

Arrikto



# From Notebook to Kubeflow Pipelines to KFServing The Data Science Odyssey

A complete data science workflow,  
from notebooks to model serving  
with Kubeflow Pipelines, KF Serving, and Kale

Stefano Fioravanzo

Arrikto

Karl Weinmeister

Google



KubeCon



CloudNativeCon

North America 2020

*Virtual*

# From Notebook to Kubeflow Pipelines to KFServing The Data Science Odyssey



Stefano Fioravanzo  
Software Engineer



Karl Weinmeister  
Cloud AI Developer Advocacy Manager

# What You'll Learn In This Session

Run hyperparameter optimization with the click of a button, and serve the best result using KF Serving and Kale

Why is this important?

- ✓ **Simplify** your HP tuning and Serving workflows using intuitive UIs
- ✓ **Accelerate** your time to production. You can now reduce the training time and the time needed from training to serving
- ✓ **Collaborate** faster, reducing the friction between the data science team and the MLOps team



*Don't forget, you can grab the slides right now at [arrik.to/kubeconBOS](https://arrik.to/kubeconBOS) as well as enter the draw to win a fabulous prize*



*Get your questions answered **live** on Twitter and LinkedIn using the three hashtags [#kubecon](#) [#ml](#) [#arrikto](#)*



The Kubeflow project is dedicated to making deployments of machine learning (ML) workflows on Kubernetes: simple, portable and scalable.

- Deployment and management of a complex ML system at scale
- Rapid experimentation
- Hyperparameter tuning
- Hybrid and multi-cloud workloads
- Continuous integration and deployment (CI/CD)

## ML tools

Chainer

Jupyter

MPI

MXNet

PyTorch

scikit-learn

TensorFlow

XGBoost

Arrikto | Google

## Kubeflow applications and scaffolding

Jupyter notebook web app and controller

Hyperparameter tuning (Katib)

Kubeflow UI

Kale

Metadata

Training operators: MPI, MXNet, PyTorch, TFJob, XGBoost

Pipelines

KFServing

PyTorch Serving

Istio

TensorFlow Serving

Argo

Seldon Core

Prometheus

Spartakus

## Platforms / clouds

Kubernetes

GCP

AWS

Azure

IBM Cloud

OpenShift

On prem

## ML tools

Chainer

Jupyter

MPI

MXNet

Arrikto | Google

PyTorch

scikit-learn

TensorFlow

XGBoost

## Kubeflow applications and scaffolding

Jupyter notebook web app and controller

Hyperparameter tuning (Katib)

Kubeflow UI

Kale

Metadata

Training operators: MPI, MXNet, PyTorch, TFJob, XGBoost

Pipelines

KFServing

PyTorch Serving

Istio

TensorFlow Serving

Argo

Seldon Core

Prometheus

Spartakus

## Platforms / clouds

Kubernetes

GCP

AWS

Azure

IBM Cloud

OpenShift

On prem

## ML tools

Chainer

Jupyter

MPI

MXNet

Arrikto | Google

PyTorch

scikit-learn

TensorFlow

XGBoost

## Kubeflow applications and scaffolding

Jupyter notebook web app and controller

Hyperparameter tuning (Katib)

PyTorch Serving

Istio

Kubeflow UI

Kale

TensorFlow Serving

Argo

Metadata

Seldon Core

Prometheus

Training operators:  
MPI, MXNet, PyTorch,  
TFJob, XGBoost

Pipelines

KFServing

Spartakus

## Platforms / clouds

Kubernetes

GCP

AWS

Azure

IBM Cloud

OpenShift

On prem



## ML tools

Chainer

Jupyter

MPI

MXNet

Arrikto | Google

PyTorch

scikit-learn

TensorFlow

XGBoost

## Kubeflow applications and scaffolding

Jupyter notebook web app and controller

Hyperparameter tuning (Katib)

PyTorch Serving

Istio

Kubeflow UI

Kale

TensorFlow Serving

Argo

Metadata

Seldon Core

Prometheus

Training operators:  
MPI, MXNet, PyTorch,  
TFJob, XGBoost

Pipelines

KFServing

Spartakus

## Platforms / clouds

Kubernetes

GCP

AWS

Azure

IBM Cloud

OpenShift

On prem

# ML workflow

Identify problem  
and collect and  
analyse data

Choose an ML  
algorithm and  
code your model

Experiment with  
data and model  
training

Tune the model  
hyperparameters

Serve the model  
for online/batch  
prediction

Jupyter  
Notebook

PyTorch

scikit-learn

TensorFlow

XGBoost

Jupyter  
Notebook

Kale

Pipelines

Katib

KFServing

NVIDIA TensorRT

PyTorch

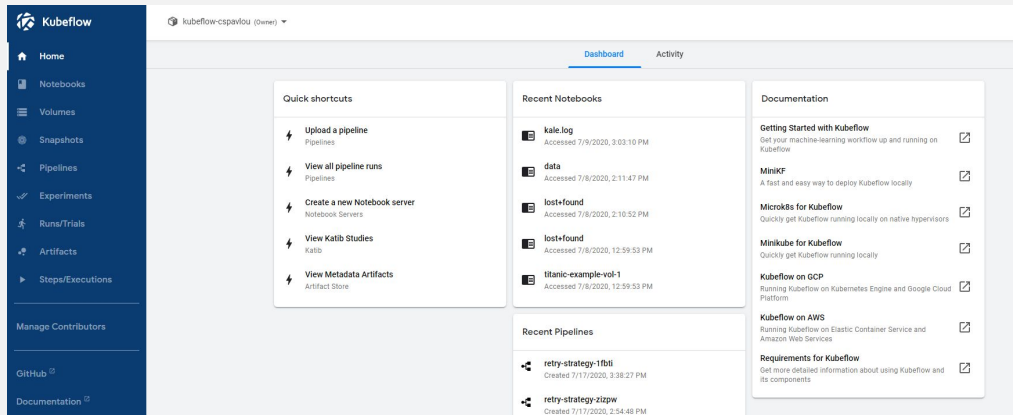
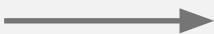
TF Serving

Seldon Core



# Interacting with Kubeflow

## User interface (UI)








kfctl CLI

kubectrl CLI






APIs and SDKs



## Quick shortcuts

-  **Upload a pipeline**  
Pipelines
-  **View all pipeline runs**  
Pipelines
-  **Create a new Notebook server**  
Notebook Servers
-  **View Katib Studies**  
Katib
-  **View Metadata Artifacts**  
Artifact Store








## Recent Notebooks

-  **kale.log**  
Accessed 7/9/2020, 3:03:10 PM
-  **data**  
Accessed 7/8/2020, 2:11:47 PM
-  **lost+found**  
Accessed 7/8/2020, 2:10:52 PM
-  **lost+found**  
Accessed 7/8/2020, 12:59:53 PM
-  **titanic-example-vol-1**  
Accessed 7/8/2020, 12:54:48 PM

## Recent Pipelines

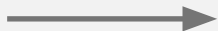
-  **retry-strategy-1fbti**  
Created 7/17/2020, 3:38:27 PM
-  **retry-strategy-zizpw**  
Created 7/17/2020, 2:54:48 PM

## Documentation

- Getting Started with Kubeflow**  
Get your machine-learning workflow up and running on Kubeflow 
- MinikF**  
A fast and easy way to deploy Kubeflow locally 
- Microk8s for Kubeflow**  
Quickly get Kubeflow running locally on native hypervisors 
- Minikube for Kubeflow**  
Quickly get Kubeflow running locally 
- Kubeflow on GCP**  
Running Kubeflow on Kubernetes Engine and Google Cloud Platform 
- Kubeflow on AWS**  
Running Kubeflow on Elastic Container Service and Amazon Web Services 
- Requirements for Kubeflow**  
Get more detailed information about using Kubeflow and its components 

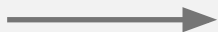
User interface (UI)

