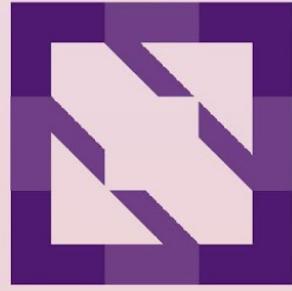




KubeCon

— North America 2023 —



CloudNativeCon





KubeCon



CloudNativeCon

North America 2023

# The Future of Interactive Data Science at Scale with Jupyter and Kubeflow

*Andrey Velichkevich - Apple  
Zachary Sailer - Apple*

# Who are we?



KubeCon



CloudNativeCon

North America 2023



Physics • Jupyter • Software



Science • Kubeflow • MLOps

# Let us clearly define what we mean...

The future of  
interactive data science  
at scale

# How (data) science is done in practice

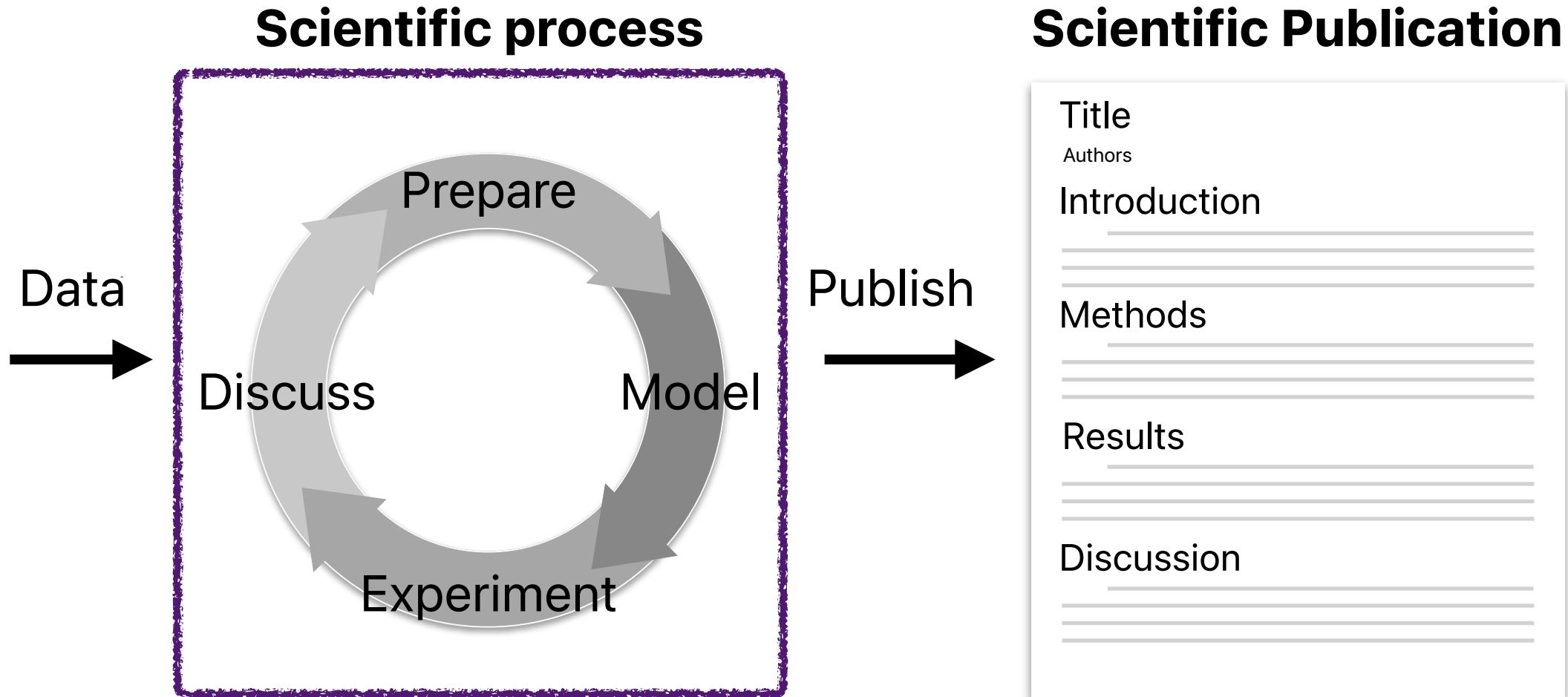


KubeCon



CloudNativeCon

North America 2023

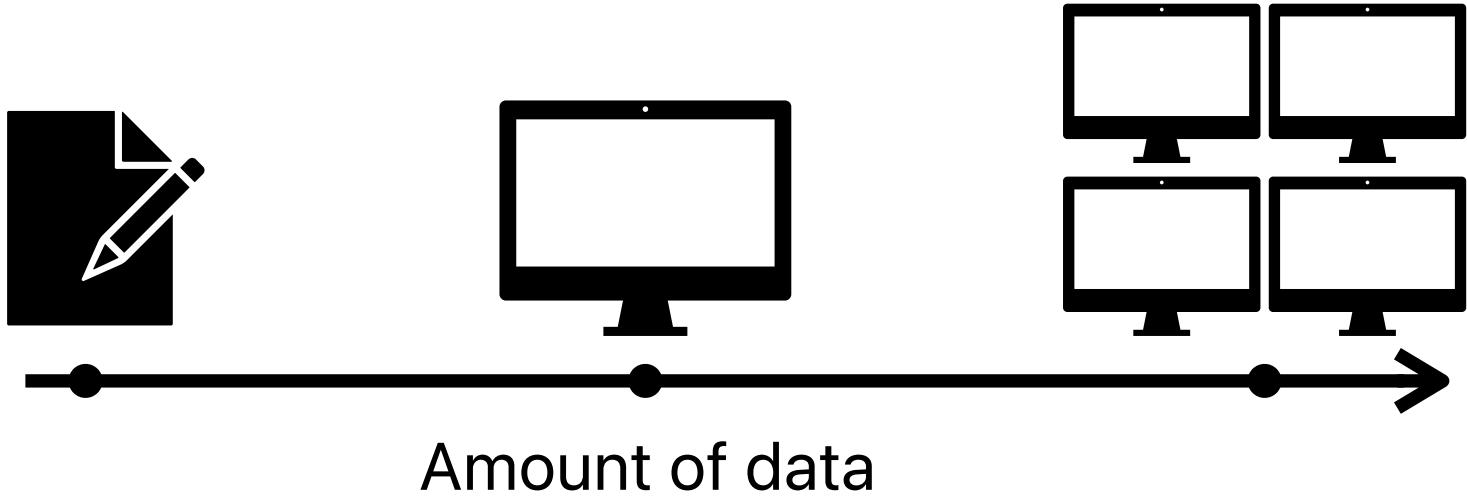


**Science was done with pencil and paper.**

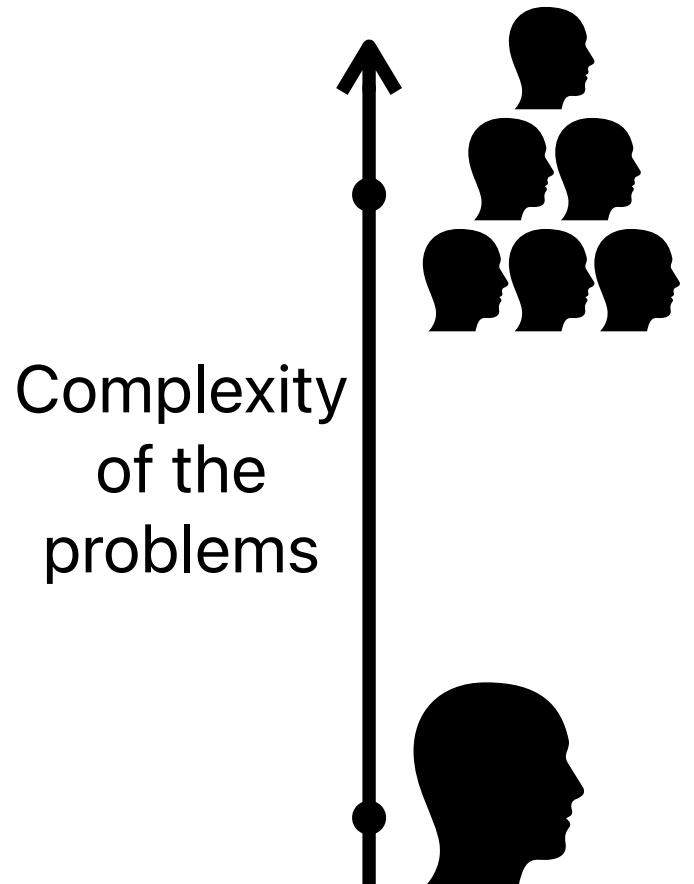
1923

