



**An open platform for the machine
learning lifecycle**

Abdulrahman Alfozan

PyData Riyadh, August 2020

Intro

Systems Engineer @ FB

- Apache Spark, Distributed Systems
- Data Engineering
- Applied Machine Learning



- Open-source distributed cluster-computing framework.
- **Unified analytics engine** for big data and machine learning.

Agenda

1. ML Pipeline Lifecycle
2. ML development challenges
3. MLflow platform
4. Demo

ML Pipeline Lifecycle

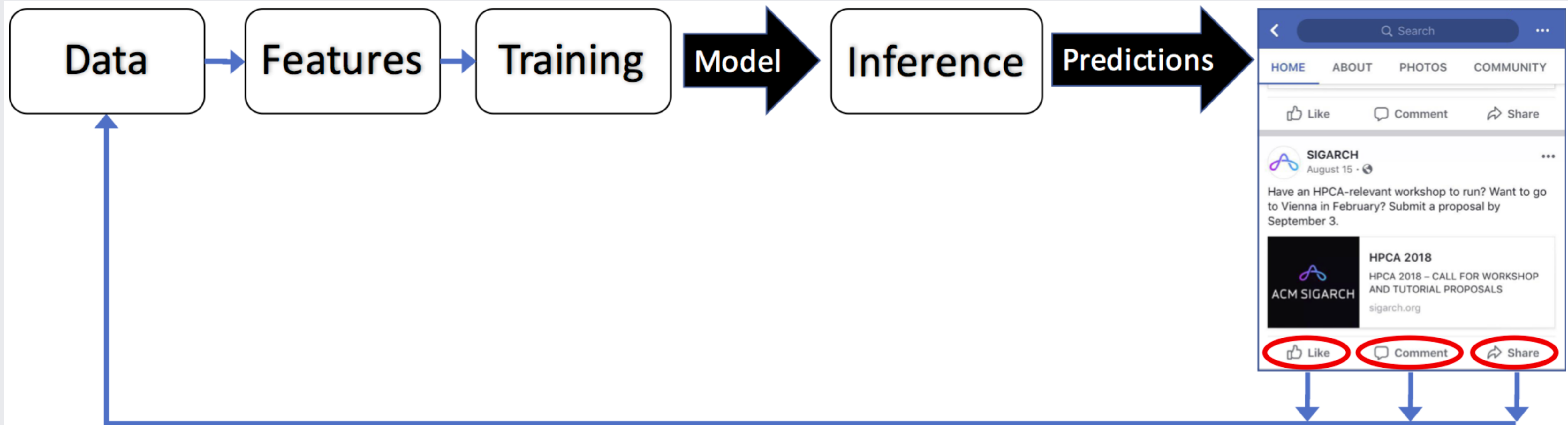
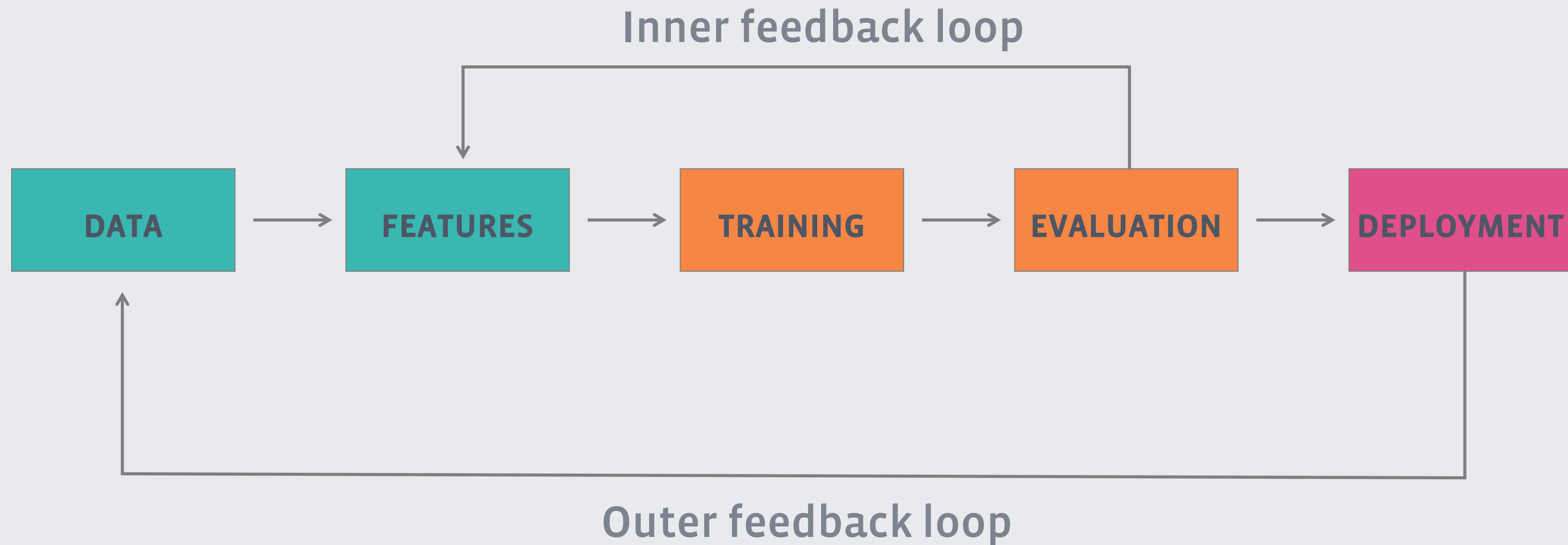


