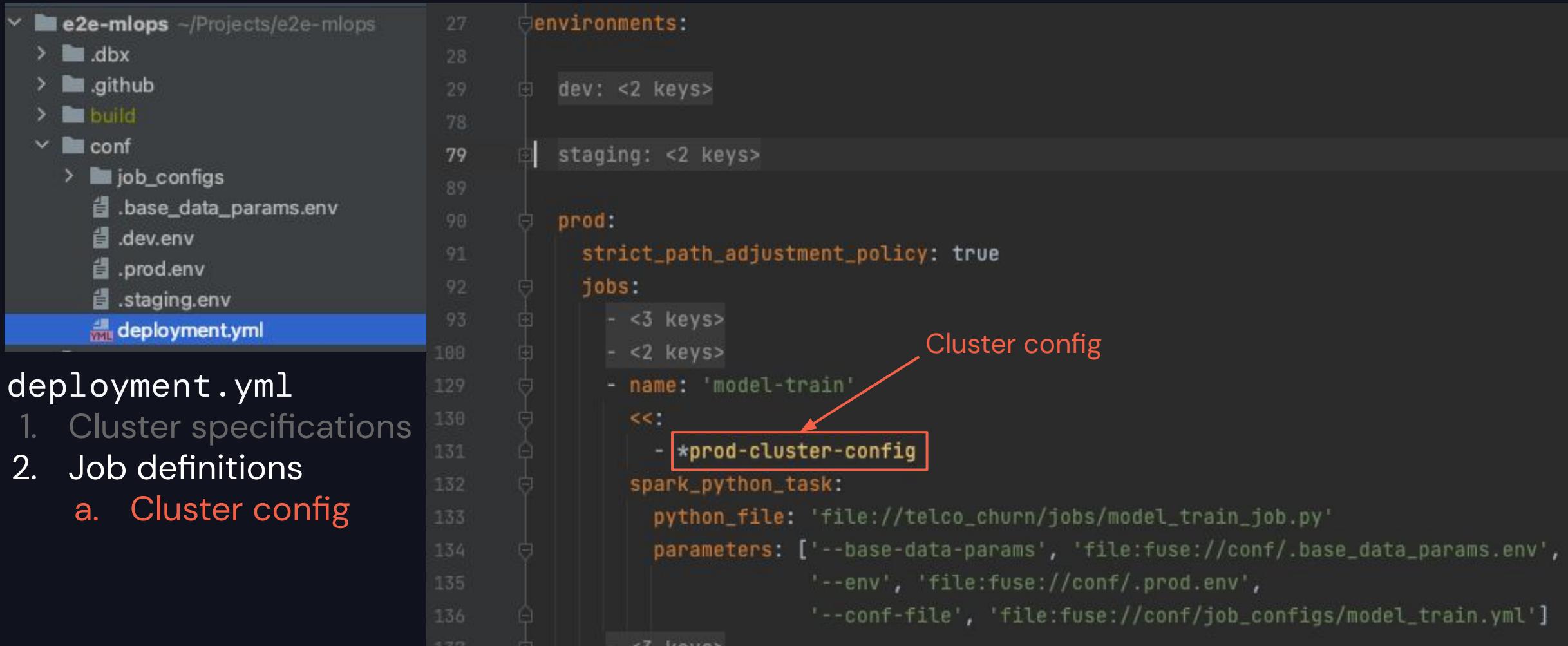
The screenshot shows a file tree on the left and the content of the `deployment.yml` file on the right. The file is part of a project named `e2e-mlops` located at `~/Projects/e2e-mlops`. The `deployment.yml` file is highlighted in blue.

```
environments:
  dev: <2 keys>
  staging: <2 keys>
prod:
  strict_path_adjustment_policy: true
  jobs:
    - <3 keys>
    - <2 keys>
    - name: 'model-train'
      <<:
        - *prod-cluster-config
      spark_python_task:
        python_file: 'file://telco_churn/jobs/model_train_job.py'
        parameters: ['--base-data-params', 'file:fuse://conf/.base_data_params.env',
                     '--env', 'file:fuse://conf/.prod.env',
                     '--conf-file', 'file:fuse://conf/job_configs/model_train.yml']
```

deployment.yml

1. Cluster specifications
2. Job definitions

Mechanics of pipeline



The image shows a file tree on the left and the content of `deployment.yml` on the right. The file tree includes `e2e-mlops`, `.dbx`, `.github`, `build`, `conf` (containing `job_configs`, `.base_data_params.env`, `.dev.env`, `.prod.env`, and `.staging.env`), and the selected `deployment.yml`.

```
environments:
  dev: <2 keys>
  staging: <2 keys>
prod:
  strict_path_adjustment_policy: true
  jobs:
    - <3 keys>
    - <2 keys>
    - name: 'model-train'
      <<:
        - *prod-cluster-config
spark_python_task:
  python_file: 'file://telco_churn/jobs/model_train_job.py'
  parameters: ['--base-data-params', 'file:fuse://conf/.base_data_params.env',
               '--env', 'file:fuse://conf/.prod.env',
               '--conf-file', 'file:fuse://conf/job_configs/model_train.yml']
```

A red arrow points from the text "Cluster config" to the line `- *prod-cluster-config`, which is highlighted with a red box.

deployment .yml

1. Cluster specifications
2. Job definitions
 - a. Cluster config

Mechanics of pipeline

The screenshot shows a file browser interface with the following project structure:

