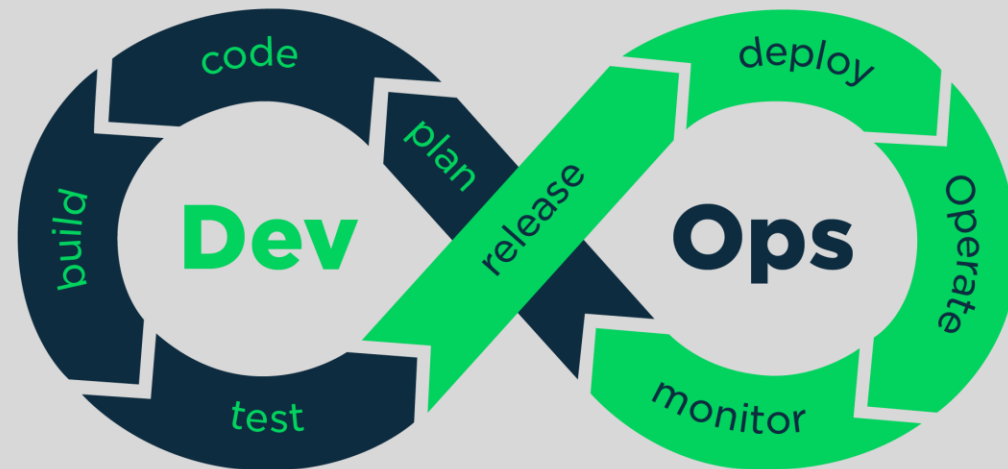# MLOPS & KUBEFLOW

Grant Stevens | University of Bristol

# WHAT IS MLOPS?

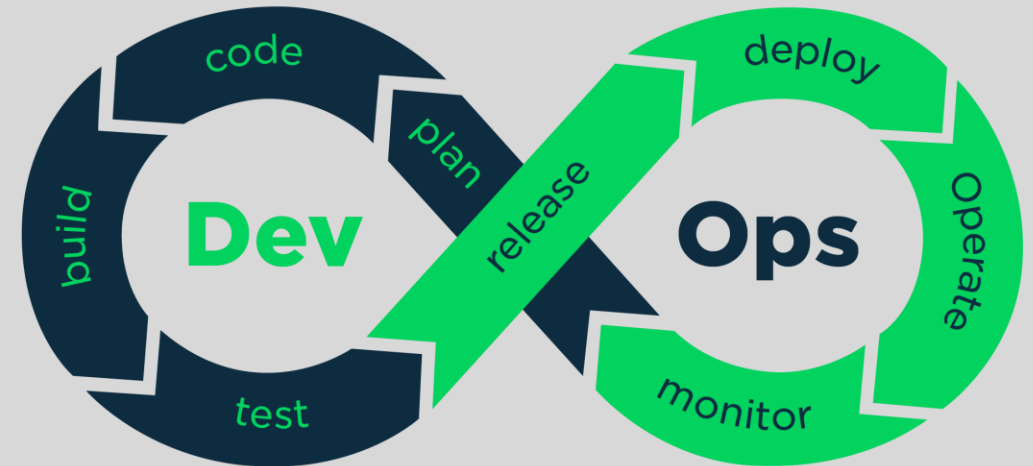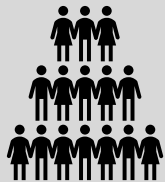# WHAT IS ~~MLOPS?~~ DEVOPS?

# What is DevOps?

DevOps is a set of practices that works to automate and integrate the processes between software development and IT teams, so they can build, test, and release software faster and more reliably.

# Increasing Scale

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

# Increasing Scale

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility
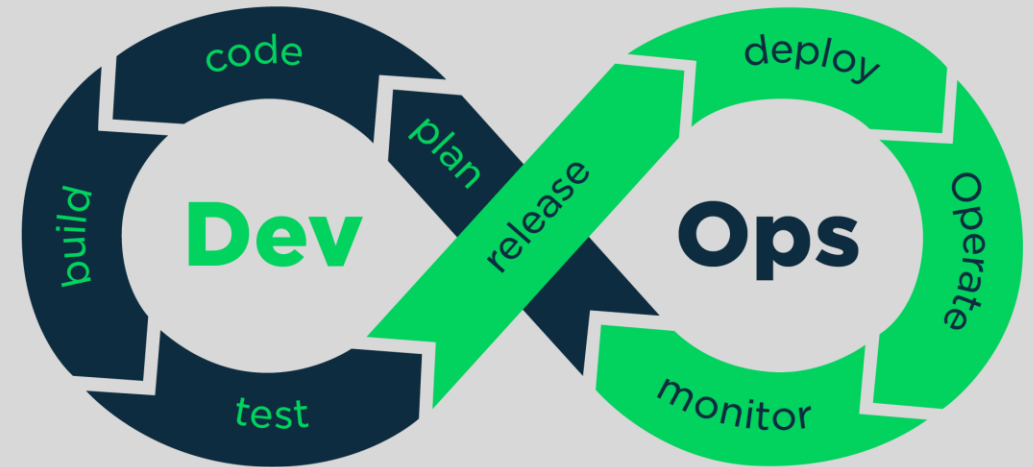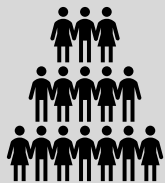
# Increasing Scale

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility

How do you develop good software at industry scale with potential hundreds of developers?

# Why Start Here?

- The key questions that the ML community are now asking were asked by the software development industry a few decades ago.

- The result of many years of creating, optimising and testing solutions has led to what we now call DevOps.

- It is important to see how these challenges were overcome and why these solutions that work well for software development, are not sufficient for the development and deployment of ML applications.
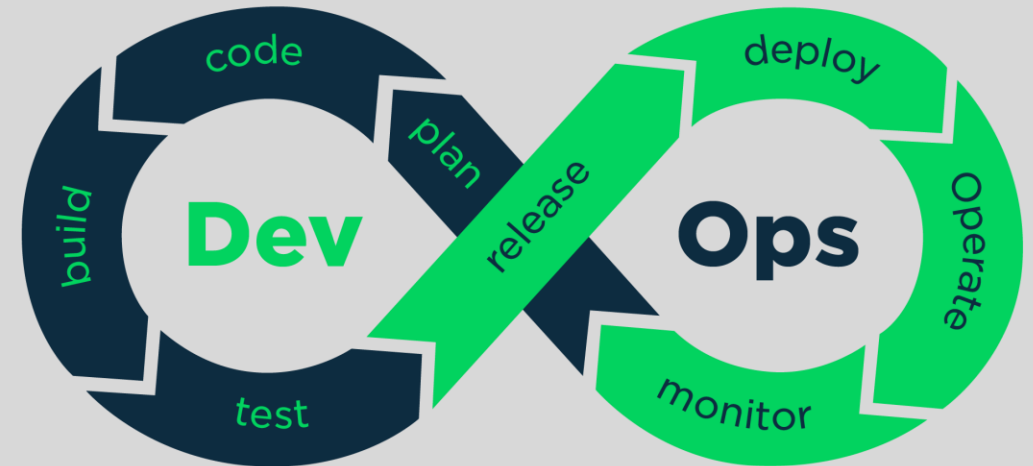
# Increasing Scale

## Software Development

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility

How do you develop and deploy good software at industry scale with potentially hundreds of developers?