Fig. 1. Example of Facebook's Machine Learning Flow and Infrastructure.

offline

online

ML Pipeline Lifecycle



Agenda

1. ML Pipeline Lifecycle
2. ML development challenges
3. MLflow platform
4. Demo

“Hidden Technical Debt in Machine Learning Systems”

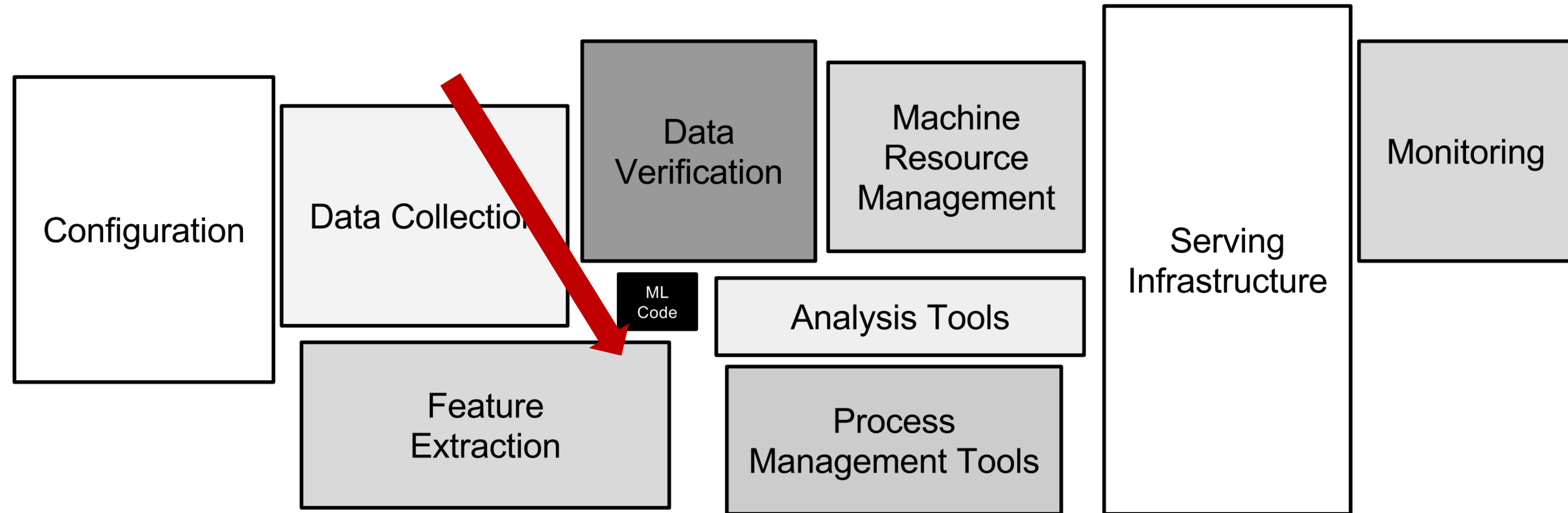


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

ML development challenges

1) Experiment Standardization and Tracking

- Reproducibility
- Parameters Tunning

ML development challenges

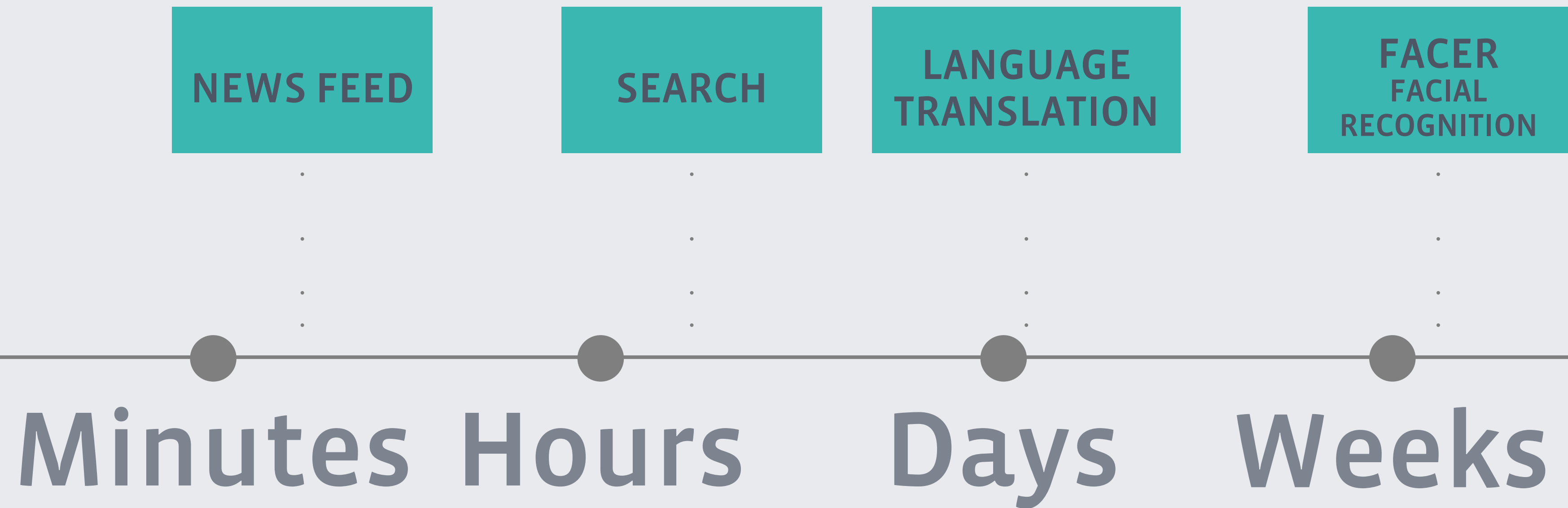
2) Model Staleness

ML pipelines need to be constantly run with new data to avoid model staleness and close feedback loop



Model Staleness

Training frequency for some model types at FB



ML development challenges

3) Productionization

- ML code is difficult to productionize
- Development environment <> production environment

Agenda

1. ML Pipeline Lifecycle
2. ML development challenges
3. MLflow platform
4. Demo



- An open platform for the machine learning lifecycle
- Python Library; runs locally and on the cloud
- Built-in UI for experiment visualization
- Logging integrations for major frameworks: scikit-learn, PyTorch, TF,...