**kfctl CLI**



```
kfctl apply -V -f ${CONFIG_URI}
```

**kubectl CLI**



```
kubectl -n kubeflow get all
```

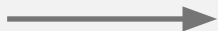
APIs and SDKs

User interface (UI)

kfctl CLI

kubectrl CLI

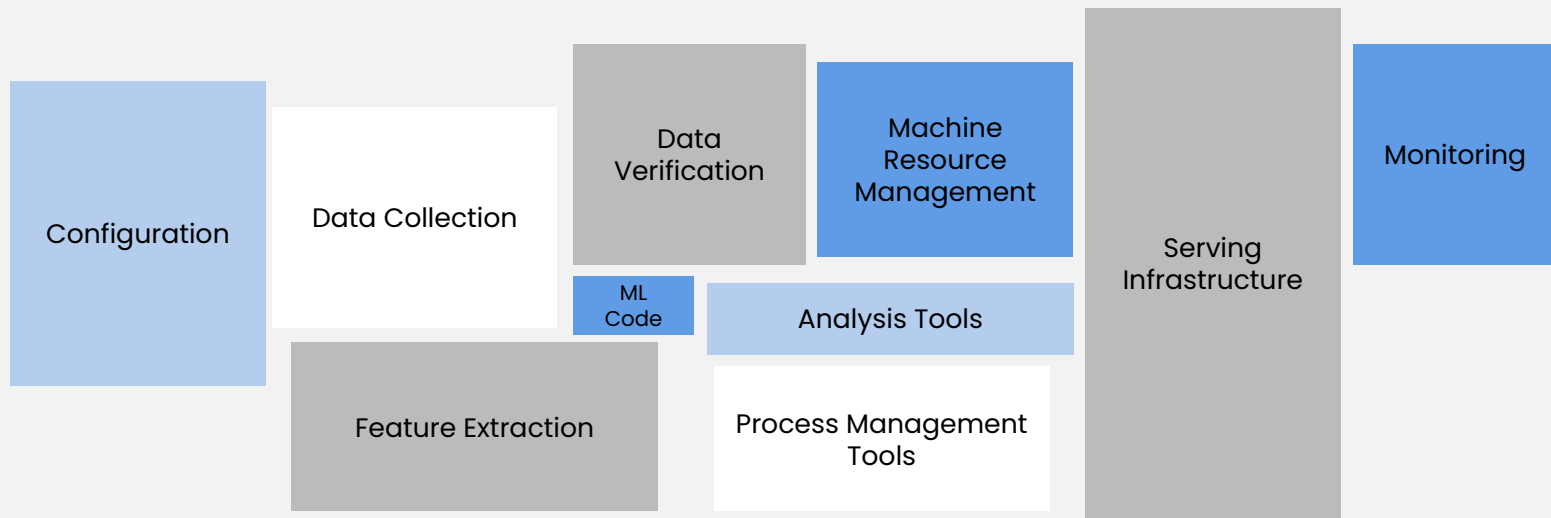
**APIs and SDKs**



Examples:

- Pipelines SDK
- Katib API
- Metadata SDK

# ML Applications are Distributed Systems

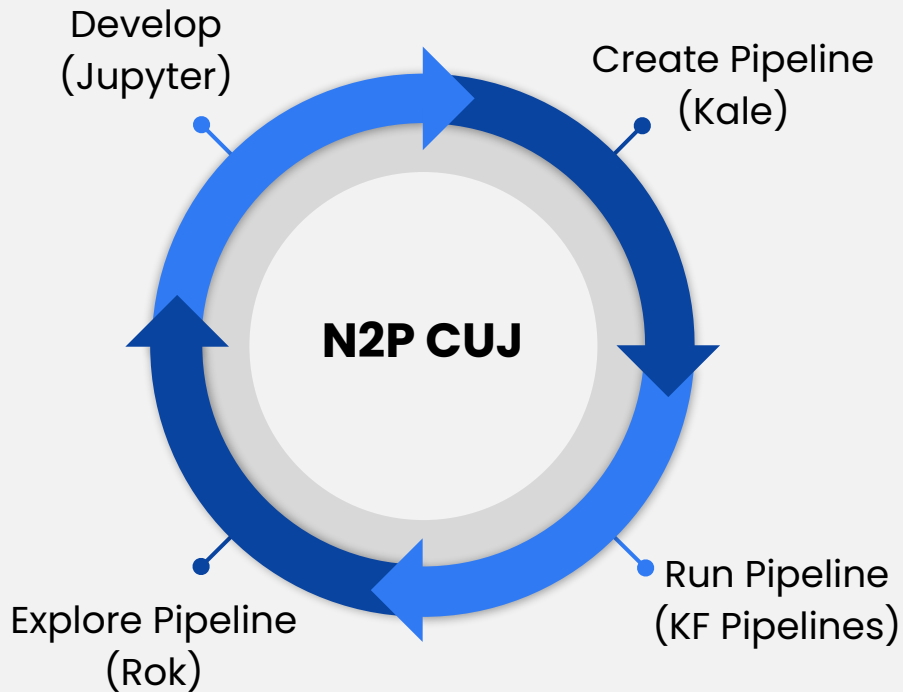


Credit: Hidden Technical Debt of Machine Learning Systems, D. Sculley, et al.

# CI/CD for ML

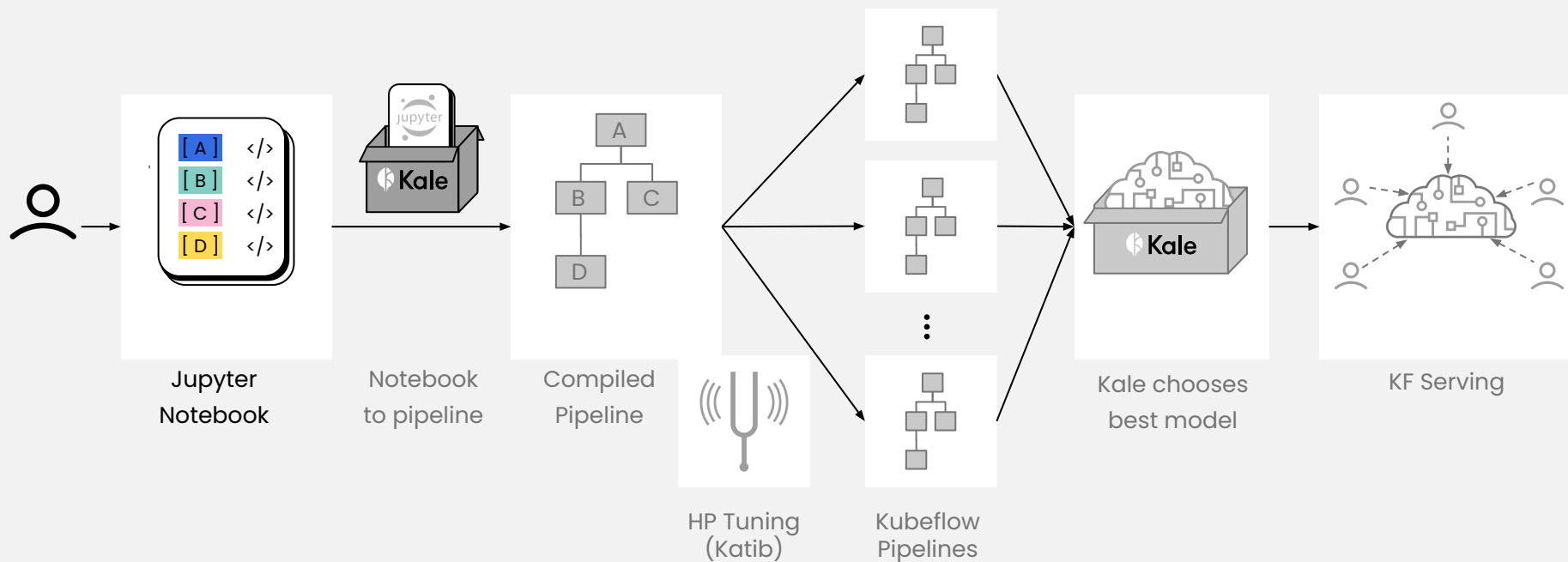
How can data scientists continually improve and validate models?