**Data is big and complex.**

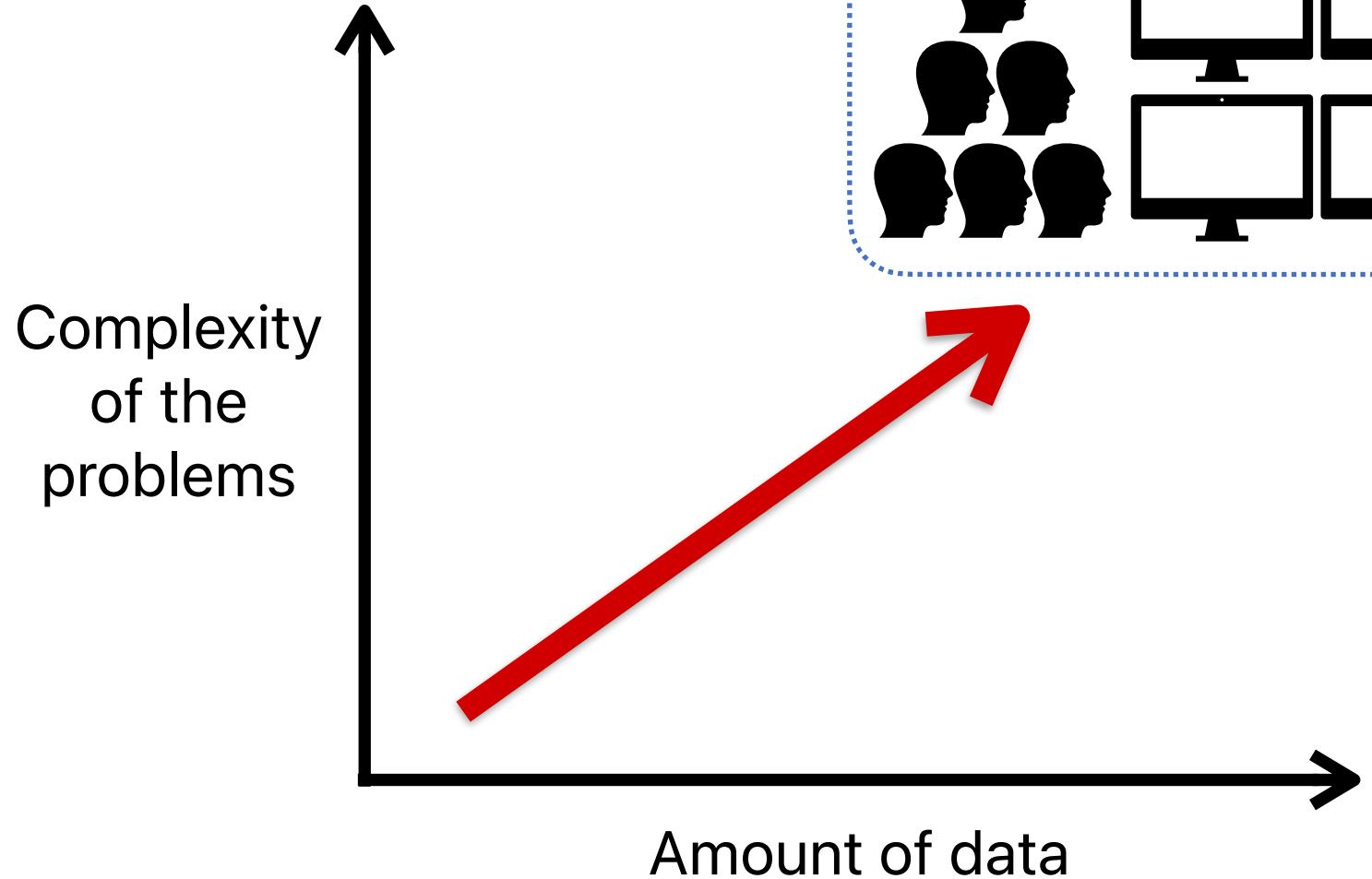
2023



2023

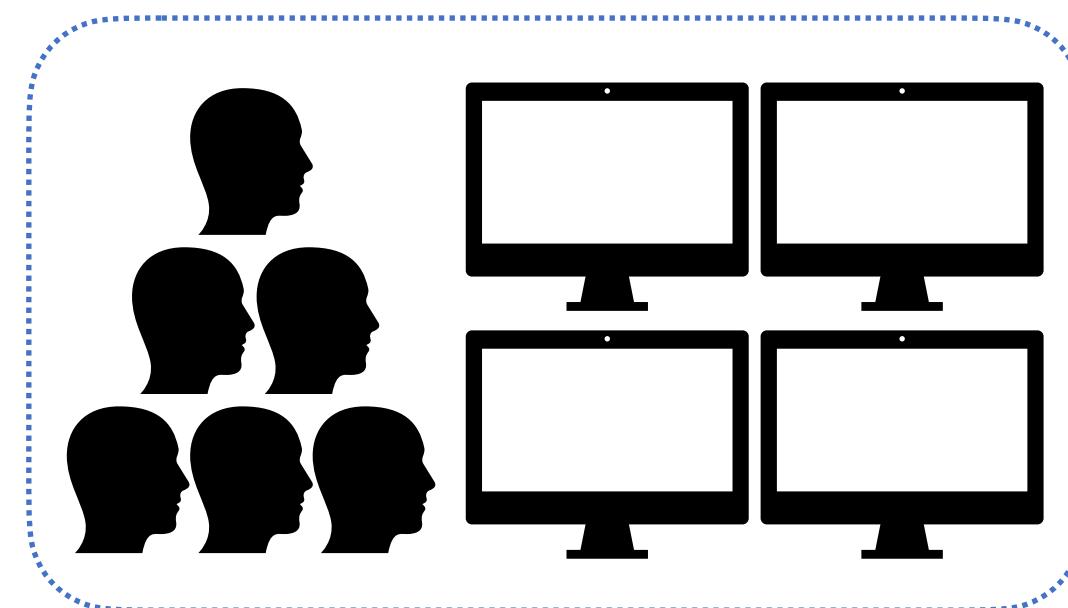


2023



2023

We need tools that promote  
*humans-in-the-loop* data science



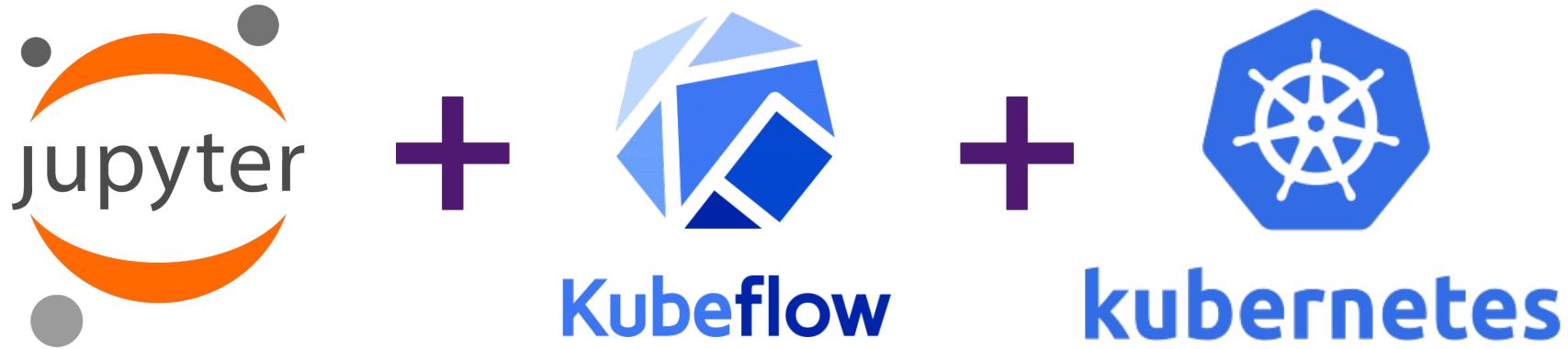
2023

# Our requirements for *interactive Data Science*

Computational  
Exploratory  
Collaborative  
Publication ready  
Built for AI/ML  
Simple to scale