<https://github.com/mlflow>



800k

monthly downloads



140+

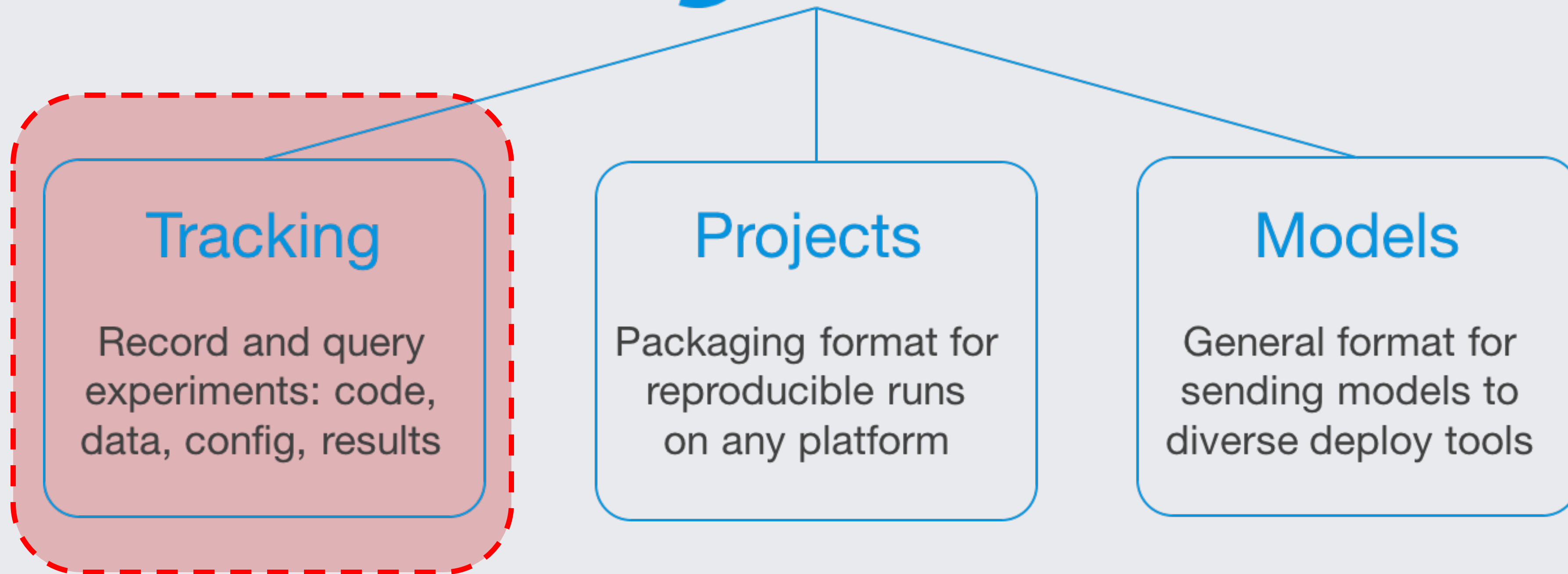
code contributors



40

contributing organizations

MLflow Components



MLflow Components