## Machine Learning

# Increasing Scale

## Software Development

- No Comments
- No Version Control
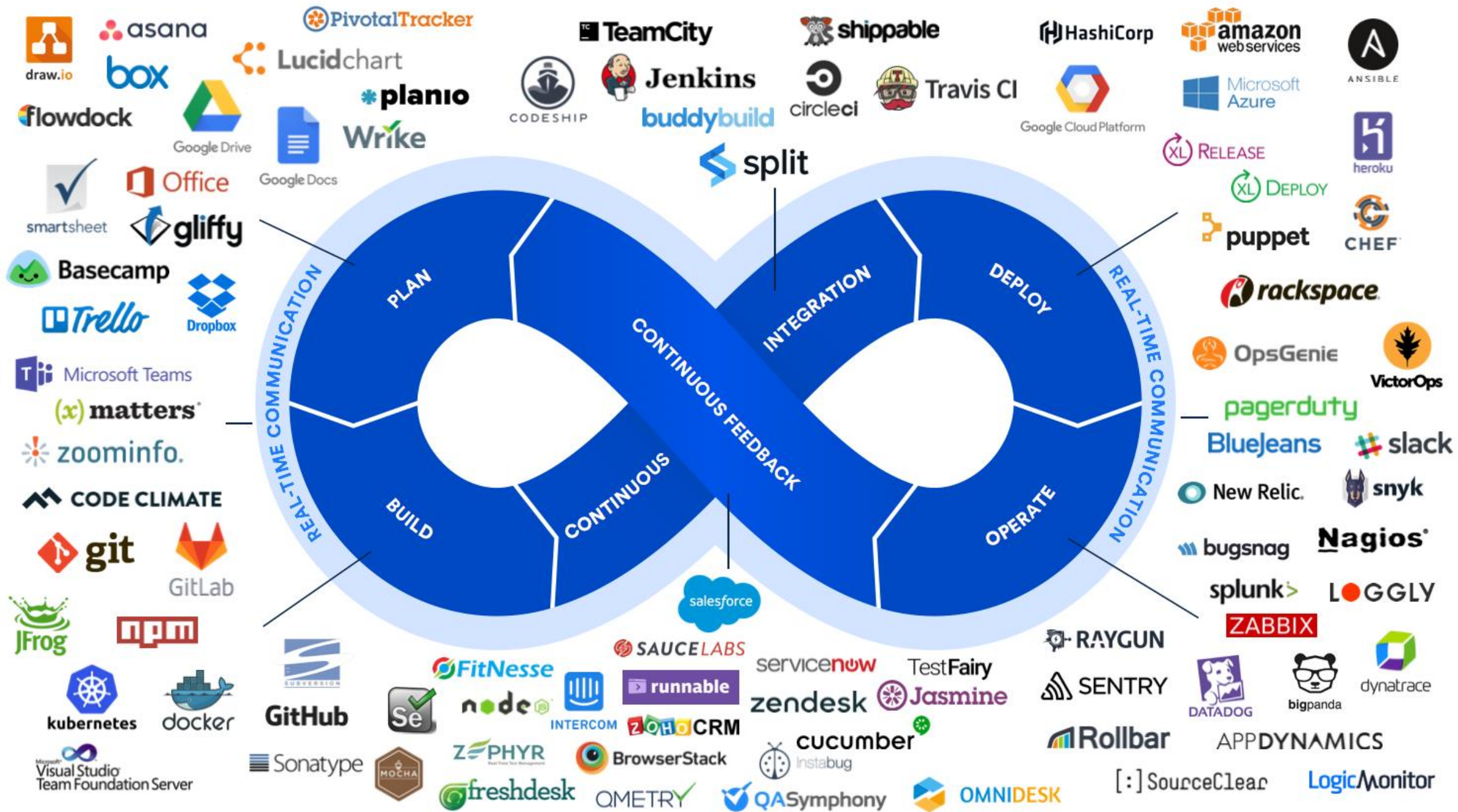- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility

How do you develop and deploy good software at industry scale with potentially hundreds of developers?

## Machine Learning

Much of the work you have done/will do over the next few years will be here.

# Increasing Scale

## Software Development

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility

How do you develop and deploy good software at industry scale with potentially hundreds of developers?

## Machine Learning

Much of the work you have done/will do over the next few years will be here.

The Group Project will give you an introductory experience doing ML in a team.

# Increasing Scale

## Software Development

- No Comments
- No Version Control
- Minimal/Local Planning
- Skip Compatibility Checks

- Communication
- Shared Codebase
- Task Assignment
- Compatibility

How do you develop and deploy good software at industry scale with potentially hundreds of developers?

## Machine Learning

Much of the work you have done/will do over the next few years will be here.

The Group Project will give you an introductory experience doing ML in a team.

How do you develop and deploy accurate and reliable ML models at industry scale with potentially hundreds of developers?

# Is MLOps just more DevOps?

# Is MLOps just more DevOps?

Software = compile(code, environment)

# Is MLOps just more DevOps?

Software = compile(code, environment)

Code → Test → Deploy

Monitor

# Is MLOps just more DevOps?

Software = compile(code, environment)

whereas

Model = train(data, params, code, environment)

Code → Test → Deploy

Monitor

# Is MLOps just more DevOps?

Software = compile(code, environment)

whereas

Model = train(data, params, code, environment)

Code → Test → Deploy

Monitor

Data runs & features   Code   Parameters

Model runs

Models / metrics → Deploy

Monitor

# Is MLOps just more DevOps?

Software = compile(code, environment)

whereas

Model = train(data, params, code, environment)

Train is:
- Sometimes non-deterministic
- Usually more expensive than compile (limiting the number of changes)
- Sometimes distributed across multiple machines (adds complexity)

Code → Test → Deploy → Monitor

Data runs & features, Code, Parameters → Model runs → Models / metrics → Deploy → Monitor

# Potential Issues

- If you have hundreds of models being trained at once, how do you choose which one to use in production?

- How do you know that multiple people aren't testing the same hyperparameters?

- If you only have capacity for training 10 models at one time, who should get priority?

- How do you share models?

- Are there the equivalent of unit tests for ML models?

- How do you effectively bridge the gap between Data Scientists, who likely won't know much of the software design, and the Software Developers who likely won't know much of the ML side of things?

And many more...

What are the key parts of
infrastructure required for ML systems?