- Develop models and pipelines in Jupyter
- Convert notebook to pipeline using Kale
- Run pipeline using Kubeflow Pipelines
- Explore and debug pipelines using Rok



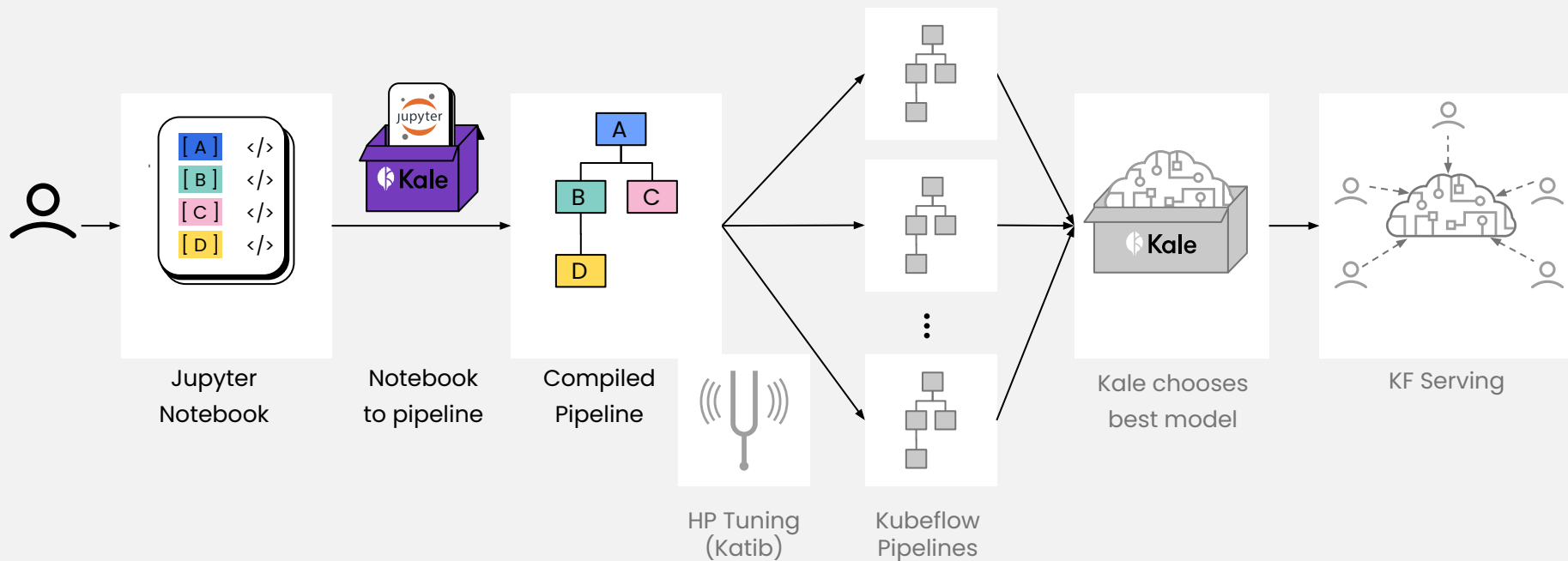


# Data Science Workflows



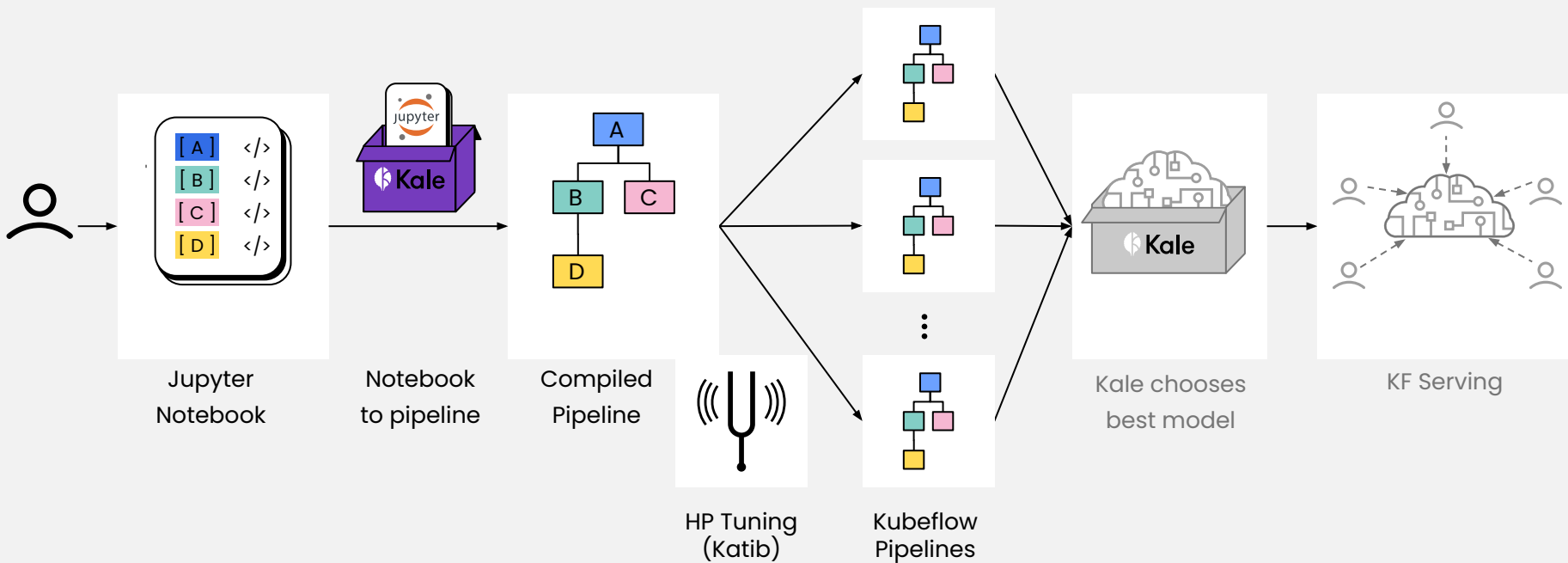
- A Step 1
- B Step 2
- C Step 3
- D Step 4

# Data Science Workflows



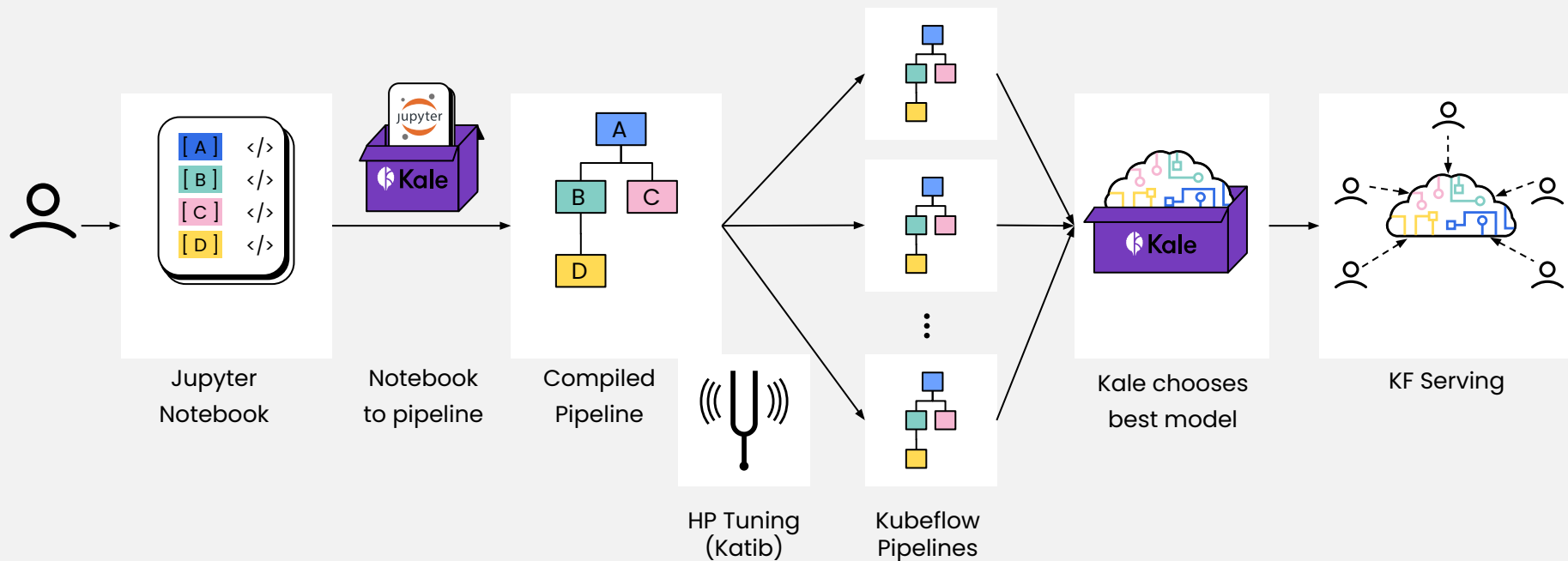
- A** Step 1
- B** Step 2
- C** Step 3
- D** Step 4