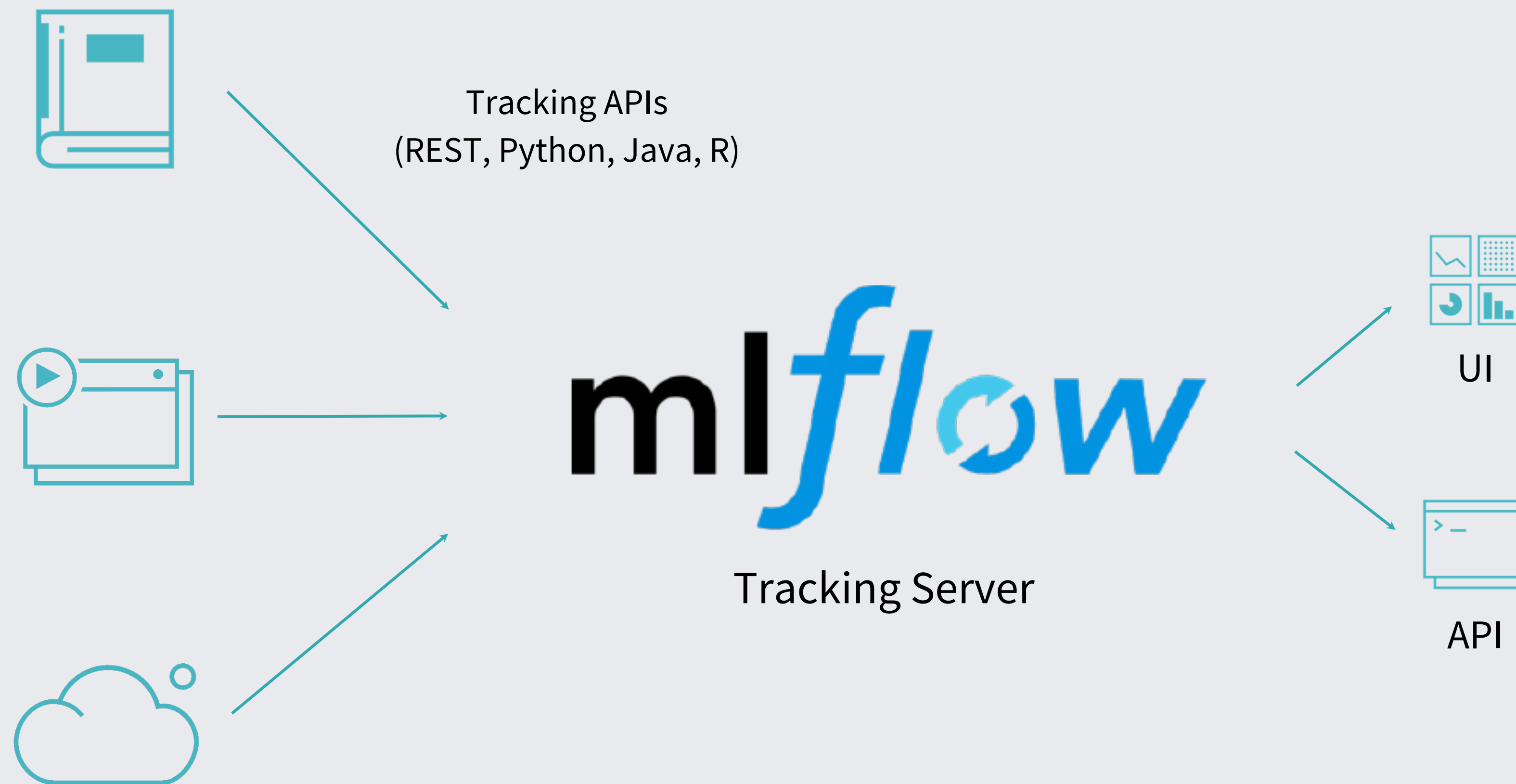
Experiment tracking

- **Hyper Parameters:** key-value inputs
- **Metrics:** numeric values (i.e. perf metrics)
- **Artifacts:** files, including data and models
- **Source:** training code