# Demo

# What you just saw



# Why Jupyter?

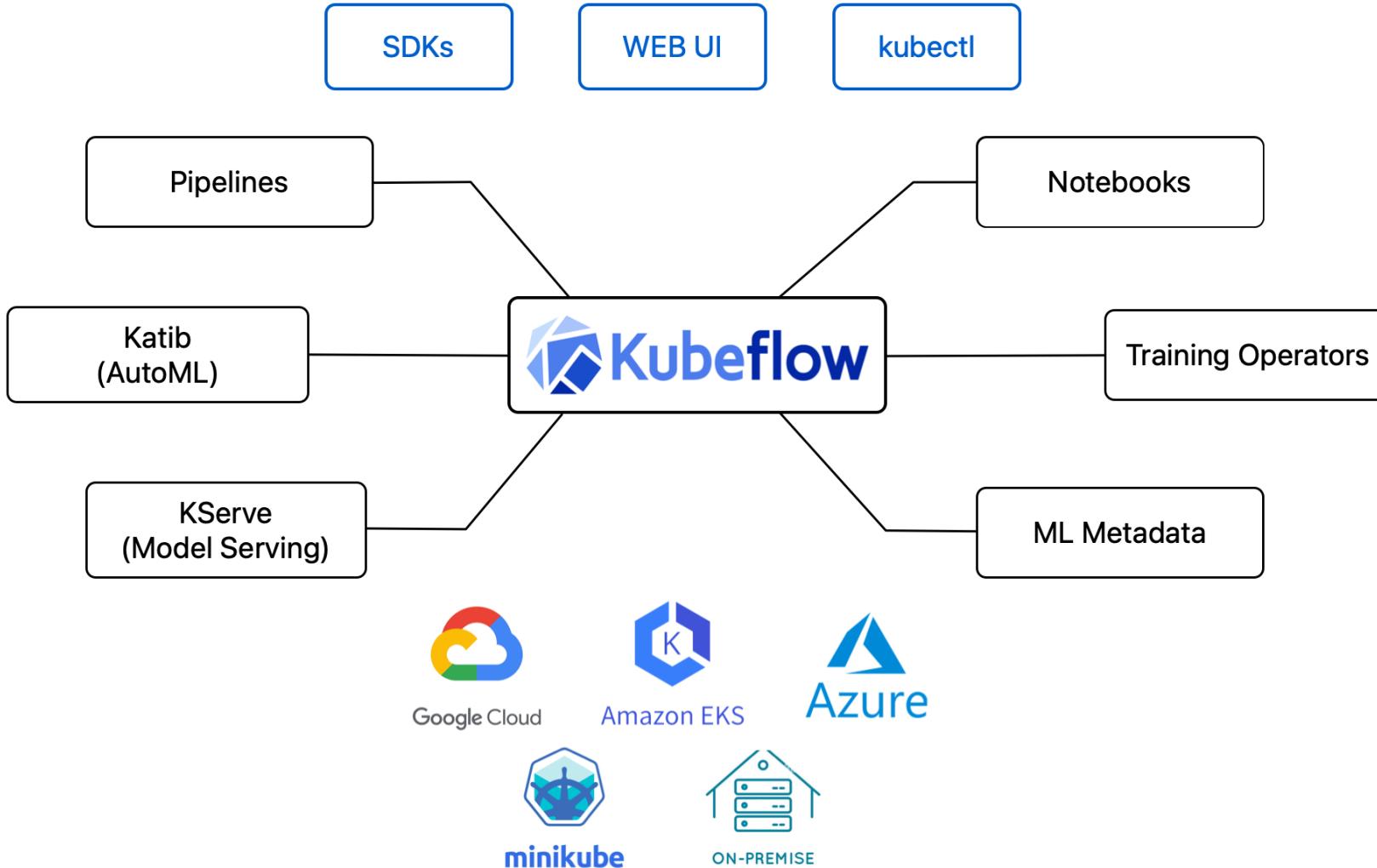
The screenshot shows a Jupyter Notebook interface with the following details:

- Title:** bayesianStatistics.ipynb
- Kernel:** Python3.8
- Code Cell 19:** Generates simulated data for Bayesian linear regression. It imports numpy and matplotlib.pyplot, defines true parameter values (alpha=1, sigma=1, beta=[1, 2.5]), creates predictor variables X1 and X2, and simulates the outcome variable Y. The code is executed at 2023-10-30 20:29:06.
- Text:** "Here is what the simulated data look like in matplotlib."
- Code Cell 25:** Plots the simulated data using matplotlib's subplots function. It creates two scatter plots: one for X1 vs Y and another for X2 vs Y. The plots are displayed below the code cell.