# Data Science Workflows



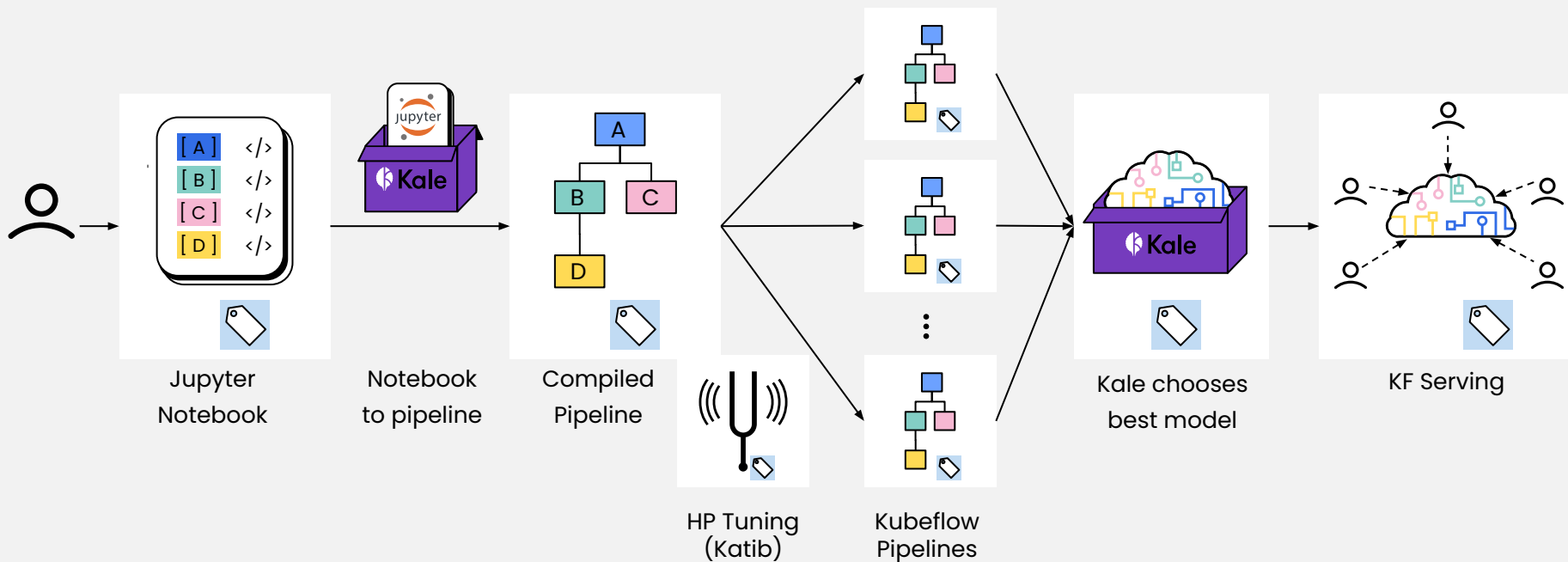
- A Step 1
- B Step 2
- C Step 3
- D Step 4