any additional information

MLflow Components

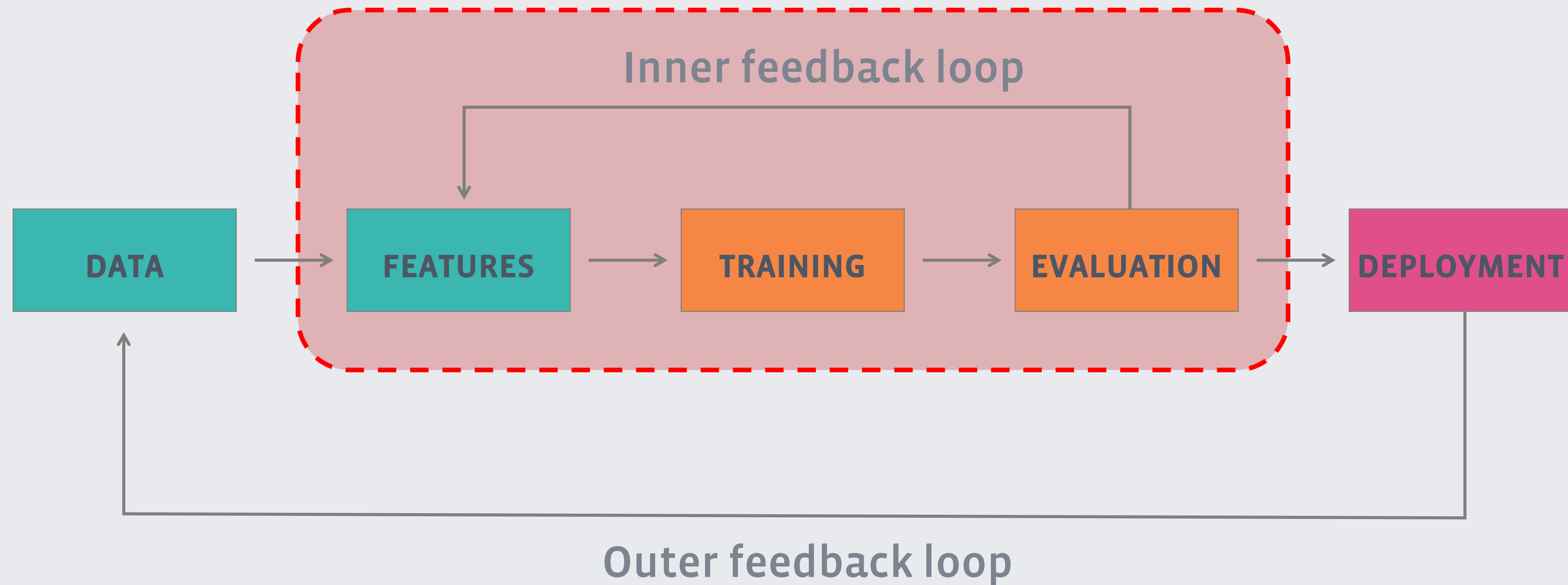
Experiment tracking



MLflow Components

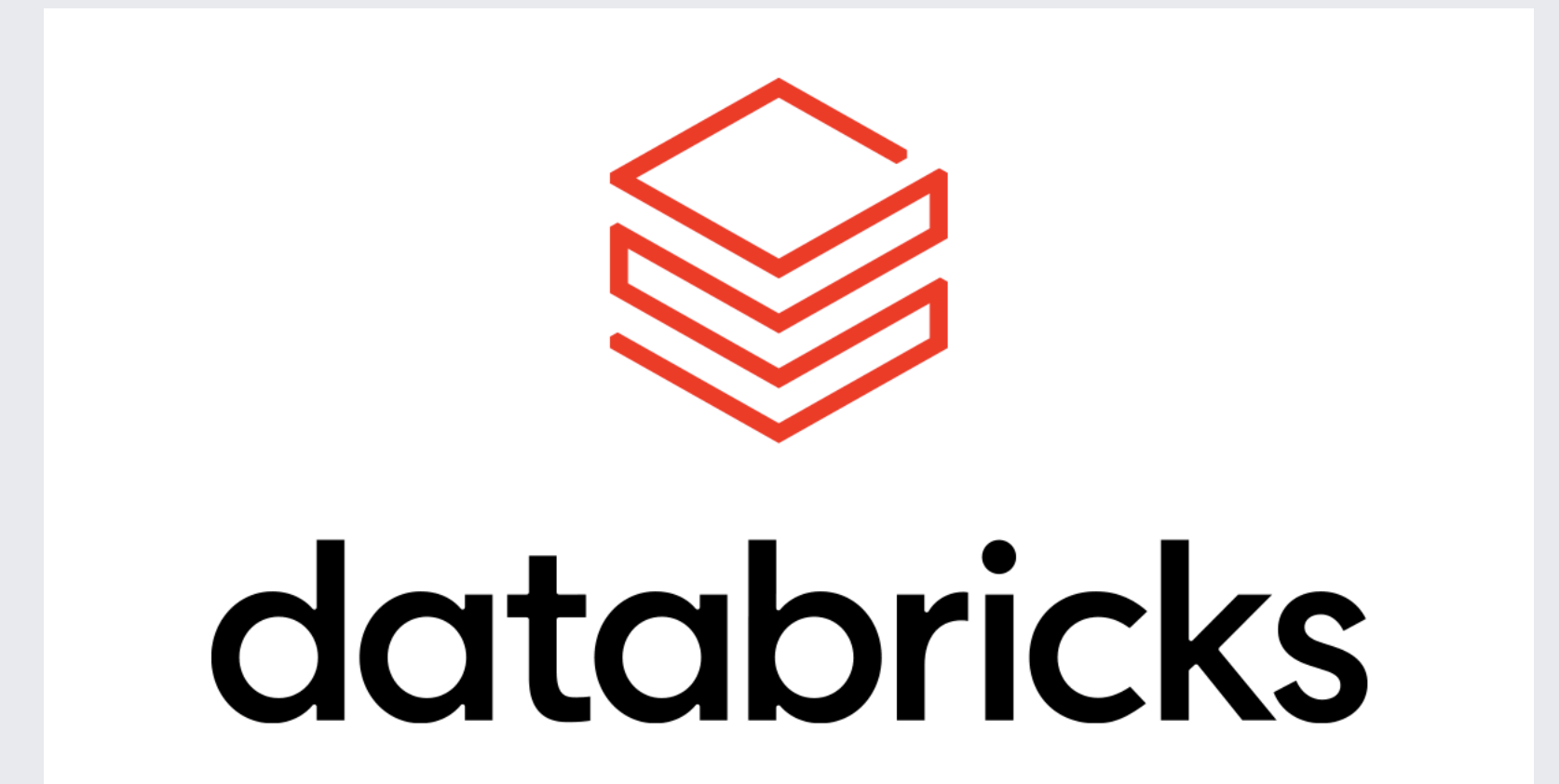


Experiment tracking



Getting Started with MLflow

- Install with `pip install mlflow`
- Find detailed tutorials at mlflow.org
- [Repo: https://github.com/mlflow](https://github.com/mlflow)
- Main contributor and maintainer:
- [Managed MLflow services offered by Databricks and Azure ML](#)





Demo

<https://github.com/alfozan/MLflow-GBRT-demo/blob/master/MLflow-GBRT-demo.ipynb>

Questions

Thank you