- ✓ Computational
- ✓ Exploratory
- ✓ Collaborative
- ✓ Publication ready

Built for AI/ML  
Simple to scale

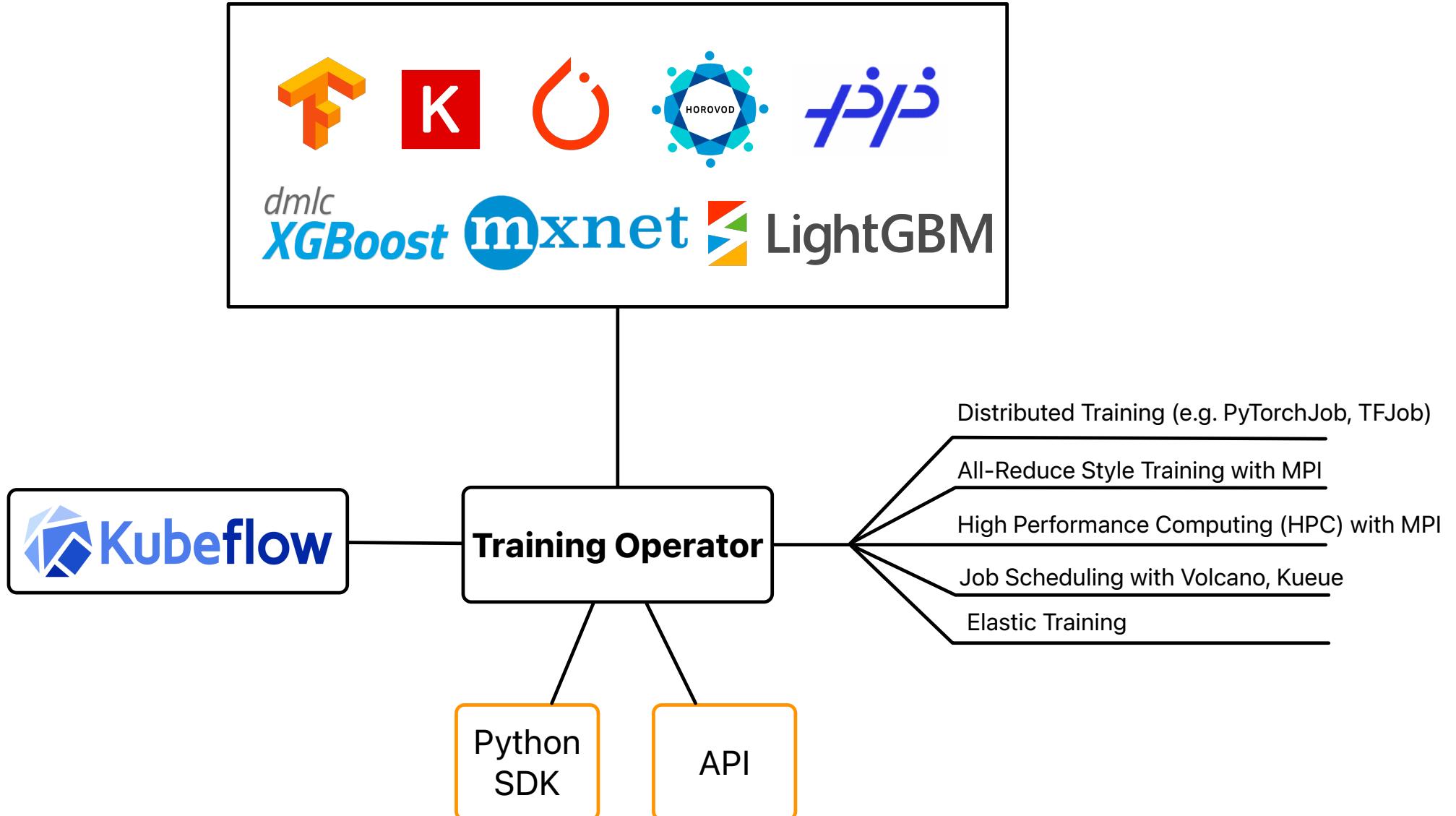
# Why Kubeflow ?