# Data Science Workflows



- A Step 1
- B Step 2
- C Step 3
- D Step 4

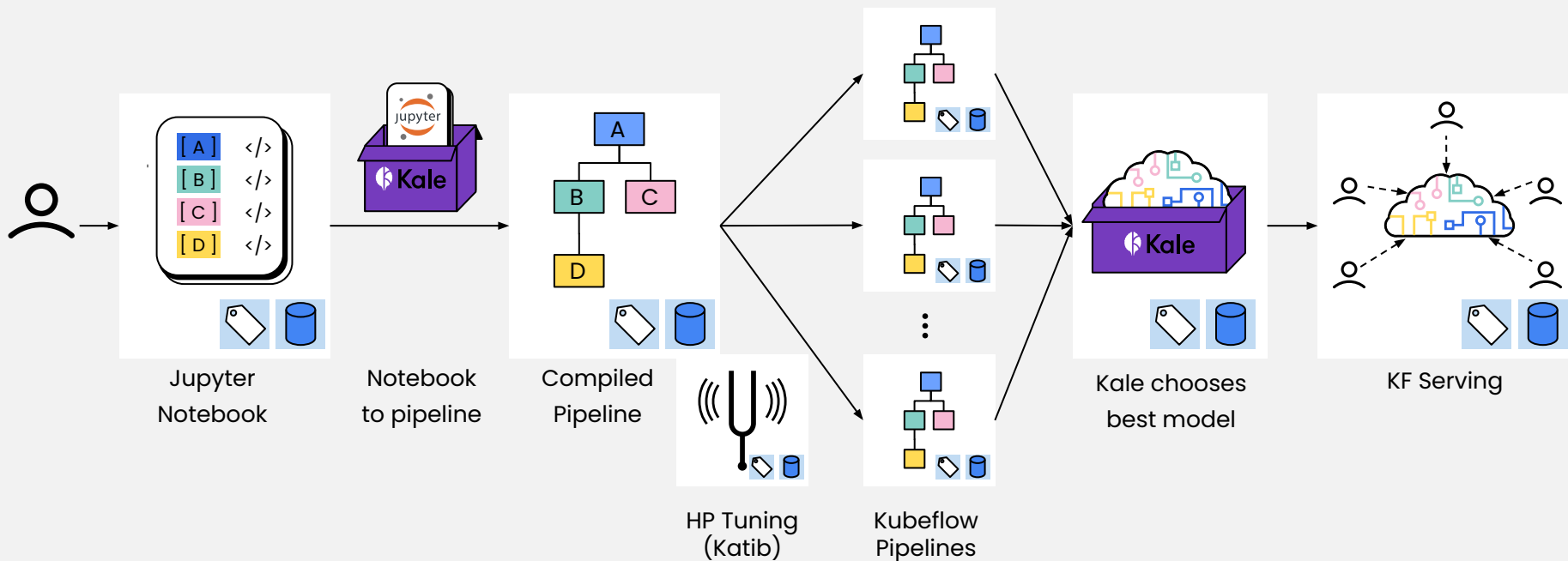
# Data Science Workflows



- A Step 1
- B Step 2
- C Step 3
- D Step 4



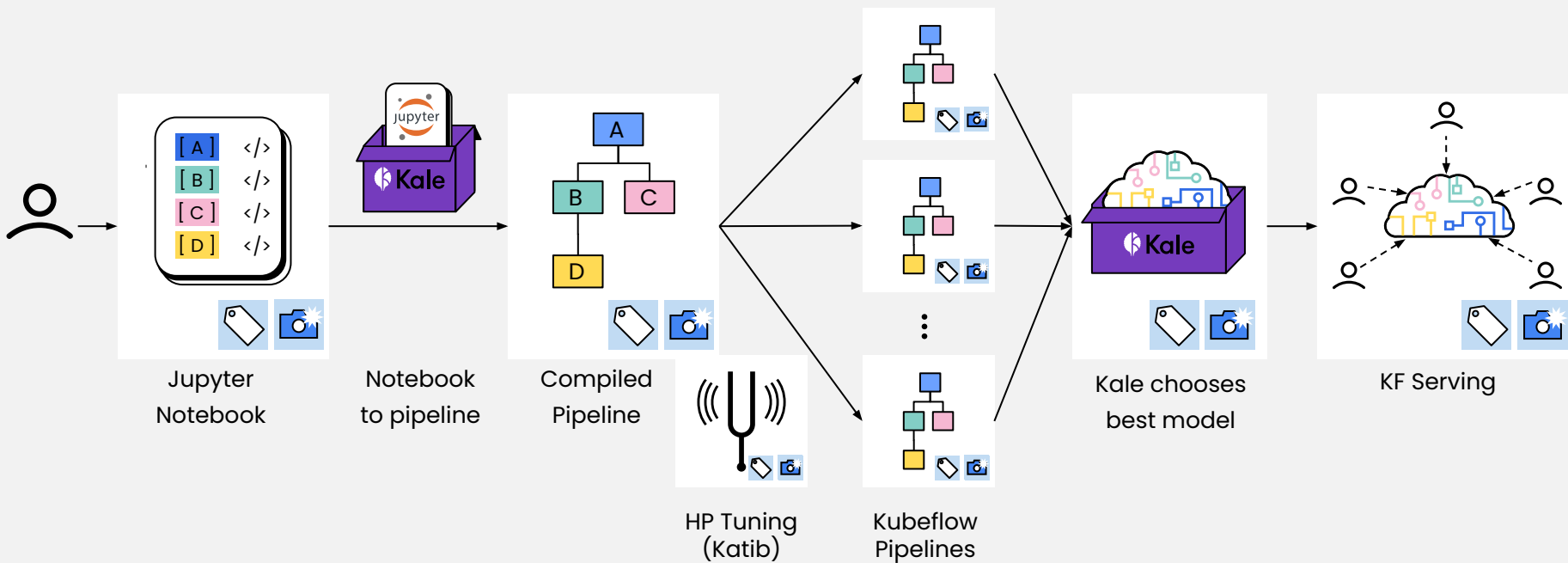
# Data Science Workflows



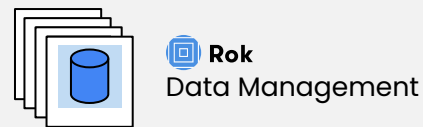
- A Step 1
- B Step 2
- C Step 3
- D Step 4



# Data Science Workflows



- A Step 1
- B Step 2
- C Step 3
- D Step 4



# Agenda



Go to [arrik.to/democ2p](https://arrik.to/democ2p) to find the Codelab with the step-by-step instructions for this tutorial

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

6

Q&A



# Agenda

1

**Install MiniKF**

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

6

Q&A



- Kubeflow on GCP, your laptop, or on-prem infrastructure in just a few minutes
- All-in-one, single-node, Kubeflow distribution
- Very easy to spin up on your own environment on-prem or in the cloud
- MiniKF = MiniKube + Kubeflow + Arrikto's Rok Data Management Platform

# Demo - Install MiniKF



# Agenda

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

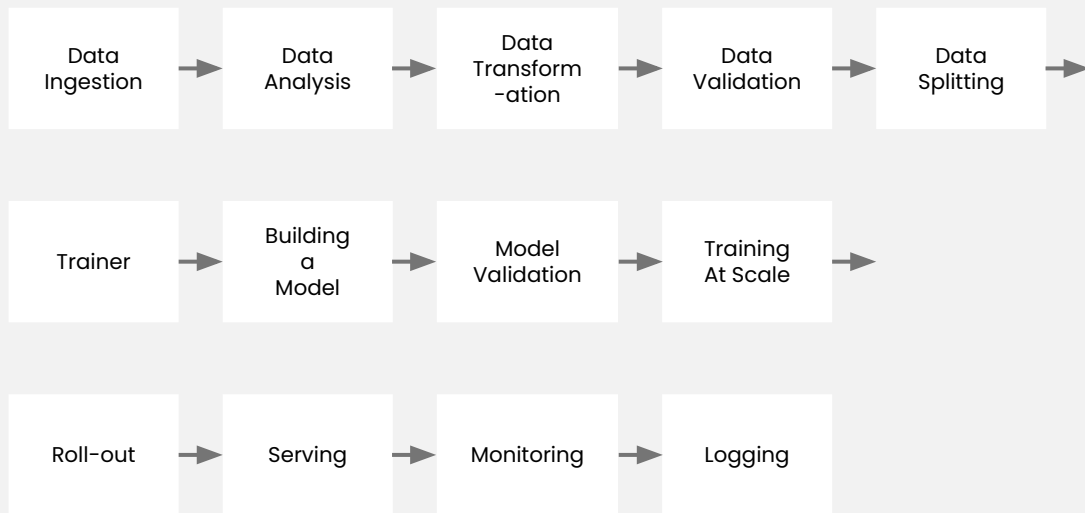
6

Q&A

**Kubeflow Pipelines** exists because Data Science and ML are inherently **pipeline processes**

This workshop will focus on two essential aspects:

- **Low barrier to entry:** deploy a Jupyter Notebook to Kubeflow Pipelines in the Cloud using a fully GUI-based approach
- **Reproducibility:** automatic data versioning to enable reproducibility and better collaboration between data scientists



**Kubeflow Pipelines** exists because Data Science and ML are inherently **pipeline processes**

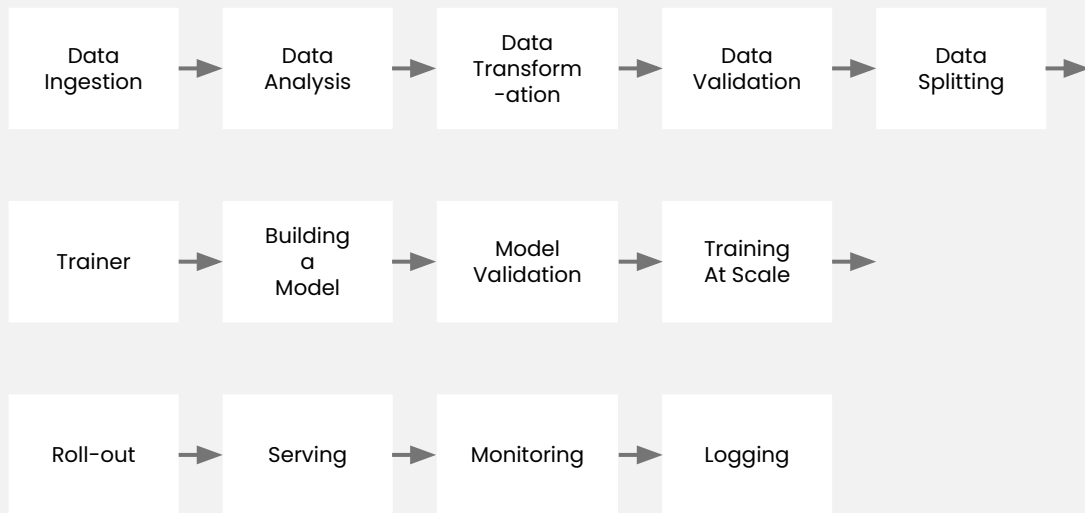
This workshop will focus on two essential aspects:

- **Low barrier to entry:** deploy a Jupyter Notebook to Kubeflow Pipelines in the Cloud using a fully GUI-based approach