- e2e-mlops (~/Projects/e2e-mlops)
 - .dbx
 - .github
 - build
 - conf
 - job_configs
 - .base_data_params.env
 - .dev.env
 - .prod.env
 - .staging.env
- deployment.yml

The `deployment.yml` file content is as follows:

```
environments:  
  dev: <2 keys>  
  staging: <2 keys>  
prod:  
  strict_path_adjustment_policy: true  
  jobs:  
    - <3 keys>  
    - <2 keys>  
    - name: 'model-train'  
      <<:  
        - *prod-cluster-config  
      spark_python_task:  
        python_file: 'file:///telco_churn/jobs/model_train_job.py'  
        parameters: ['--base-data-params', 'file:/fuse://conf/.base_data_params.env',  
                     '--env', 'file:/fuse://conf/.prod.env',  
                     '--conf-file', 'file:/fuse://conf/job_configs/model_train.yml']
```

Annotations:

- A red box highlights the `*prod-cluster-config` entry under the `model-train` job.
- A blue box highlights the `spark_python_task` section.
- A red arrow points from the text "Cluster config" to the `*prod-cluster-config` entry.
- A blue arrow points from the text "Job entry file" to the `spark_python_task` section.

deployment .yml

1. Cluster specifications
2. Job definitions
 - a. Cluster config
 - b. Job entry file

Mechanics of pipeline

The image shows a file tree on the left and the content of `deployment.yml` on the right. The file tree includes `e2e-mlops`, `.dbx`, `.github`, `build`, `conf` (containing `job_configs` with files `.base_data_params.env`, `.dev.env`, `.prod.env`, `.staging.env`), and the `deployment.yml` file itself.

```
environments:
  dev: <2 keys>
  staging: <2 keys>
  prod:
    strict_path_adjustment_policy: true
    jobs:
      - <3 keys>
      - <2 keys>
      - name: 'model-train'
        <<:
          - *prod-cluster-config
spark_python_task:
  python_file: 'file:///telco_churn/jobs/model_train_job.py'
  parameters: ['--base-data-params', 'file:/fuse://conf/.base_data_params.env',
               '--env', 'file:/fuse://conf/.prod.env',
               '--conf-file', 'file:/fuse://conf/job_configs/model_train.yml']
```

Annotations on the right side of the code:

- Cluster config**: Points to the `*prod-cluster-config` placeholder in the `spark_python_task` section.
- Job entry file**: Points to the `model-train` job definition.
- Environment specific variables**: Points to the `parameters` section of the `spark_python_task`.

deployment .yml

1. Cluster specifications
2. Job definitions
 - a. Cluster config
 - b. Job entry file
 - c. Environment vars

Mechanics of pipeline

The image shows a file browser interface with a sidebar on the left displaying the project structure:

- e2e-mlops ~/Projects/e2e-mlops
 - .dbx
 - .github
 - build
 - conf
 - job_configs
 - .base_data_params.env
 - .dev.env
 - .prod.env
 - .staging.env
- deployment.yml

The deployment.yml file is open in the main pane, showing its YAML content:

```
environments:  
  dev: <2 keys>  
  staging: <2 keys>  
prod:  
  strict_path_adjustment_policy: true  
  jobs:  
    - <3 keys>  
    - <2 keys>  
    - name: 'model-train'  
      <<:  
        - *prod-cluster-config  
      spark_python_task:  
        python_file: 'file:///telco_churn/jobs/model_train_job.py'  
        parameters: ['--base-data-params', 'file:/fuse://conf/.base_data_params.env',  
                     '--env', 'file:/fuse://conf/.prod.env',  
                     '--conf-file', 'file:/fuse://conf/job_configs/model_train.yml']
```

Annotations on the right side of the code highlight specific parts:

- Cluster config**: Points to the `prod` section.
- Job entry file**: Points to the `model-train` job definition.
- Environment specific variables**: Points to the `parameters` section.
- Config file**: Points to the `.base_data_params.env` file in the sidebar.

Example: model-train

Job entry file: telco_churn/jobs/model_train_job.py

The screenshot shows a code editor with a file tree on the left and the content of `model_train_job.py` on the right. The file tree shows the project structure: `e2e-mlops` contains `.dbx`, `.github`, `build`, `conf`, `dist`, `telco_churn` (which contains `jobs` with files: `__init__.py`, `demo_setup_job.py`, `feature_table_creator_job.py`, `model_deployment_job.py`, `model_inference_job.py`, `model_train_job.py` (selected), and `sample_job.py`). The code editor window shows the following Python code:

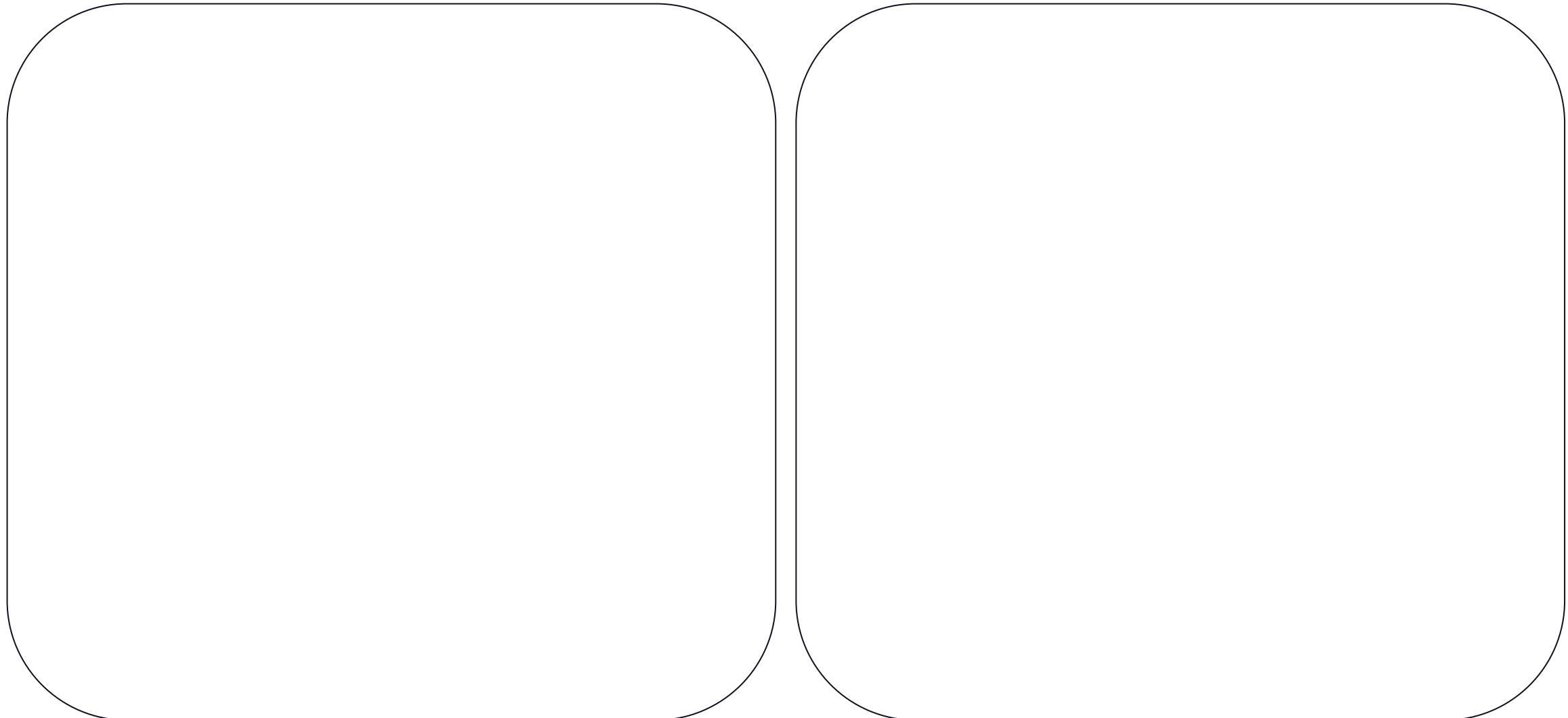
```
class ModelTrainJob(Workload):
    def _get_mlflow_tracking_cfg(self):
        try:
            experiment_id = self.env_vars['model_'
        except KeyError:
            experiment_id = None
        try:
            experiment_path = self.env_vars['mode'

```

Each job class does the following:

- Loads the environment specific variables
- Loads the config file
- Launches a Databricks job when executed

Demo Workflow



Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard

Data Scientist

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard

2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard
2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes
3. Create pull request
 - Unit tests
 - Integration tests
 - If tests pass, change is merged into main (staging) branch

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard
2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes
3. Create pull request
 - Unit tests
 - Integration tests
 - If tests pass, change is merged into main (staging) branch

ML Engineer

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard

2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes

3. Create pull request
 - Unit tests
 - Integration tests
 - If tests pass, change is merged into main (staging) branch

4. Cut release branch ML Engineer
 - Tag a new release version
 - On pushing tag, release pipeline triggered

Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard
2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes
3. Create pull request
 - Unit tests
 - Integration tests
 - If tests pass, change is merged into main (staging) branch

4. Cut release branch ML Engineer
 - Tag a new release version
 - On pushing tag, release pipeline triggered
5. Release pipeline deploys workflows to prod
 - model-train
 - model-deployment
 - model-inference-batch