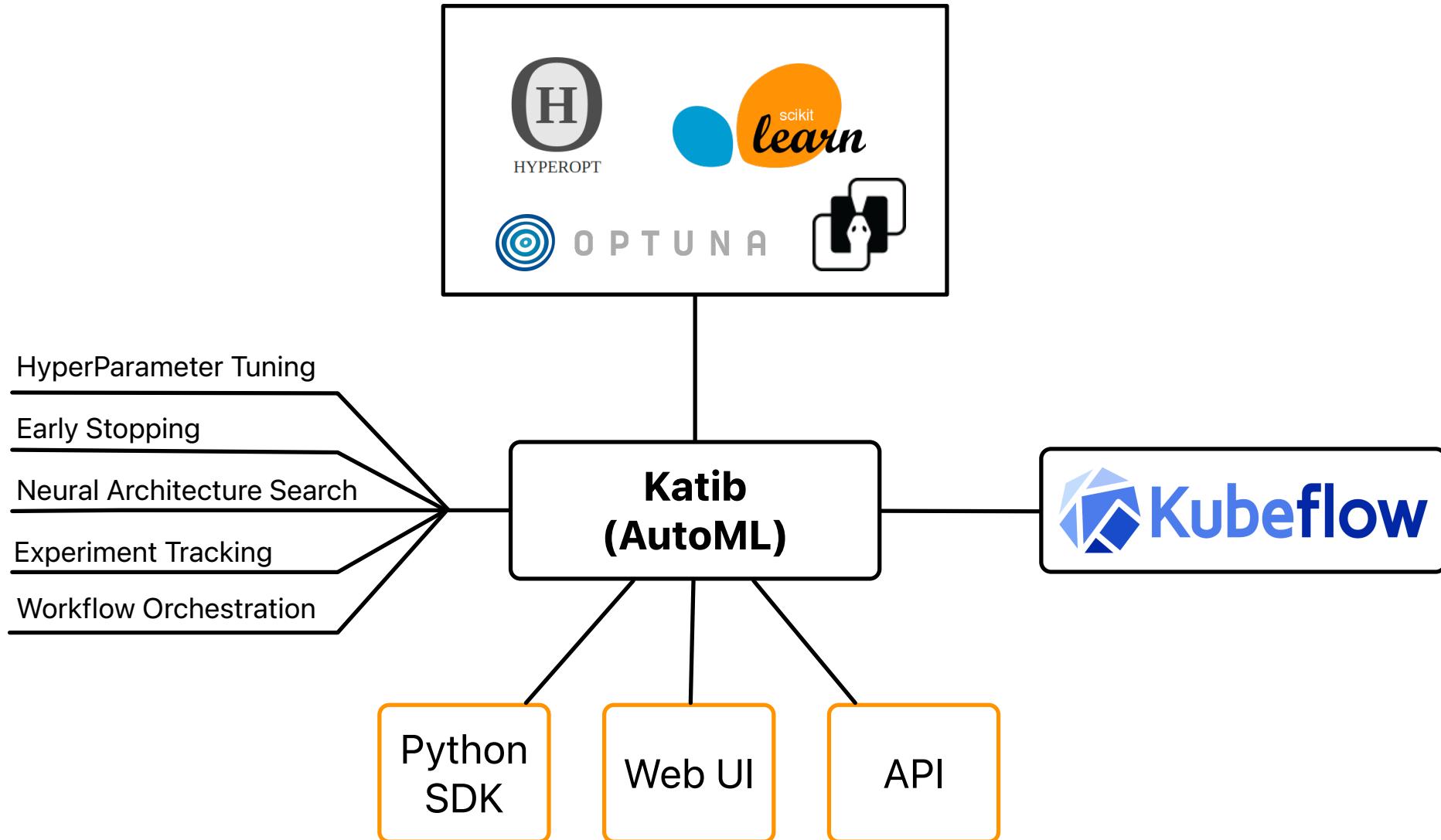
Computational  
Exploratory  
Collaborative  
Publication ready

- ✓ Built for AI/ML
- ✓ Simple to scale

# Kubeflow Overview - Training Operator



# Kubeflow Overview - Katib



# How does Kubeflow simplify scale?

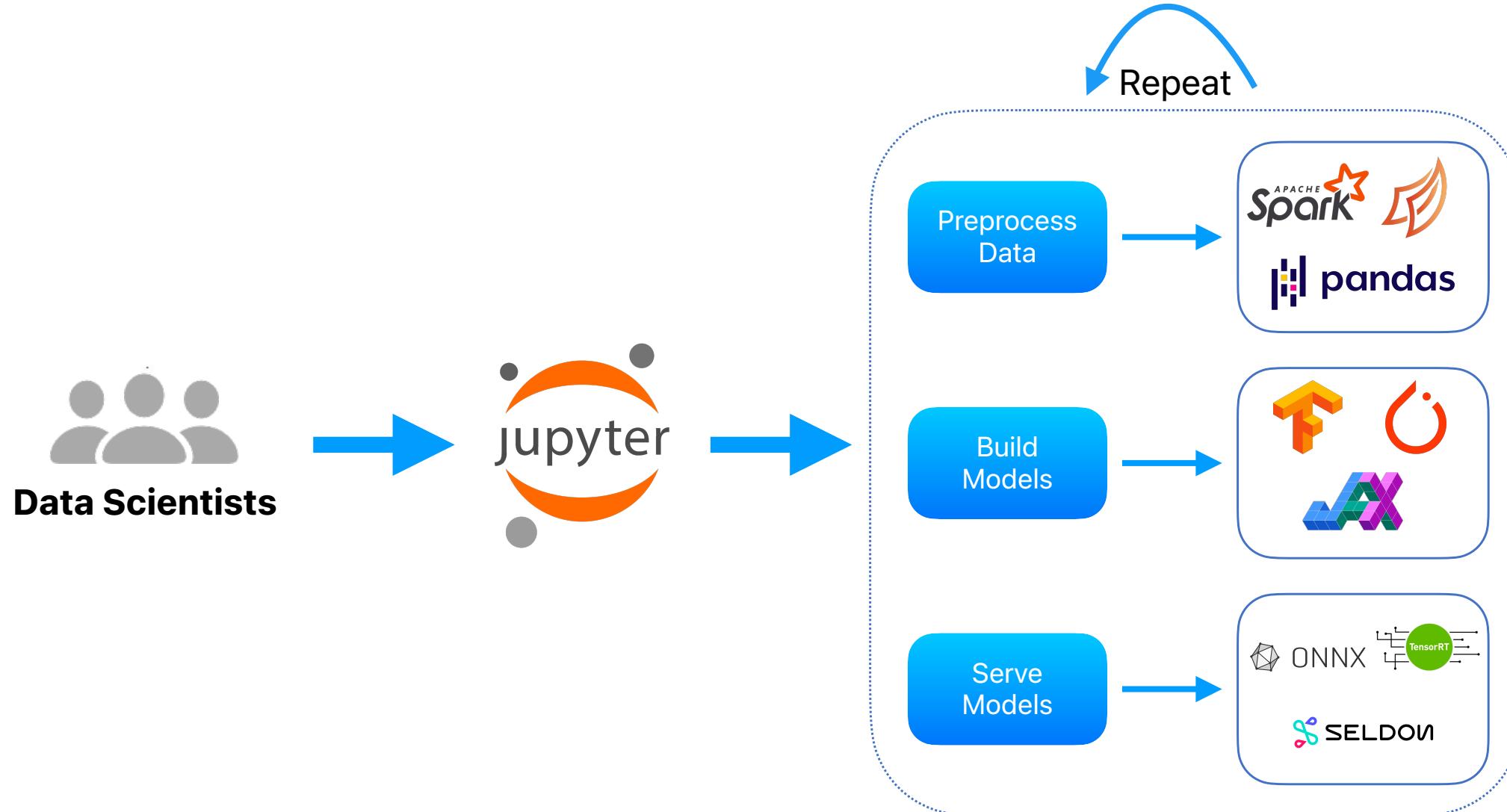
Computational  
Exploratory  
Collaborative  
Publication ready