- **Reproducibility:** automatic data versioning to enable reproducibility and better collaboration between data scientists

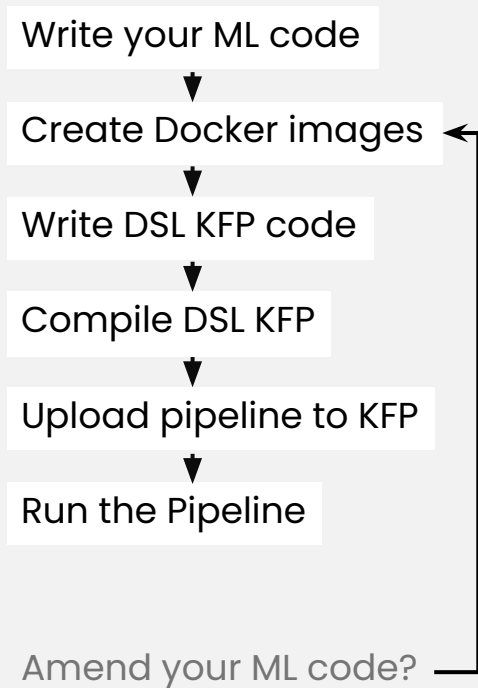
**Arrikto**



# Benefits of running a Notebook as a Pipeline

- The steps of the workflow are clearly defined
- Parallelization & isolation
  - Hyperparameter tuning
- Data versioning
- Different infrastructure requirements
  - Different hardware (GPU/CPU)

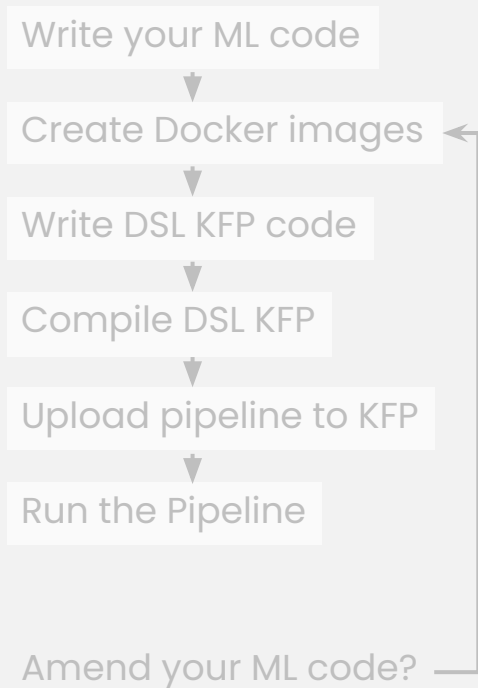
## Before



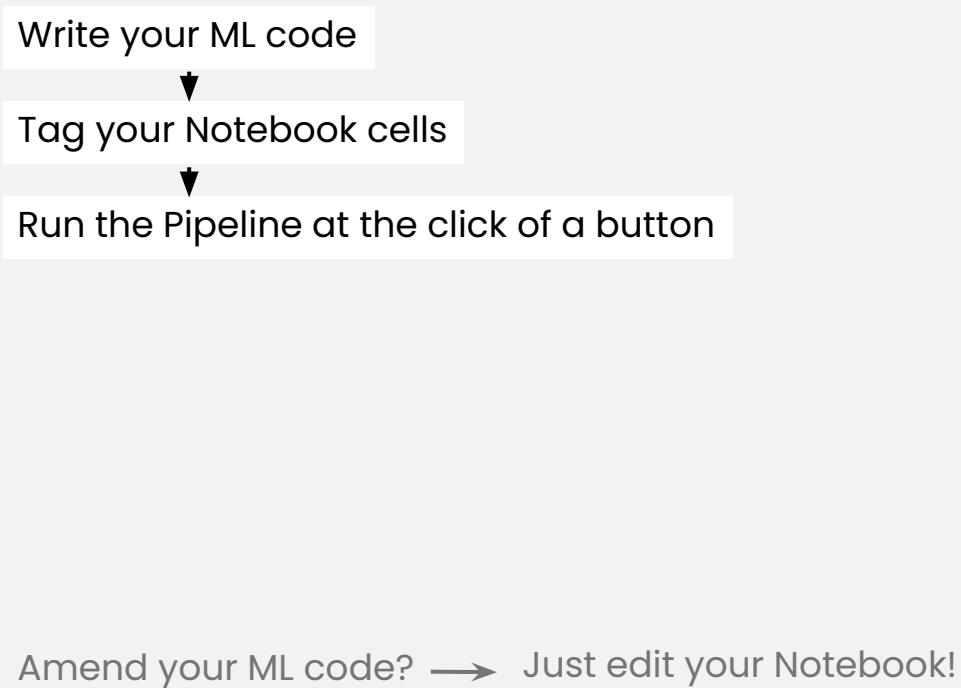


# Workflow

## Before

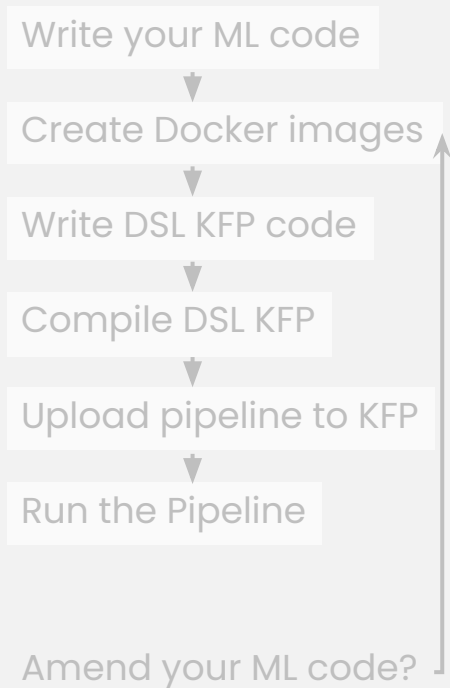


## After

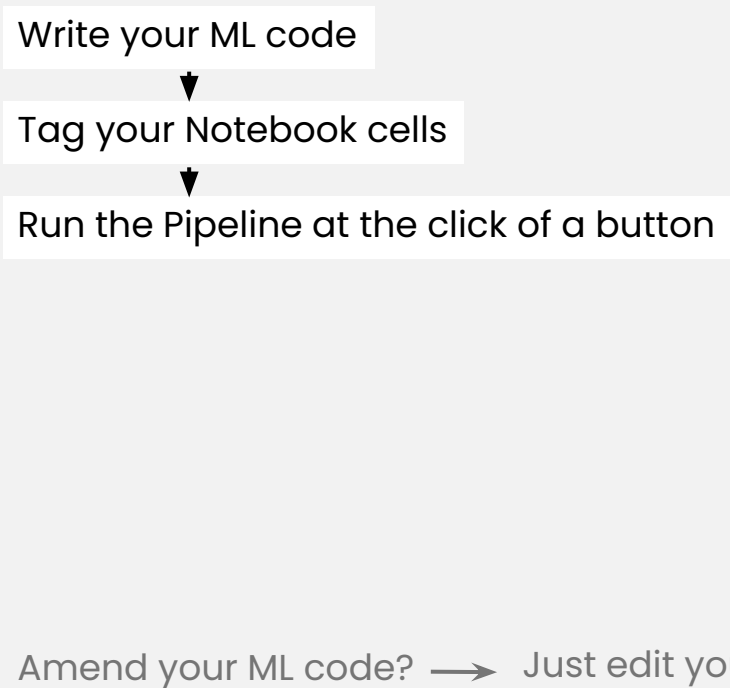


# Workflow

## Before



## After



A Data Scientist can now reduce the time taken to write ML code and run a pipeline by 70%.

That means you can now run 3x as many experiments as you did before.