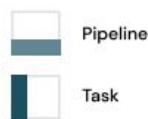
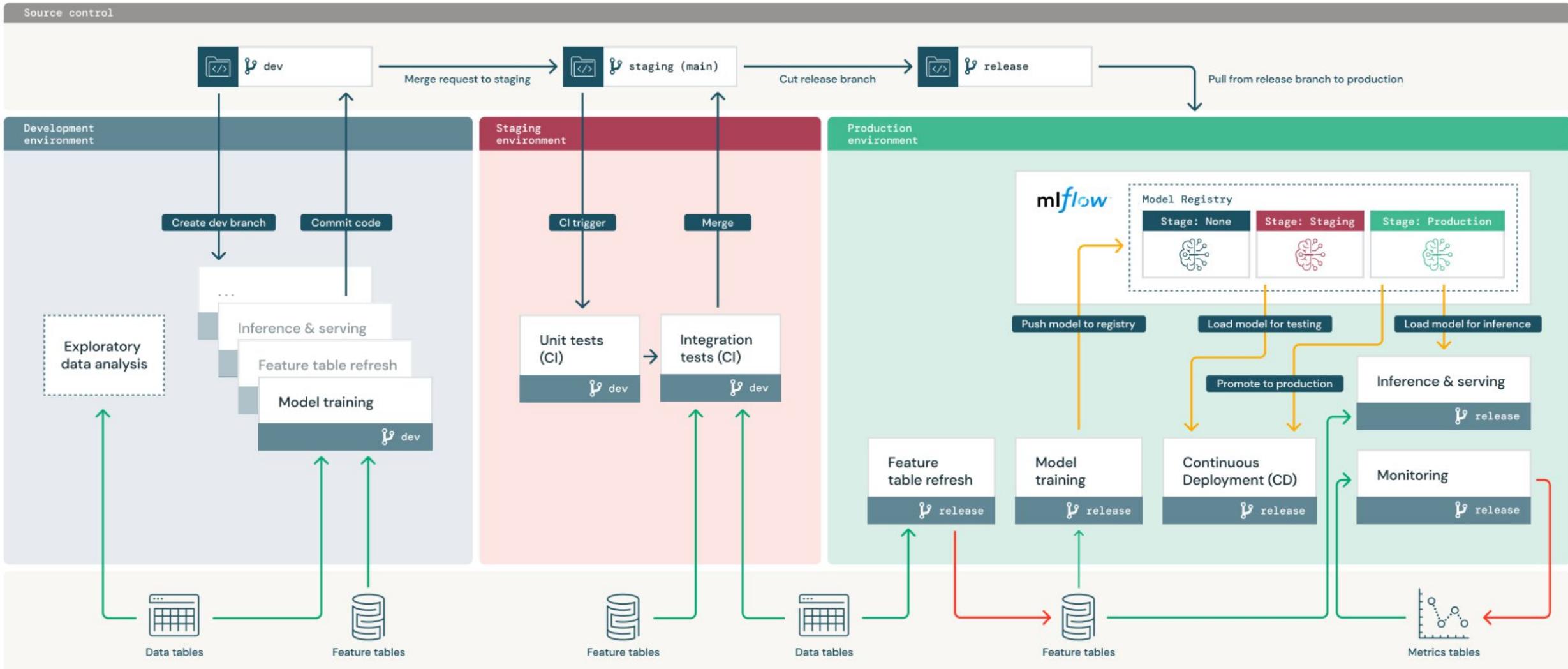
Demo Workflow

1. Start at steady state
 - Model in production
 - Feature tables
 - Existing pipelines in prod environment
 - Monitoring dashboard
2. Model update Data Scientist
 - Model performance has degraded
 - Update model train config file
 - Push code changes
3. Create pull request
 - Unit tests
 - Integration tests
 - If tests pass, change is merged into main (staging) branch

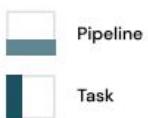
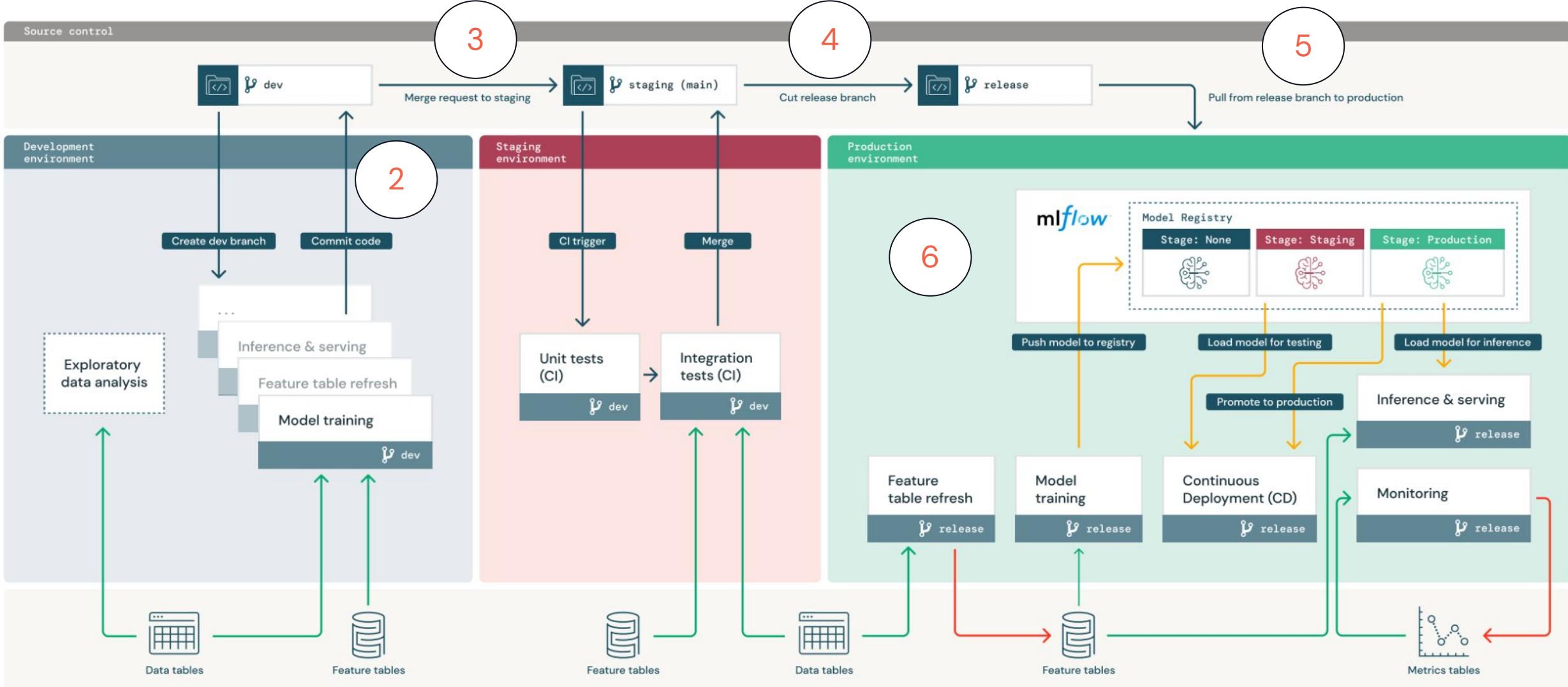
4. Cut release branch ML Engineer
 - Tag a new release version
 - On pushing tag, release pipeline triggered
5. Release pipeline deploys workflows to prod
 - model-train
 - model-deployment
 - model-inference-batch
6. Manually trigger or automate workflows