Data
Collection

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

Feature Extraction

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

ML Code

Feature Extraction

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

ML Code

Analysis Tools

Feature Extraction

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Data Collection

Data Verification

Machine Resource Management

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Configuration

Data Collection

Data Verification

Machine Resource Management

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Configuration

Data Collection

Data Verification

Machine Resource Management

Serving Infrastructure

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Configuration

Data
Collection

Data
Verification

Machine
Resource
Management

Monitoring

ML Code

Analysis Tools

Serving
Infrastructure

Feature
Extraction

Process
Management
Tools

Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.
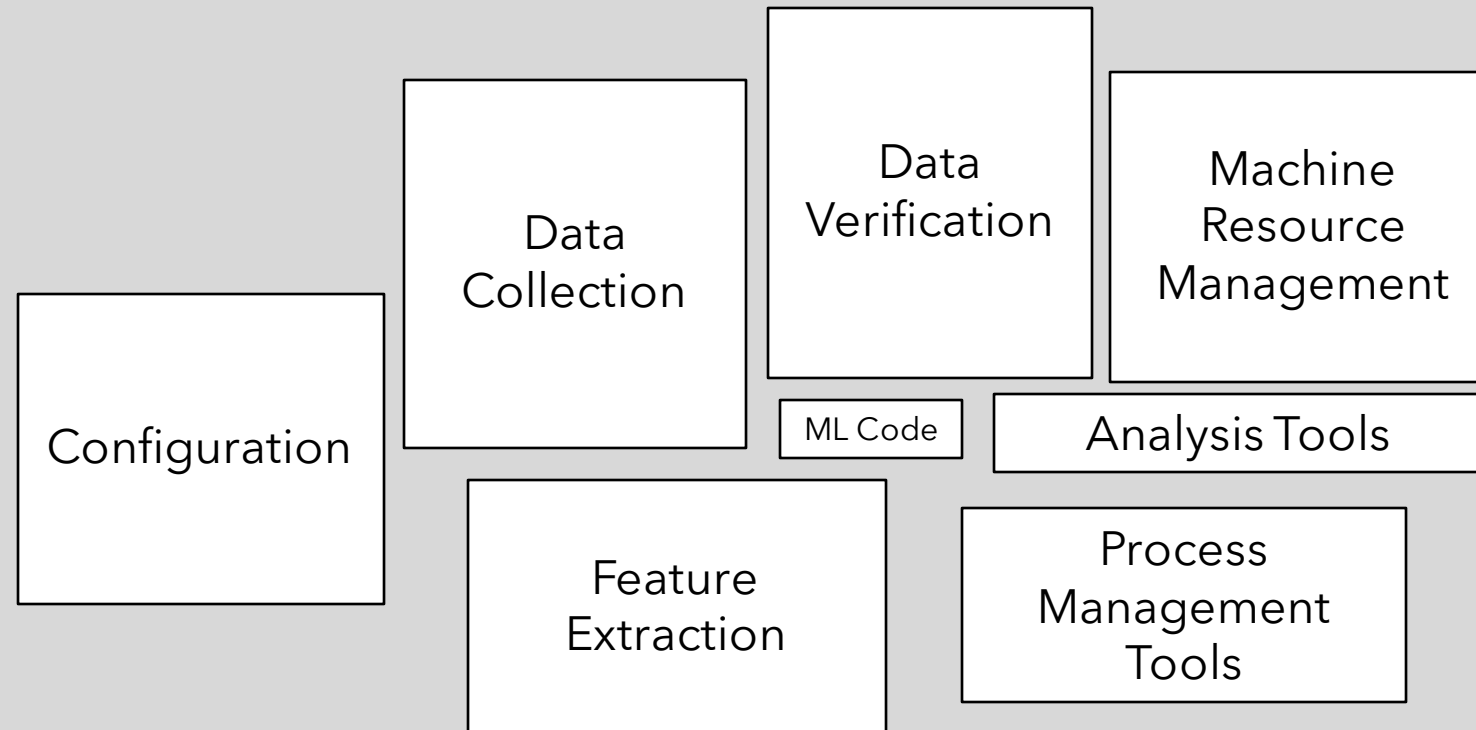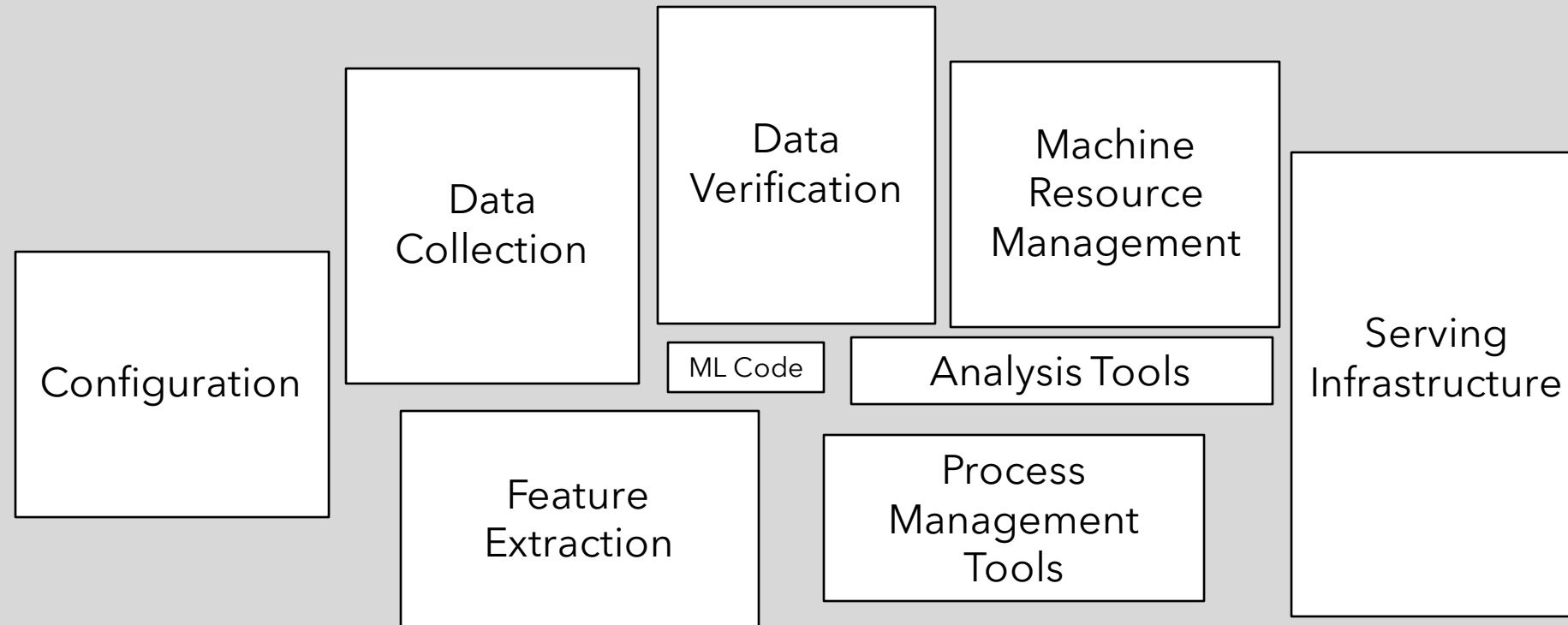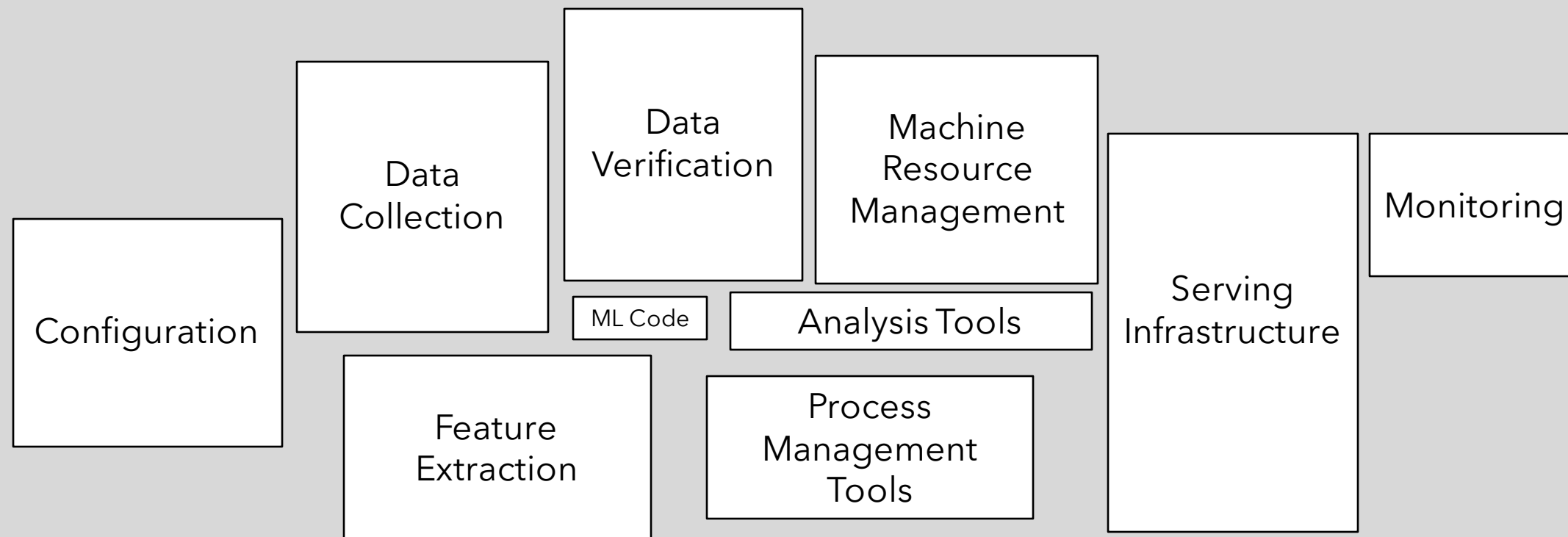
# Kubeflow