- ✓ Built for AI/ML
- ✓ Simple to scale

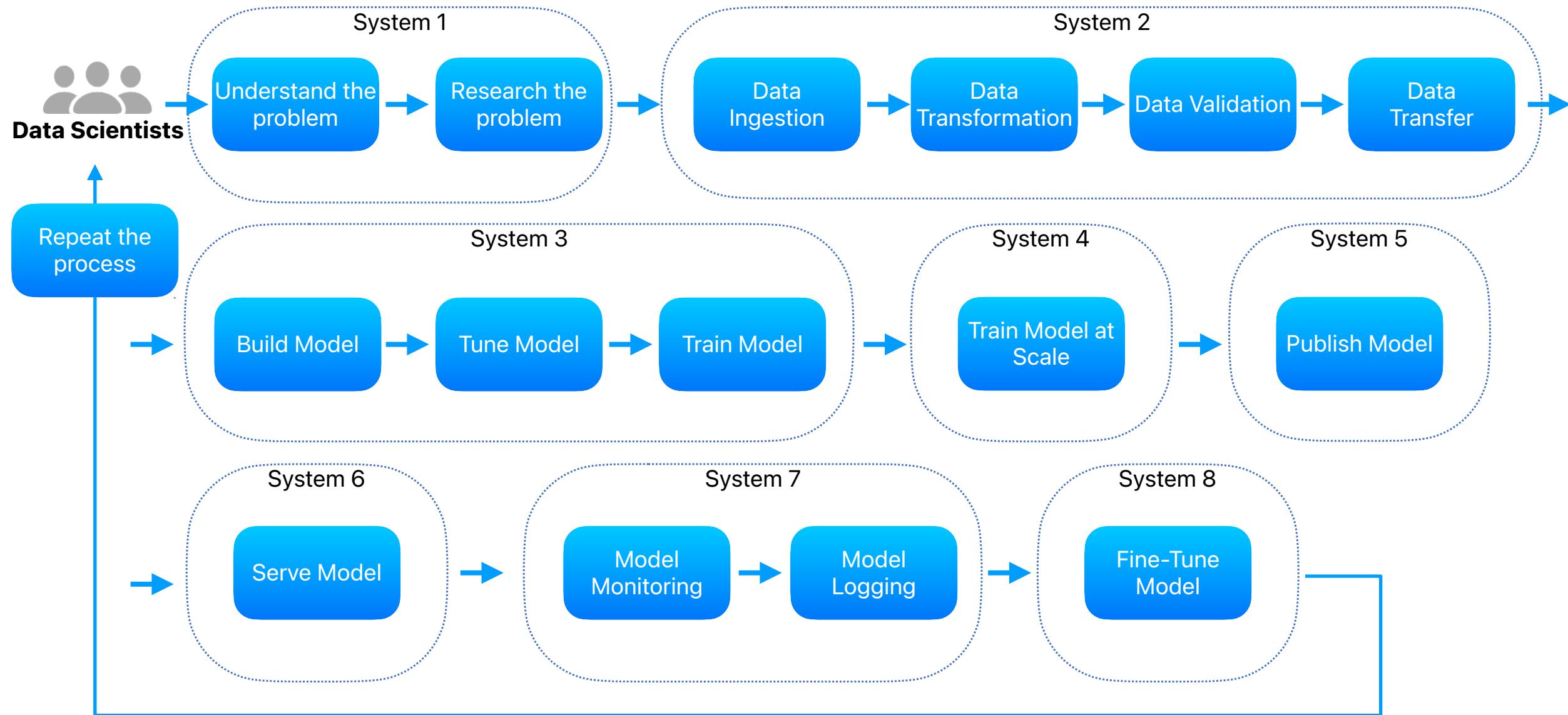


**Kubeflow**

# Lifecycle of ML - Desired

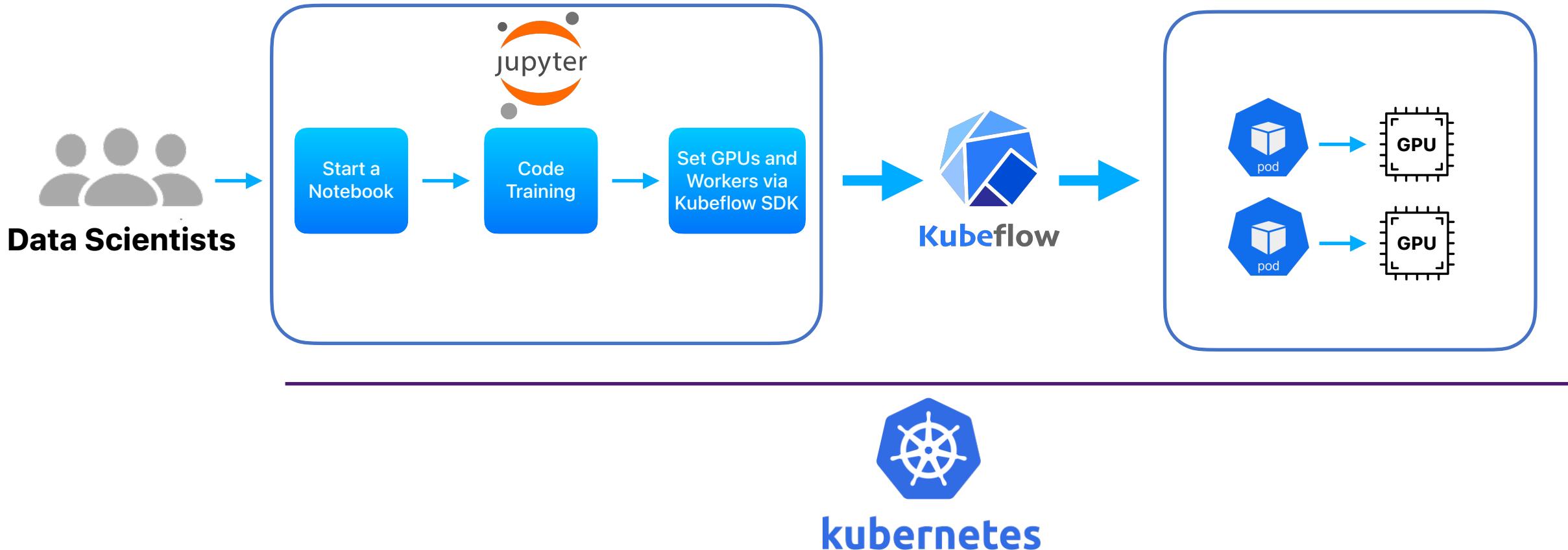


# Lifecycle of ML - Reality



Adding machines should be  
**effortless**

# Simple to Scale ML Experience



# Simple Experiment for Training

```
[5]: def train_model():
    import torch
    ...
    # Set dist strategy.
    torch.distributed.init_process_group(
        backend="nccl", rank=RANK, world_size=WORLD_SIZE
    )
    Distributor = nn.parallel.DistributedDataParallel
    model = Distributor(model)

    # Train model
    model.train()
```

```
[6]: from kubeflow.training import TrainingClient

job_name = "train-job"

TrainingClient().create_job(
    name=job_name,
    func=train_model,
    num_worker_replicas=100,
    resources_per_worker={"gpu": "1"},
)
```

```
[7]: TrainingClient().get_job_logs(name=job_name)
```

Last executed at 2023-01-18 12:52:45 in 63ms

```
2023-01-18T12:44:42Z INFO      Start training for RANK: 0. WORLD_SIZE: 4
2023-01-18T12:44:42Z INFO      Train Epoch: 0 [0/60000 (0%)] loss=2.3032
2023-01-18T12:44:42Z INFO      Reducer buckets have been rebuilt in this iteration.
2023-01-18T12:44:43Z INFO      Train Epoch: 0 [320/60000 (1%)] loss=2.2963