Demo

Demo workflow

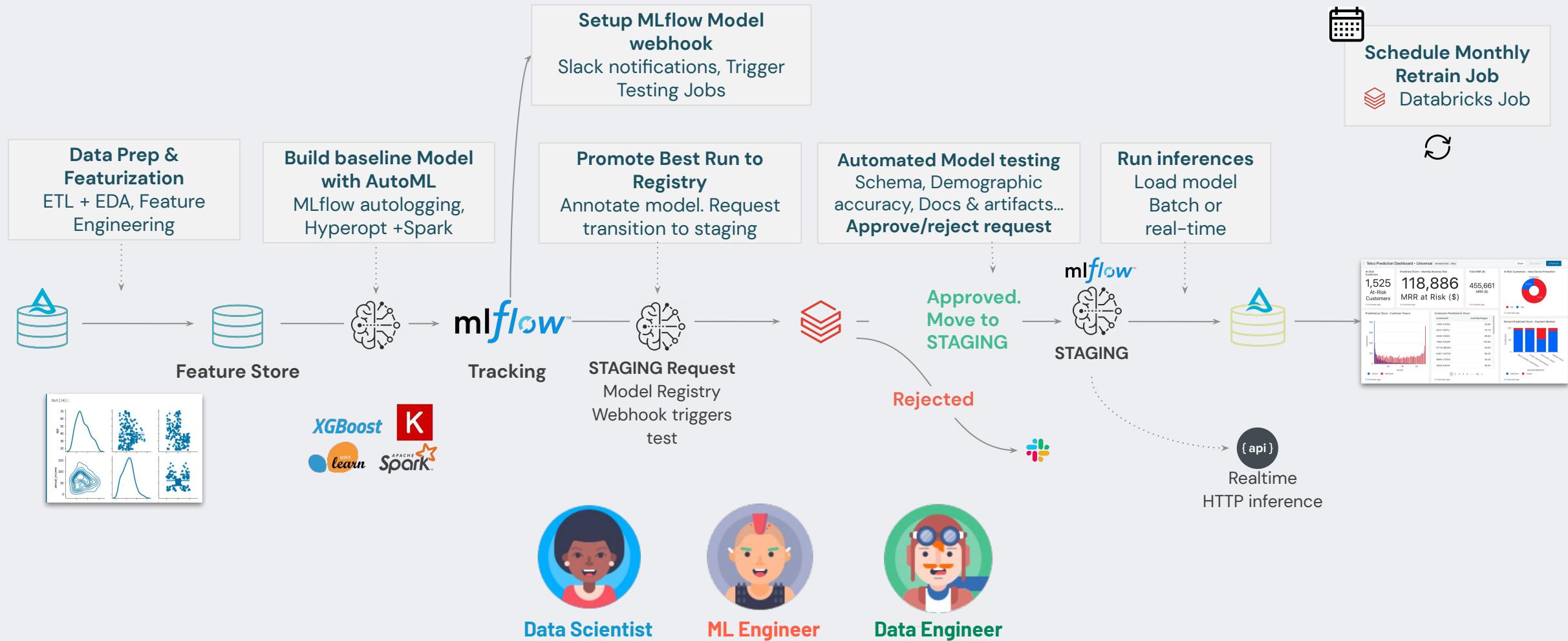


Demo workflow



Model Management

End to End Model Management



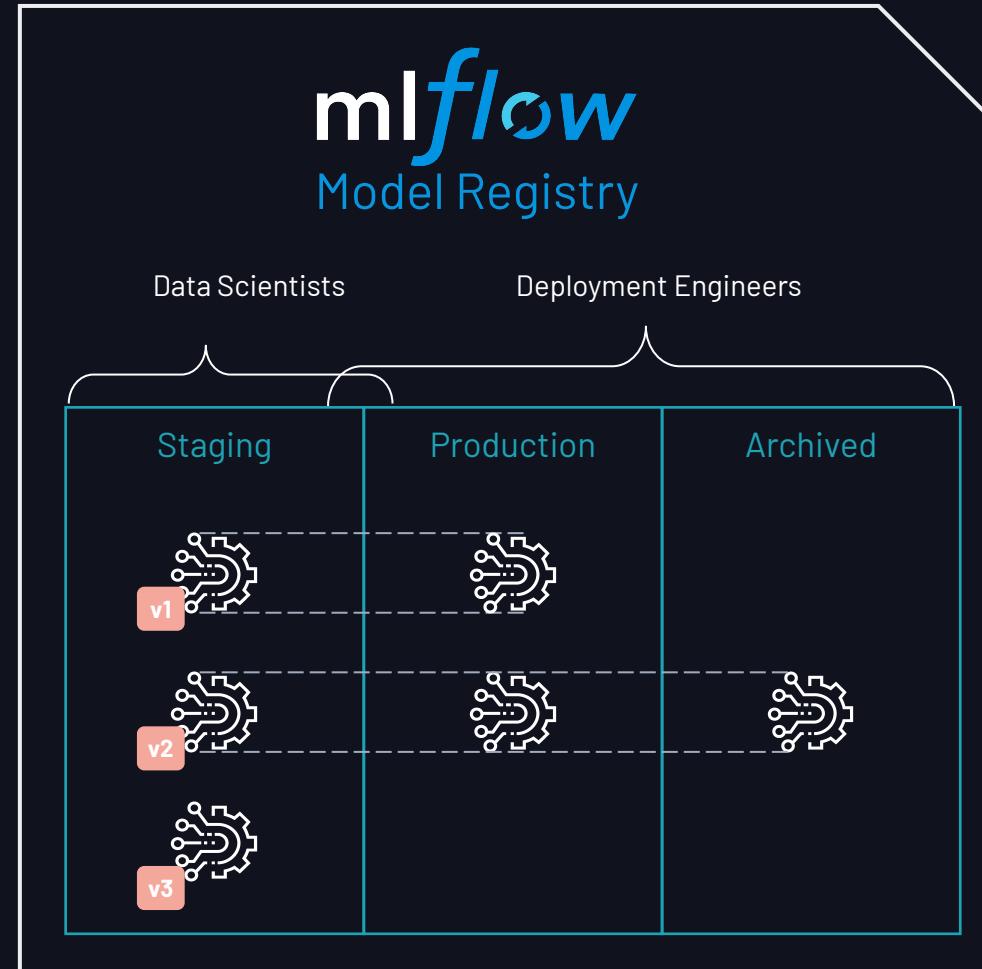
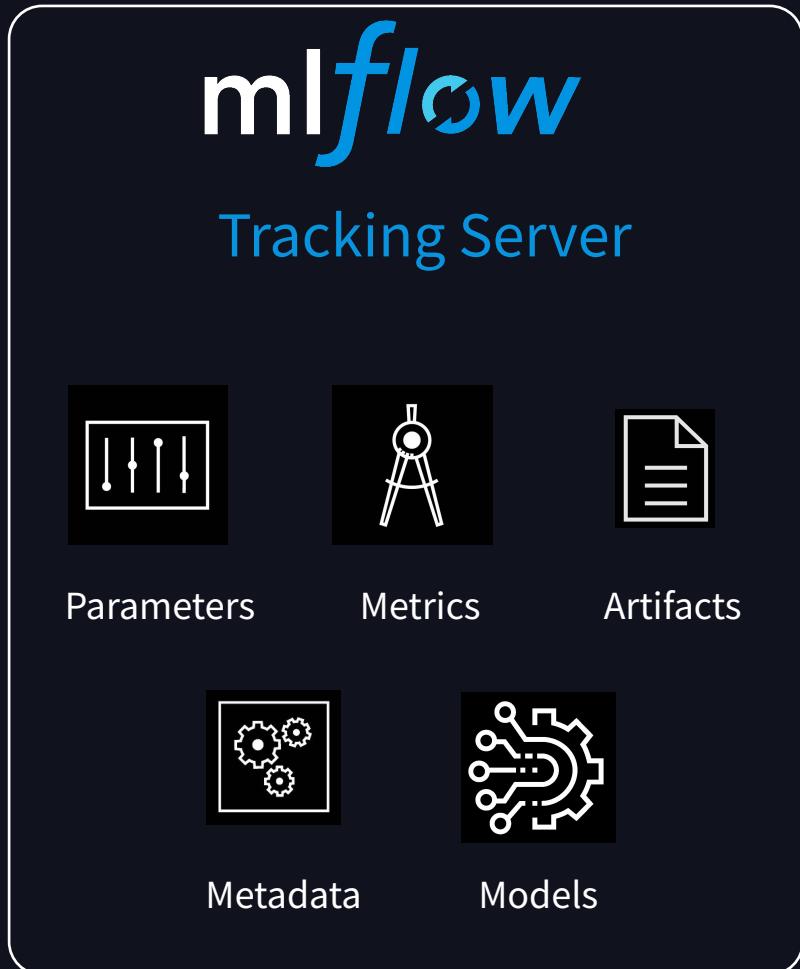
Data Transformation

Within the model or prior to model training? It depends!

| | Pros | Cons |
|--------------------------------|--|---|
| Prior to Model Training | <ul style="list-style-type: none">- Compute once and reuse- Leverages full dataset | <ul style="list-style-type: none">- Increases production footprint- Slower to iterate |
| Within the Model | <ul style="list-style-type: none">- Easier to iterate- Smaller production footprint | <ul style="list-style-type: none">- Model latency depends upon transformation overheads- Data visibility |

mlflow Model Registry

VISION: Centralized and collaborative model lifecycle management



Model Management

Use one central place to collaboratively manage ML models, from experimentation to online testing and production.