Kubeflow is an open source project that contains a curated set of compatible tools and frameworks specific for ML.

Its built on top of Kubernetes, allowing it to run consistently across different environments.

Kubeflow is built around 3 principles:
- Composability (you choose what works for you)
- Portability (Run any part of your workflow wherever you are running Kubeflow)
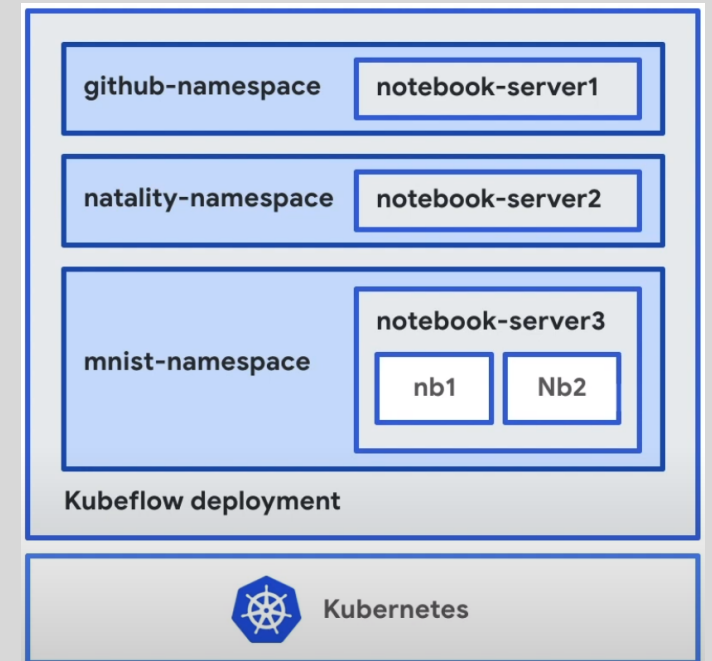- Scalability (Your project can access more resources when needed and release them when not)

**Kubeflow**

# Jupyter Notebooks

Kubeflow has Jupyter Notebooks built into the system.

Using the notebooks on Kubeflow (rather than locally) allows for the following benefits:
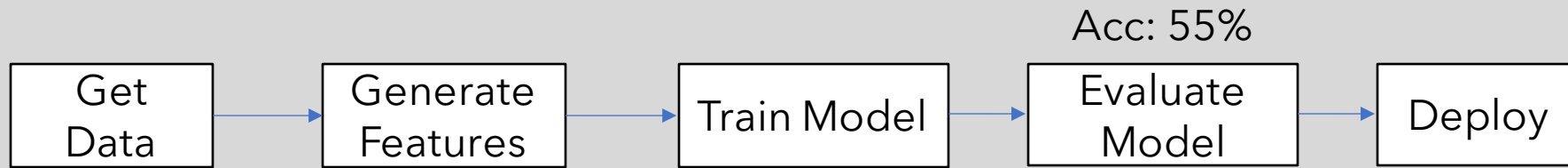
- Integration with the other Kubeflow components

- Access control and authentication (both for developers and users)

- Automated resource allocation

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.

Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

Acc: 55%

| Get Data | → | Generate Features | → | Train Model | → | Evaluate Model | → | Deploy |

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.

Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

Acc: 55%

```
Get Data  →  Generate Features  →  Train Model  →  Evaluate Model  →  Deploy
```

Review Inaccuracies

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.
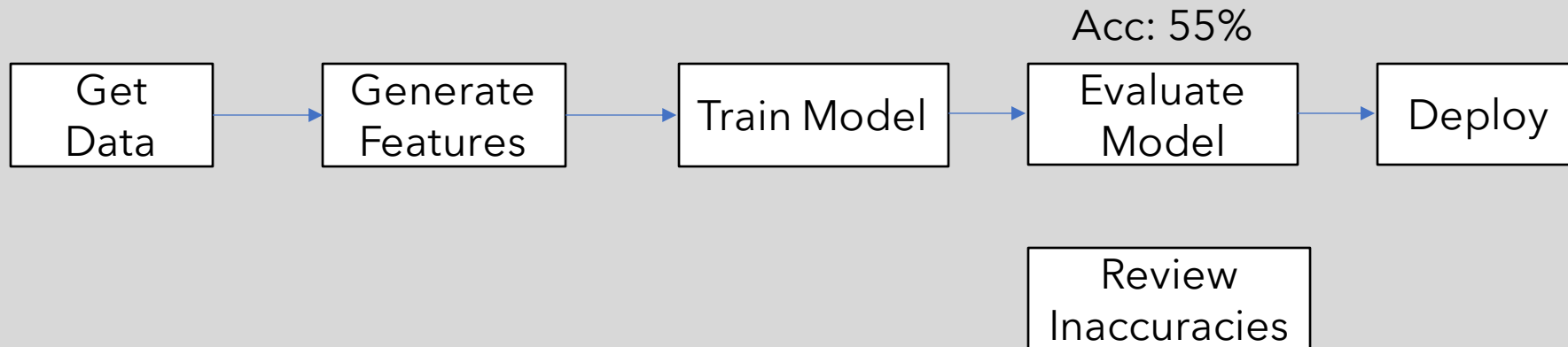
Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

```
Get Data → Generate Features → Train Model → Evaluate Model (Acc: 55%) → Deploy
                                                    ↓
                                            Review Inaccuracies
```