What that really means is that you can deliver work faster to the business and drive more revenue

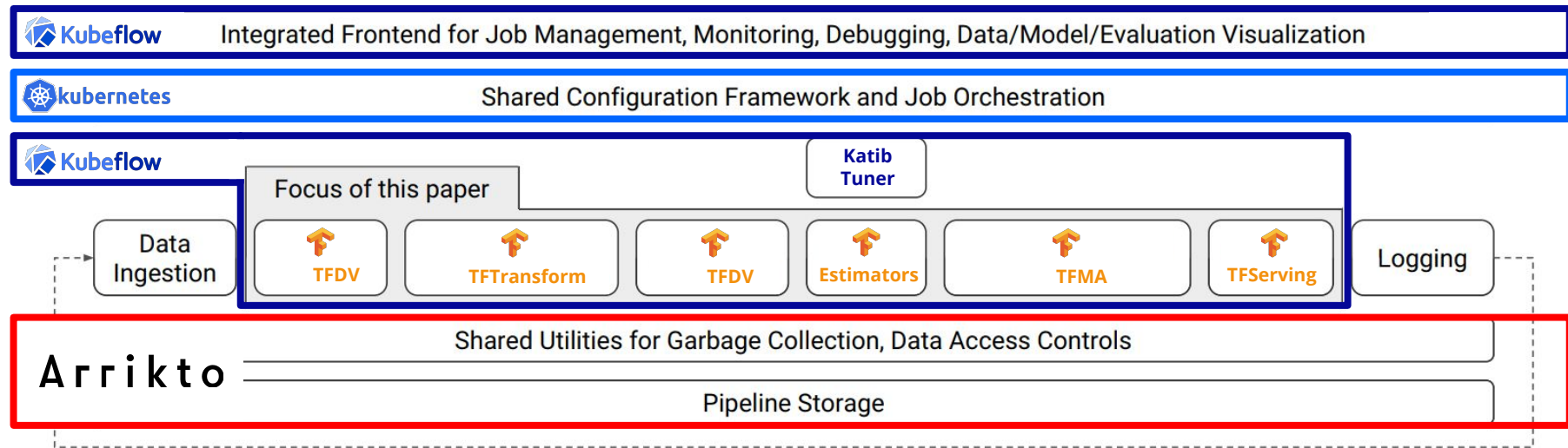
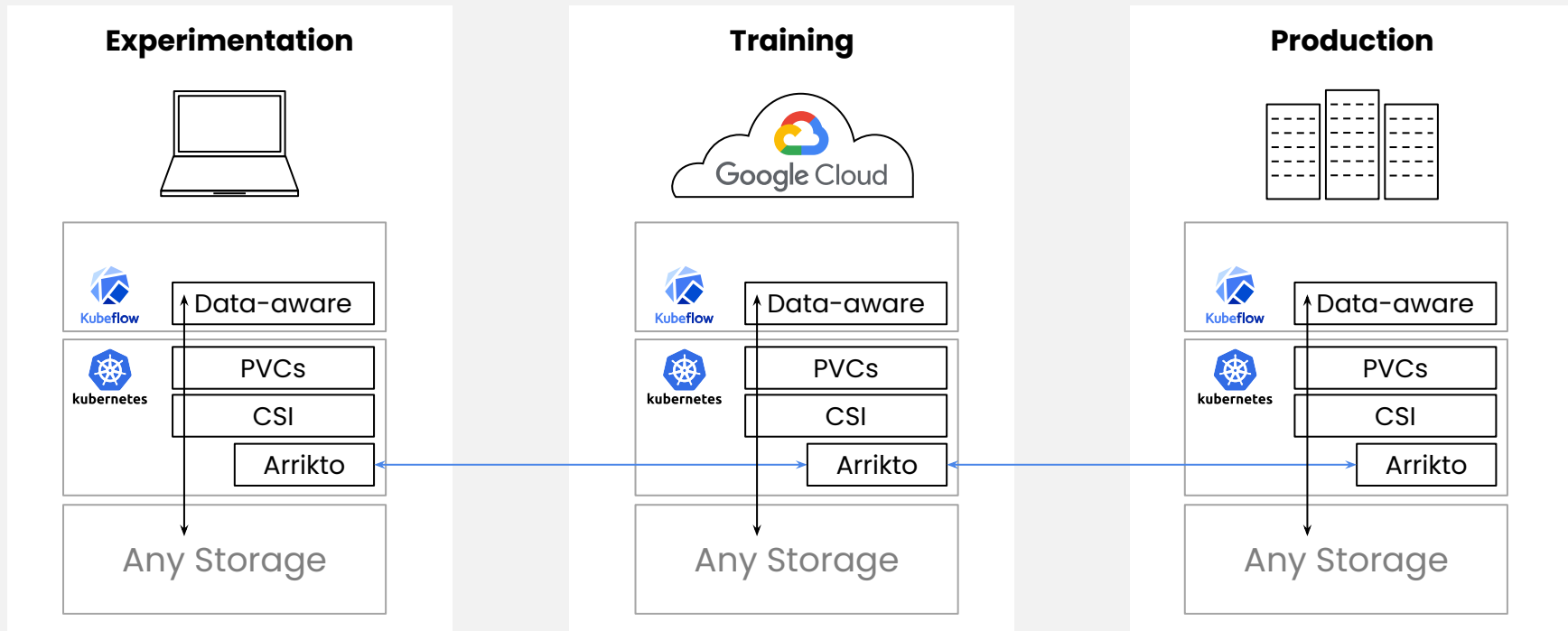


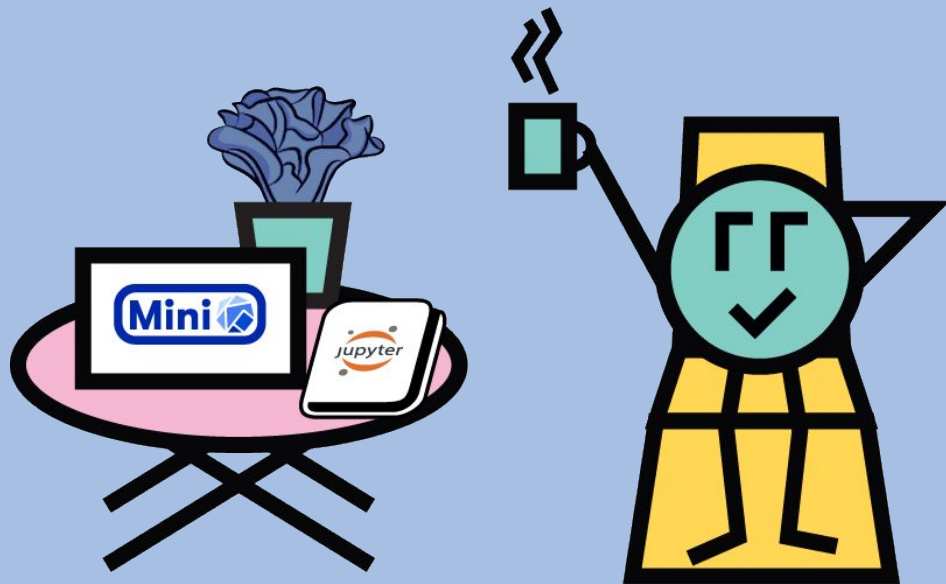
Figure 1: High-level component overview of a machine learning platform.