Provides:

- Automatic models versioning
- Integration with governance workflows
- Automatic/manual approval of stage transition
- Full visibility into model lifecycle from experimentation to production

| Registered Models | | | | | |
|--|----------------|---------|------------|---------------------|--|
| i Share and serve machine learning models. Learn more | | | | | |
| Create Model | | | | | |
| Name | Latest Version | Staging | Production | Last Modified | |
| airbnb | Version 1 | - | Version 1 | 2021-06-16 10:16:39 | |
| airbnb_hawaii | Version 1 | - | Version 1 | 2022-02-16 11:31:03 | |
| airbnb_model_monitoring_chengyin | Version 1 | - | Version 1 | 2022-02-26 14:33:09 | |
| airbnb_rf_model_39dbc2a3ca | Version 2 | - | - | 2021-07-14 19:34:17 | |
| airbnb_rf_model_3f55aa8cd6 | Version 2 | - | - | 2021-08-04 09:05:10 | |

Model Registry Benefits



One Collaborative Hub: The Model Registry provides a central hub for making models discoverable, improving **collaboration** and **knowledge sharing** across the organization.



Manage the entire Model Lifecycle (MLOps): The Model Registry provides lifecycle management for models from experimentation to deployment, improving **reliability** and robustness of the model deployment process.



Visibility and Governance: The Model Registry provides full visibility into the deployment stage of all models, who requested and approved changes, allowing for full governance and auditability.

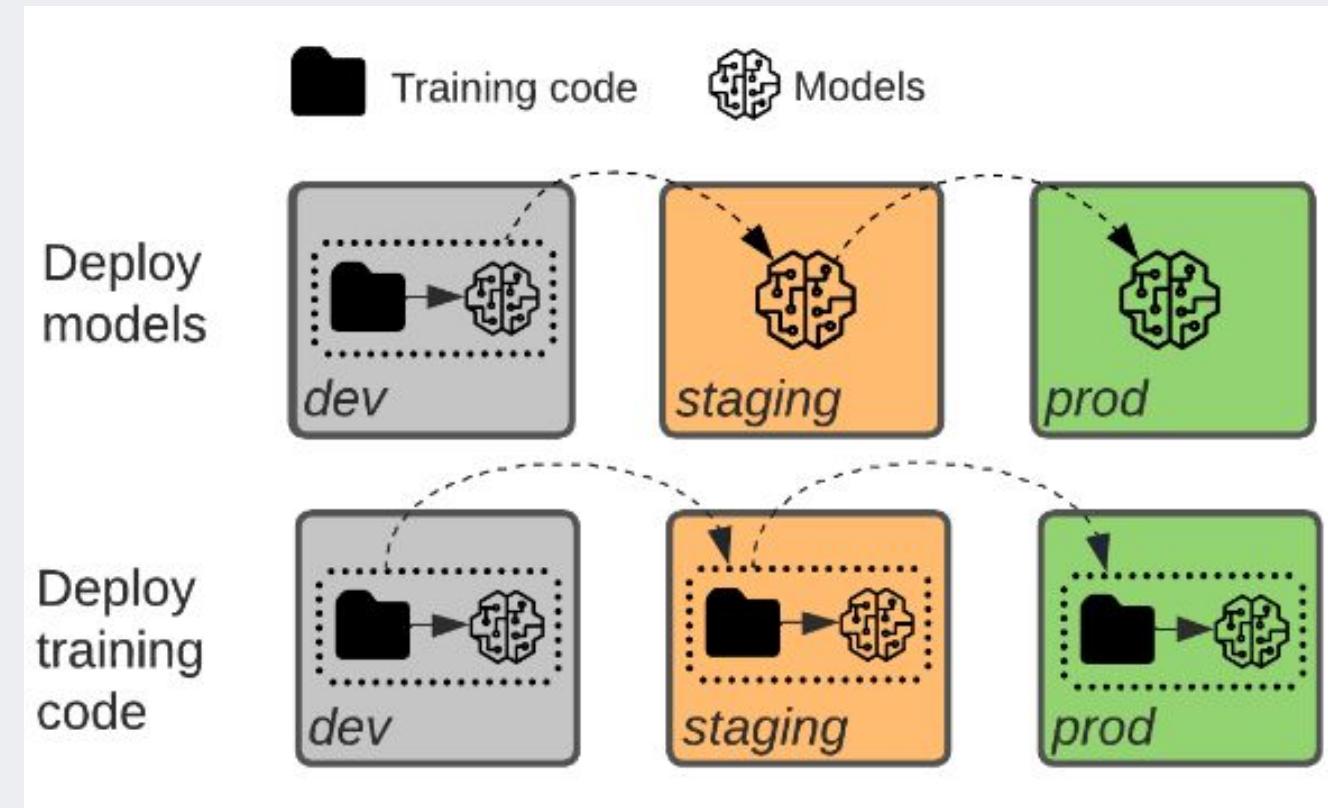
Model Deployment Patterns

Deploy models:

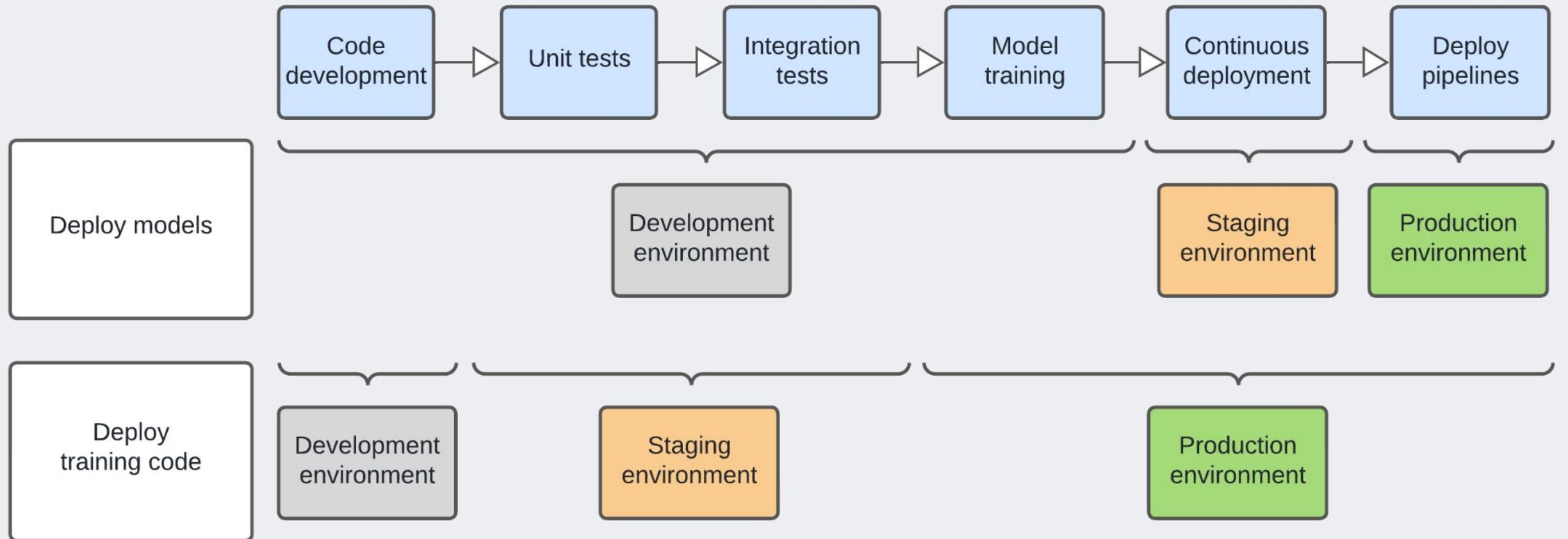
the model artifact is generated by training code in the development environment. This artifact is then tested in staging for compliance and performance, before finally being deployed into production.

Deploy training code:

the code to train models is developed in the dev environment, and this code is moved to staging and then production.



Contrast for Code Life Cycle



Comparisons (Process)

| | | Deploy models | Deploy training code |
|----------------|----------------|--|--|
| Process | Dev | Develop training code. Develop ancillary code. ¹ Train model on prod data. → Promote model & ancillary code. | Develop training code. Develop ancillary code. → Promote code. |
| | Staging | Test model & ancillary code. → Promote model & ancillary code. | Train model on data subset. Test ancillary code. → Promote code. |
| | Prod | Deploy model. Deploy ancillary pipelines. | Train model on prod data. Test model. Deploy model. Deploy ancillary pipelines. |

| | | Deploy models | Deploy training code |
|--------------------|------------------------------------|--|---|
| Trade-offs | Automation | ⬇️ Does not support automated retraining in locked-down env. | ⬆️ Supports automated retraining in locked-down env. |
| | Data access control | ⬇️ Dev env needs read access to prod training data. | ⬆️ Only prod env needs read access to prod training data. |
| | Reproducible models | ⬇️ Less Eng control over training env, so harder to ensure reproducibility. | ⬆️ Eng control over training env, which helps to simplify reproducibility. |
| | Data science familiarity | ⬆️ DS team builds and can directly test models in their dev env. | ⬇️ DS team must learn to write & hand off modular code to Eng. |
| | Support for large projects | ⬇️ This pattern does not force the DS team to use modular code for model training, and it has less iterative testing. | ⬆️ This pattern forces the DS team to use modular code and iterative testing, which helps with coordination and development in larger projects. |
| | Eng setup & maintenance | ⬆️ Has the simplest setup, with less CI/CD infra required. | ⬇️ Requires CI/CD infra for unit and integration tests, even for one-off models. |
| When to use | | Use this pattern when your model is a one-off or when model training is very expensive. Use when dev, staging and prod are not strictly separated envs. | Use this pattern by default. Use when dev, staging and prod are strictly separated envs. |

Getting Started

Resources

Full whitepaper: “The Big Book of MLOps”

bit.ly/big-book-of-mlops

End-to-end MLOps demo repository

bit.ly/e2e-mlops-databricks

Catch up on other DAIS talks on MLOps.

Databricks product-related talks include:

- [Enable Production ML with Databricks Feature Store](#) (Wed @ 10:45am)
- [MLflow Pipelines: Accelerating MLOps from Development to Production](#) (Wed @ 11am)
- [The Databricks Notebook: Front door of the Lakehouse](#) (Wed @ 5:30pm)

Questions?

DATA+AI
SUMMIT 2022

Thank you