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.
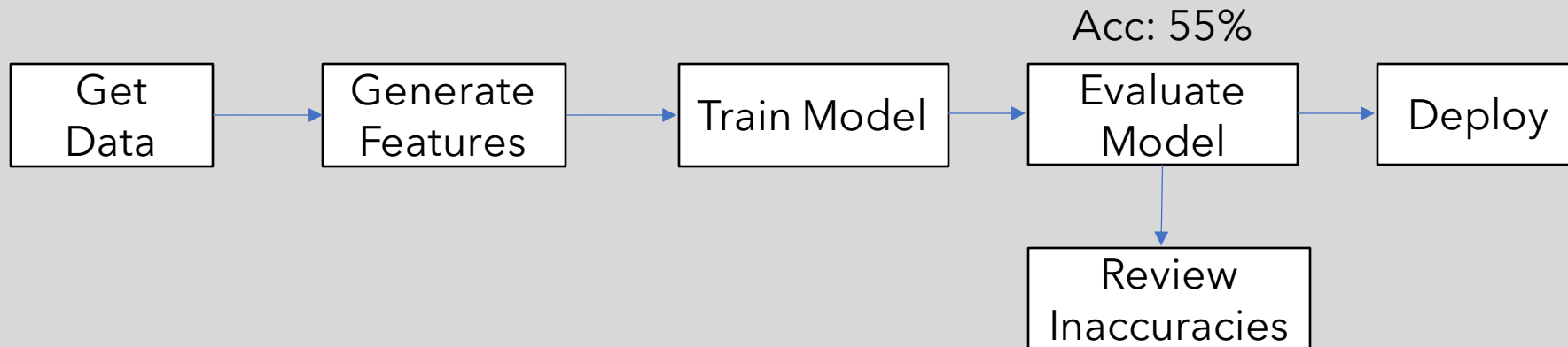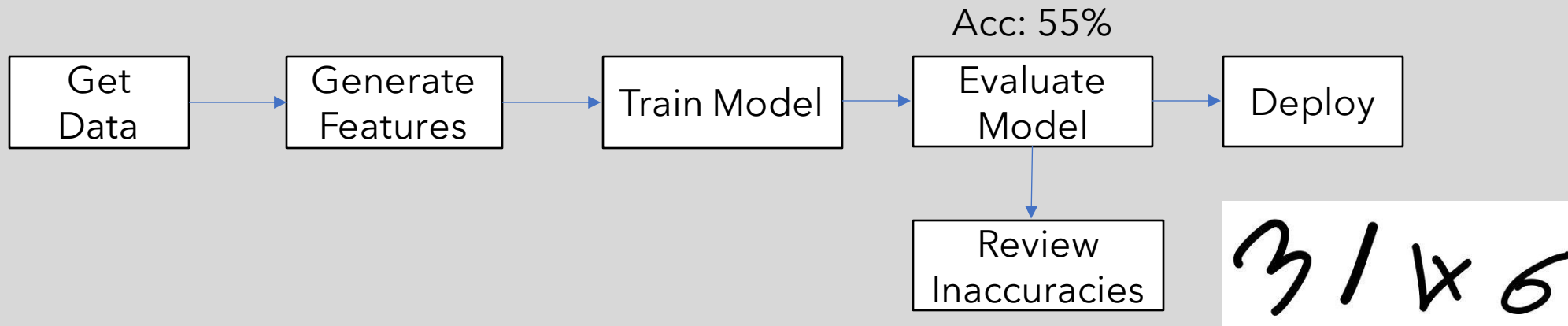
Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

Acc: 55%

```
┌──────────┐     ┌──────────┐     ┌──────────────┐     ┌──────────┐     ┌──────────┐
│   Get    │ ──> │ Generate │ ──> │ Train Model  │ ──> │ Evaluate │ ──> │  Deploy  │
│   Data   │     │ Features │     │              │     │  Model   │     │          │
└──────────┘     └──────────┘     └──────────────┘     └──────────┘     └──────────┘
                                                             │
                                                             v
                                                      ┌──────────────┐
                                                      │    Review    │
                                                      │ Inaccuracies │
                                                      └──────────────┘
```

31x6

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.
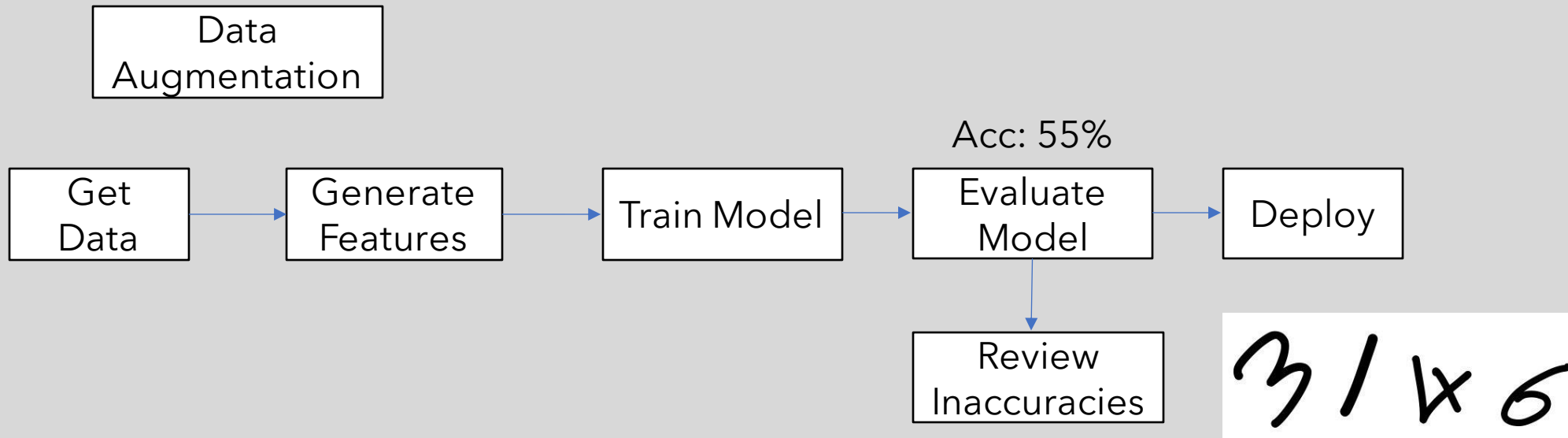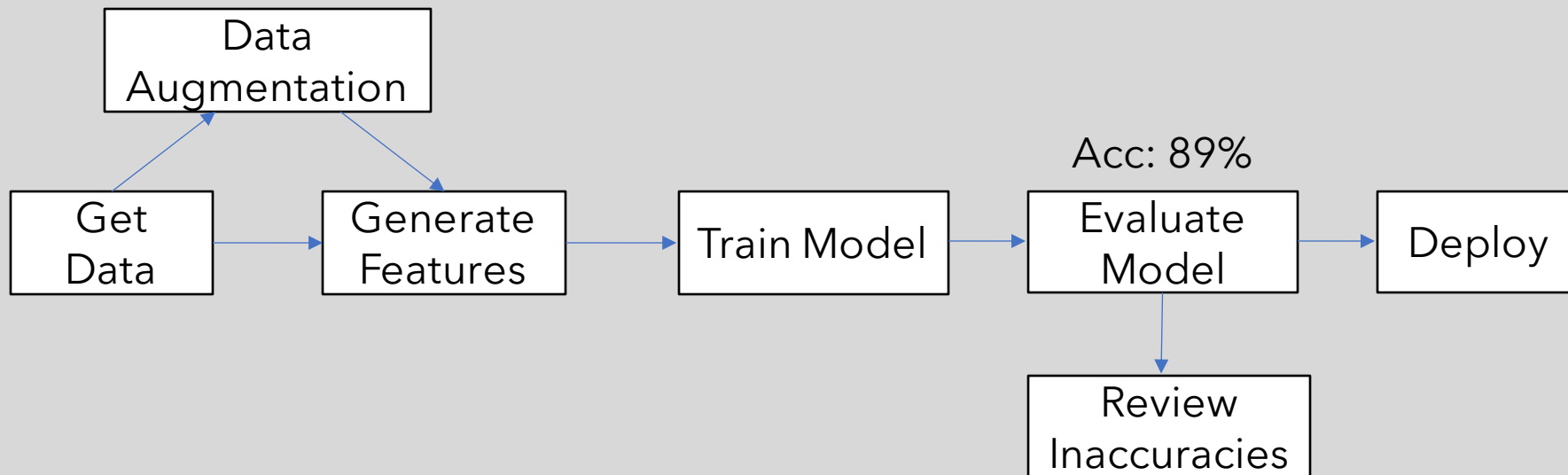
Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