Data Versioning, Packaging, and Sharing

Across teams and cloud boundaries for complete Reproducibility, Provenance, and Portability



# Demo - Notebook to Pipelines



# Agenda

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

**Notebook to Katib  
User Journey**

4

Notebook to Serving  
User Journey

5

Summary

6

Q&A

# Hyperparameter optimization

## The two ways of life

- Change the parameters manually
- Use Katib



K a t i b

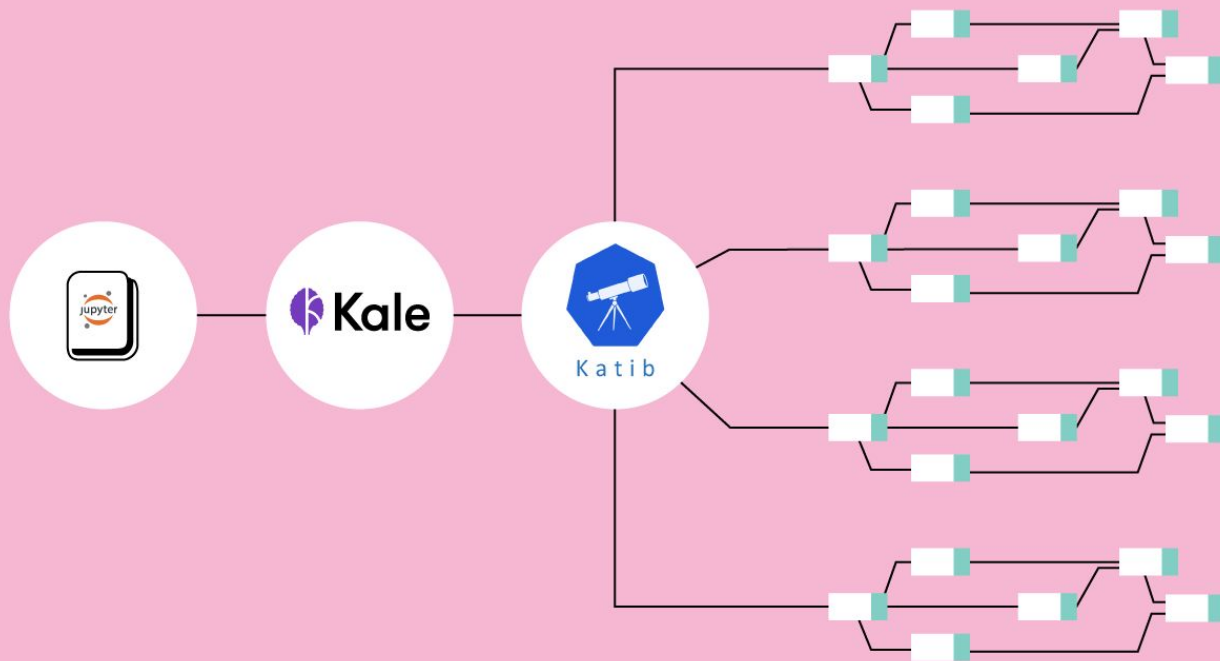
Katib is a Kubernetes-based system for Hyperparameter Tuning and Neural Architecture Search. It supports a number of ML frameworks, including TensorFlow, Apache MXNet, PyTorch, XGBoost, and others.



## Combining the N2P CUJ with Katib

- Configure parameters, search algorithm, and objectives using a GUI
- Start HP tuning with the click of a button
- Reproducibility of every pipeline and every step
- Run Katib Trials as Pipelines
- Complete visibility of every different Katib Trial
- Caching for faster computation

# Demo - Notebook to Katib



# Agenda

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

6

Q&A

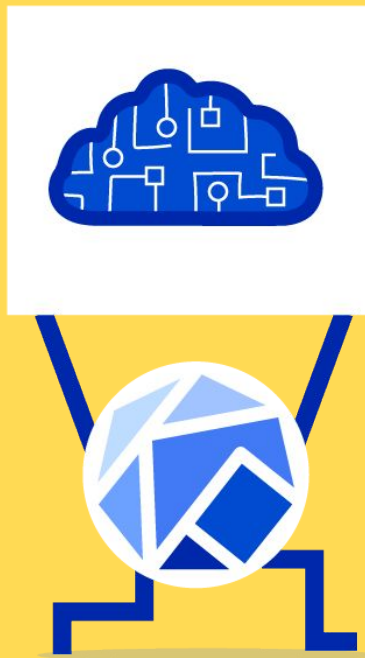
KFServing enables serverless inferencing on Kubernetes and provides performant, high abstraction interfaces for common machine learning (ML) frameworks like TensorFlow, XGBoost, scikit-learn, PyTorch, and ONNX to solve production model serving use cases.

# Serving from a notebook

Kale provides a simple to use API to serve a model

- Choose the best Trial of a HP Tuning experiment
- Restore a notebook from a Rok snapshot
- Create and deploy InferenceService CRs with a convenient API
- No need to build new Docker images
- Run predictions directly from the notebook

# Demo - Notebook to Serving



# Agenda

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

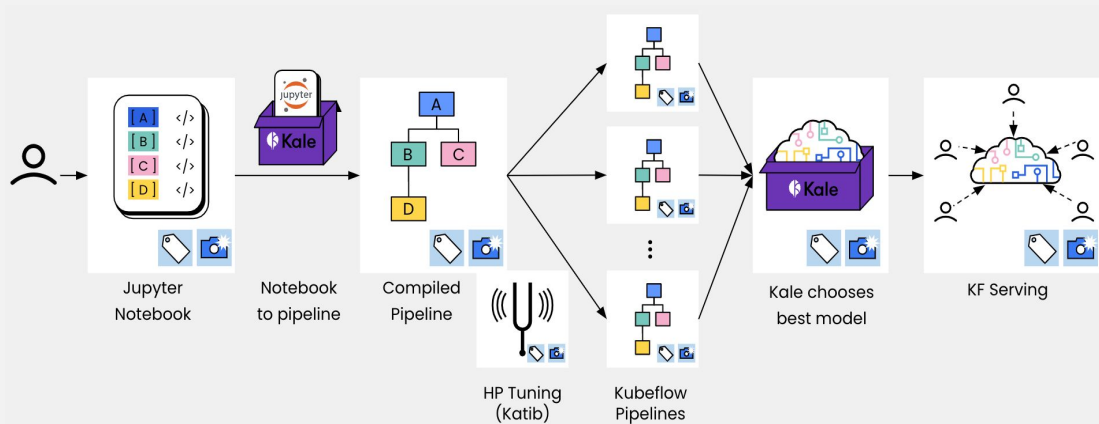
6

Q&A

# Summary

What you have learned during this tutorial:

- Run a pipeline-based hyperparameter tuning workflow starting from your Jupyter Notebook
- Use Kale as a workflow tool to orchestrate Katib and KubeFlow Pipelines experiments
- Run and monitor model servers directly from your notebook
- Use the new and intuitive HP Tuning and Models UI
- Navigate between KubeFlow UIs, across linked entities tracked by MLMD
- **Simplify** your ML workflows using intuitive UIs
- Exploit the caching feature so that you **accelerate** your pipeline runs
- **Collaborate** faster and more easily, and have complete visibility of your training





# Just a small sample of community contributions

- Jupyter manager UI
- Pipelines volume support
- MiniKF
- Auth with Istio + Dex
- On-premise installation
- Linux Kernel

# Community

## Kubeflow is open

- Open community
- Open design
- Open source
- Open to ideas

## Get involved

- [github.com/kubeflow](https://github.com/kubeflow)
- [kubeflow.slack.com](https://kubeflow.slack.com)
- [@kubeflow](https://twitter.com/kubeflow)
- [kubeflow-discuss@googlegroups.com](mailto:kubeflow-discuss@googlegroups.com)
- Community call on Tuesdays



# Thank you, team!



Ilias Katsakioris,  
Dimitris Pouloupoulos,  
Kimonas Sotirchos,  
Apostolos Plakias,  
Konstantinos Palaiologos,  
Chris Pavlou

# Thank You



More Info

[cloud.google.com](https://cloud.google.com)



More Info

[arrik.to/kubeconBOS](https://arrik.to/kubeconBOS)



[google](https://twitter.com/google)



[arrikto](https://twitter.com/arrikto)



[linkedin.com/in/karlweinmeister](https://linkedin.com/in/karlweinmeister)



[linkedin.com/in/stefanofioravanzo](https://linkedin.com/in/stefanofioravanzo)



Email Address:

**kweinmeister@google.com**



Email Address:

**stefano@arrikto.com**

# Agenda

1

Install MiniKF

2

Notebook to Pipelines  
User Journey

3

Notebook to Katib  
User Journey

4

Notebook to Serving  
User Journey

5

Summary

6

Q&A