Data Augmentation

Get Data → Generate Features → Train Model → Acc: 55% Evaluate Model → Deploy

Evaluate Model → Review Inaccuracies

3/x6

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.
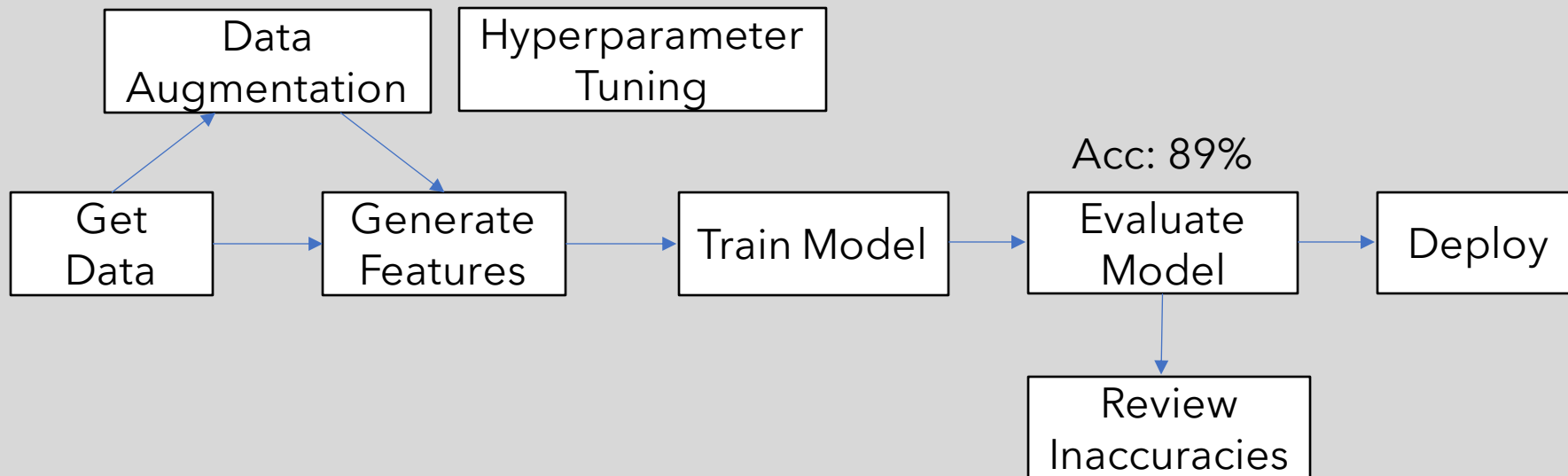
Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.

Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

# Pipelines

Typical machine learning workflows involve multiple steps, which can become complicated to keep track of when they are arranged in multiple scripts or notebooks.
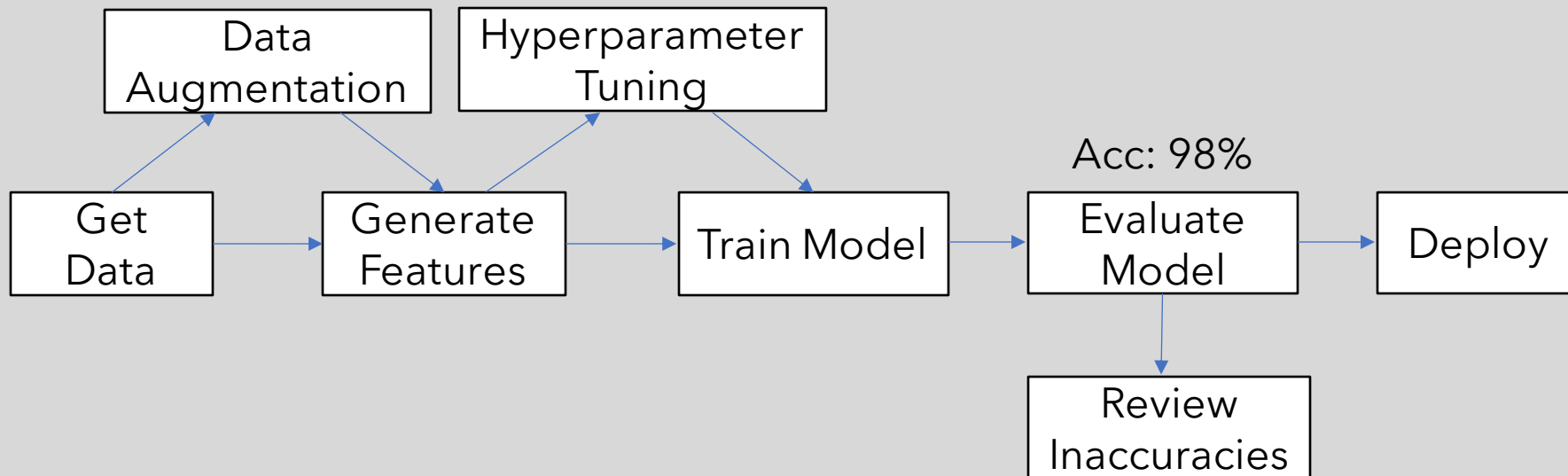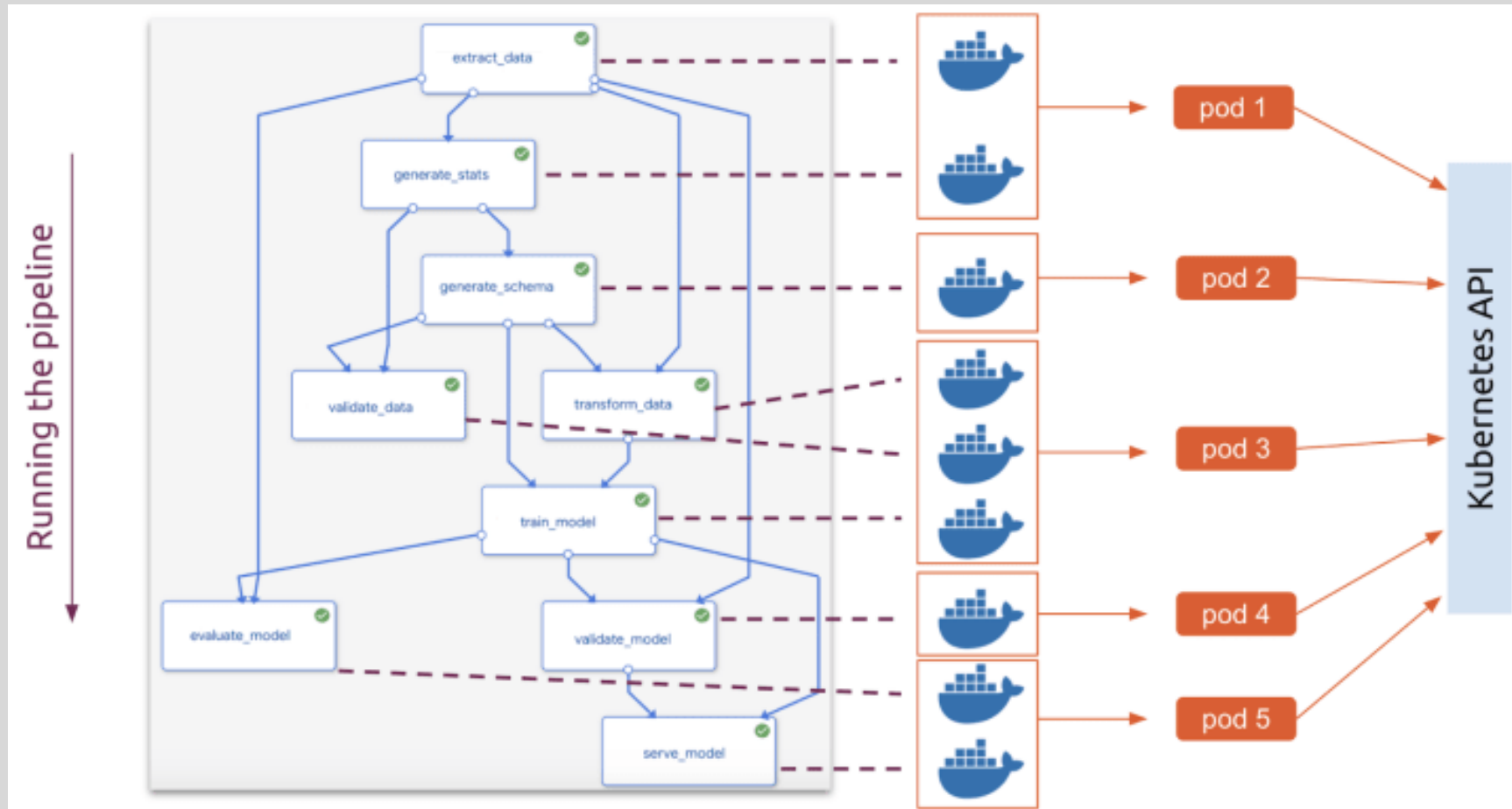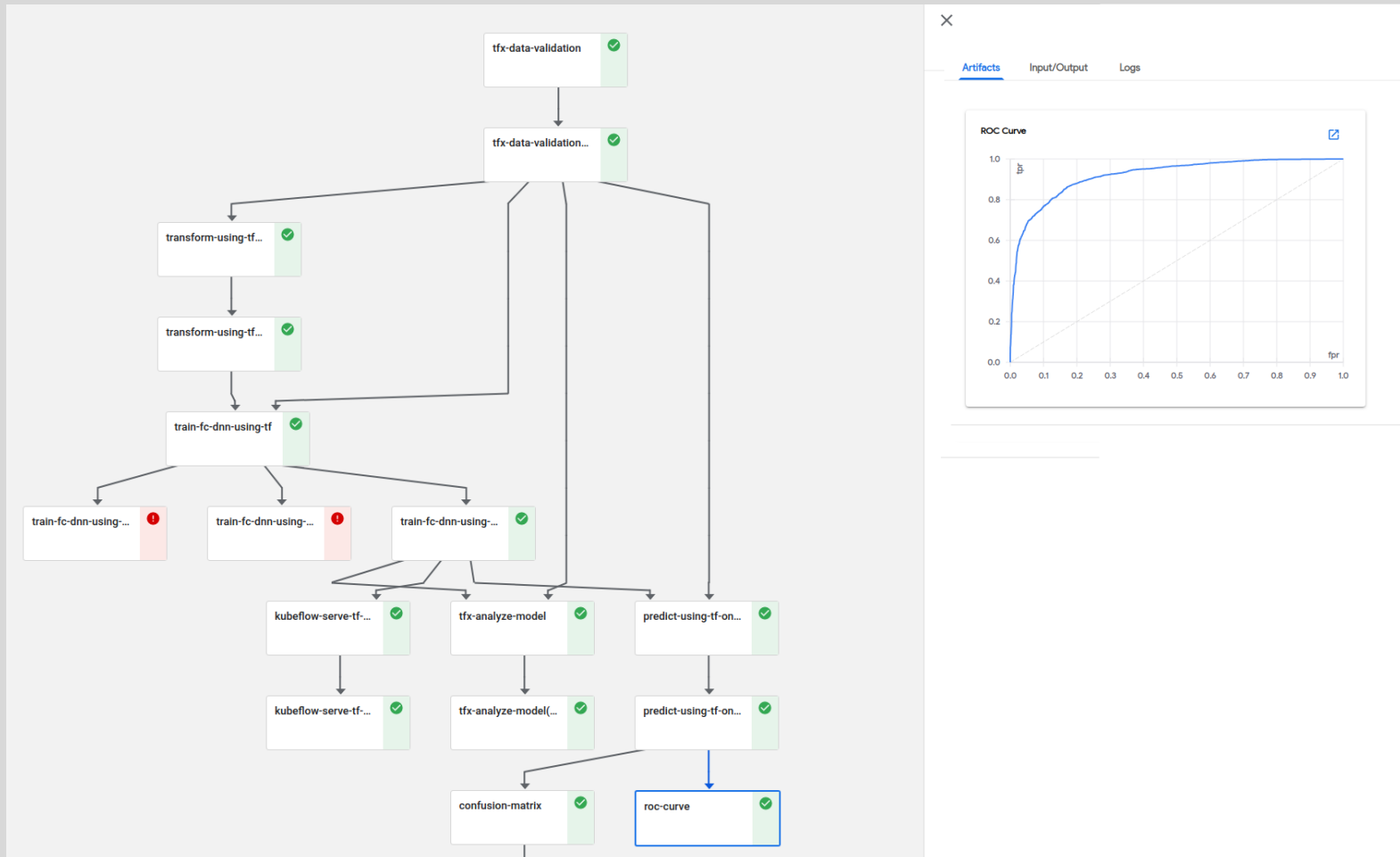
Kubeflow Pipelines allows developers to codify their ML workflows so that they are easily composable, shareable and reproducible.

| | |
|---|---|
| Data Augmentation | Hyperparameter Tuning |

Acc: 98%

| Get Data | Generate Features | Train Model | Evaluate Model | Deploy |
|---|---|---|---|---|

Review Inaccuracies

# Pipelines

# Pipelines

# Katib

Katib is Kubeflow's built in hyperparameter tuner.

Given that hyperparameter tuning locally is done through a big for loop, using Kubeflow can make the whole process more optimised.

```
spec:
  parallelTrialCount: 3
  maxTrialCount: 12
  maxFailedTrialCount: 3
  objective:
    type: maximize
    goal: 0.99
    objectiveMetricName: accuracy_1
  algorithm:
    algorithmName: random
```

How many at one time

Total number of tests

Stop after this many fails
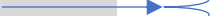
Metric to improve on and auto-stop criteria

Algorithm for choosing next value to test (grid search, Bayesian optimisation, …)

# Katib

Katib is Kubeflow's built in hyperparameter tuner.

Given that hyperparameter tuning locally is done through a big for loop, using Kubeflow can make the whole process more optimised.
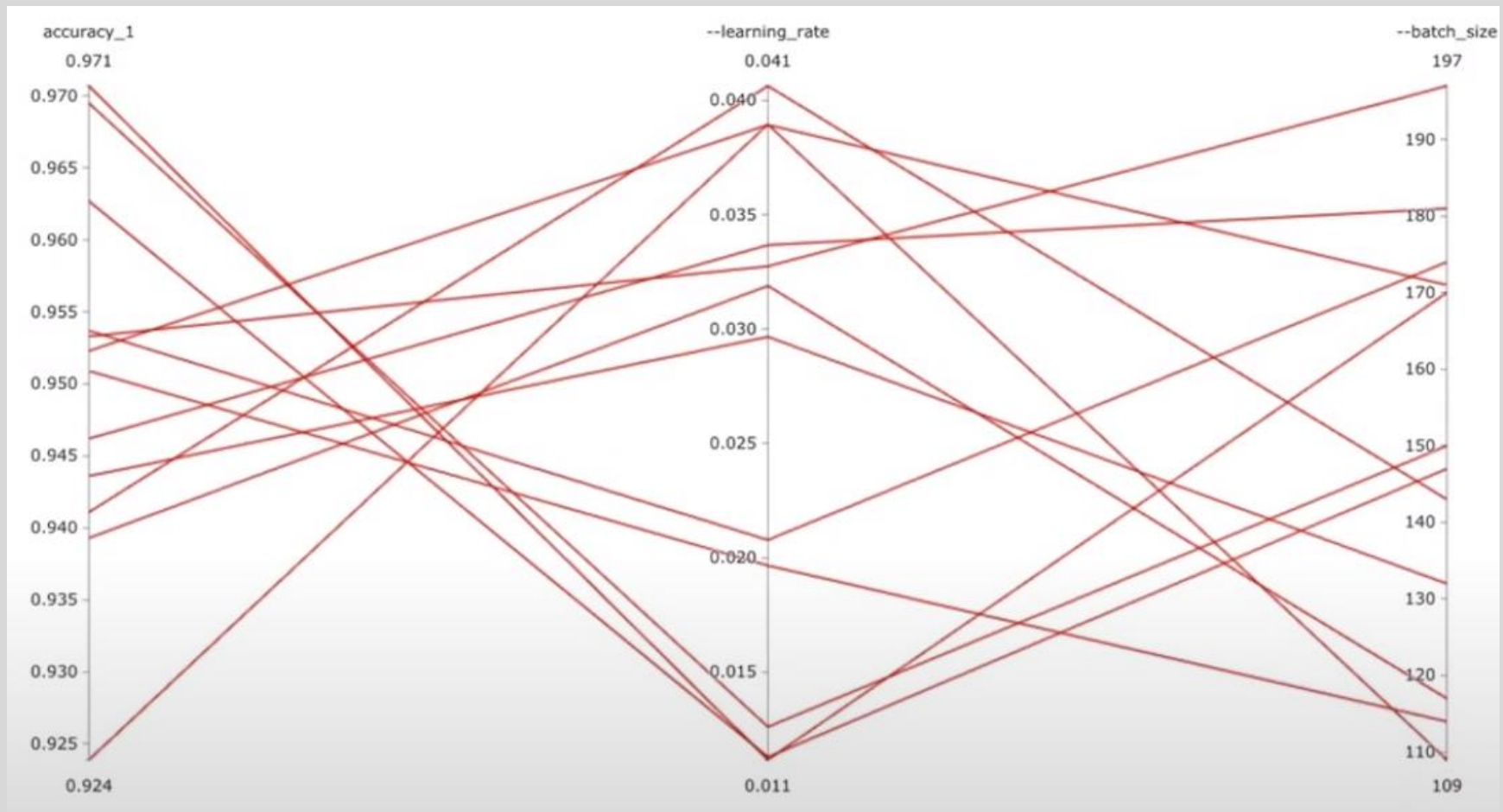
Choose hyperparameter and range to search

```
parameters:
    - name: --learning_rate
      parameterType: double
      feasibleSpace:
        min: "0.01"
        max: "0.05"

    - name: --batch_size
      parameterType: int
      feasibleSpace:
        min: "100"
        max: "200"
```

# Katib

# Extras for your Pipelines

○ There are a number of external tools that can be added to your Kubeflow pipeline to improve your ML models.

| Method | Models | Explanations | Classification | Regression | Tabular | Text | Images |
|---|---|---|---|---|---|---|---|
| ALE | BB | global | ✔ | ✔ | ✔ | | |
| Anchors | BB | local | ✔ | | ✔ | ✔ | ✔ |
| CEM | BB* TF/Keras | local | ✔ | | ✔ | | ✔ |
| Counterfactuals | BB* TF/Keras | local | ✔ | | ✔ | | ✔ |
| Prototype Counterfactuals | BB* TF/Keras | local | ✔ | | ✔ | | ✔ |
| Integrated Gradients | TF/Keras | local | ✔ | ✔ | ✔ | ✔ | ✔ |
| Kernel SHAP | BB | local global | ✔ | ✔ | ✔ | | |
| Tree SHAP | WB | local global | ✔ | ✔ | ✔ | | |

## Model Confidence

These algorithms provide **instance-specific** scores measuring the model confidence for making a particular prediction.

| Method | Models | Classification | Regression | Tabular | Text | Images | Categorical Features | Train set required |
|---|---|---|---|---|---|---|---|---|
| Trust Scores | BB | ✔ | | ✔ | ✔ (1) | ✔ (2) | | Yes |
| Linearity Measure | BB | ✔ | ✔ | ✔ | | ✔ | | Optional |

https://github.com/seldonio/alibi

# Extras for your Pipelines

| Method | Models | Explanations | Classification | Regression | Tabular | Text | Images |
|--------|--------|--------------|----------------|------------|---------|------|--------|
| Anchors | BB | local | ✓ | | ✓ | ✓ | ✓ |



(a) Original image    (b) Anchor for "beagle"    (c) Images where Inception predicts $P(\text{beagle}) > 90\%$

https://github.com/seldonio/alibi

Time to open up Kubeflow