2023-01-18T12:44:44Z INFO      Train Epoch: 0 [640/60000 (1%)] loss=2.2893
2023-01-18T12:44:45Z INFO      Train Epoch: 0 [960/60000 (2%)] loss=2.2707
2023-01-18T12:44:46Z INFO      Train Epoch: 0 [1280/60000 (2%)] loss=2.2700
2023-01-18T12:44:47Z INFO      Train Epoch: 0 [1600/60000 (3%)] loss=2.2697
2023-01-18T12:44:48Z INFO      Train Epoch: 0 [1920/60000 (3%)] loss=2.2424
2023-01-18T12:44:49Z INFO      Train Epoch: 0 [2240/60000 (4%)] loss=2.1990
2023-01-18T12:44:50Z INFO      Train Epoch: 0 [2560/60000 (4%)] loss=2.1941
2023-01-18T12:44:51Z INFO      Train Epoch: 0 [2880/60000 (5%)] loss=2.1618
```

# Simple Experiment for Tuning

```
[1]: def train_model(parameters):
    import tensorflow

    ...
    model.fit(
        dataset,
        epochs=int(parameters["num_epochs"]),
        steps_per_epoch=10
    )

[2]: import kubeflow.katib as katib

job_name = "tune-job"

katib.KatibClient().tune(
    name=job_name,
    objective=train_model,
    parameters={
        "num_epochs": katib.search.int(min=6, max=7)
    },
    objective_metric_name="Accuracy",
    max_trial_count=100,
    parallel_trial_count=10,
    resources_per_trial={"gpu": "1"},
)
```

```
[3]: BEST_TRIAL = katib.KatibClient().get_optimal_hyperparameters(
    name=job_name
)

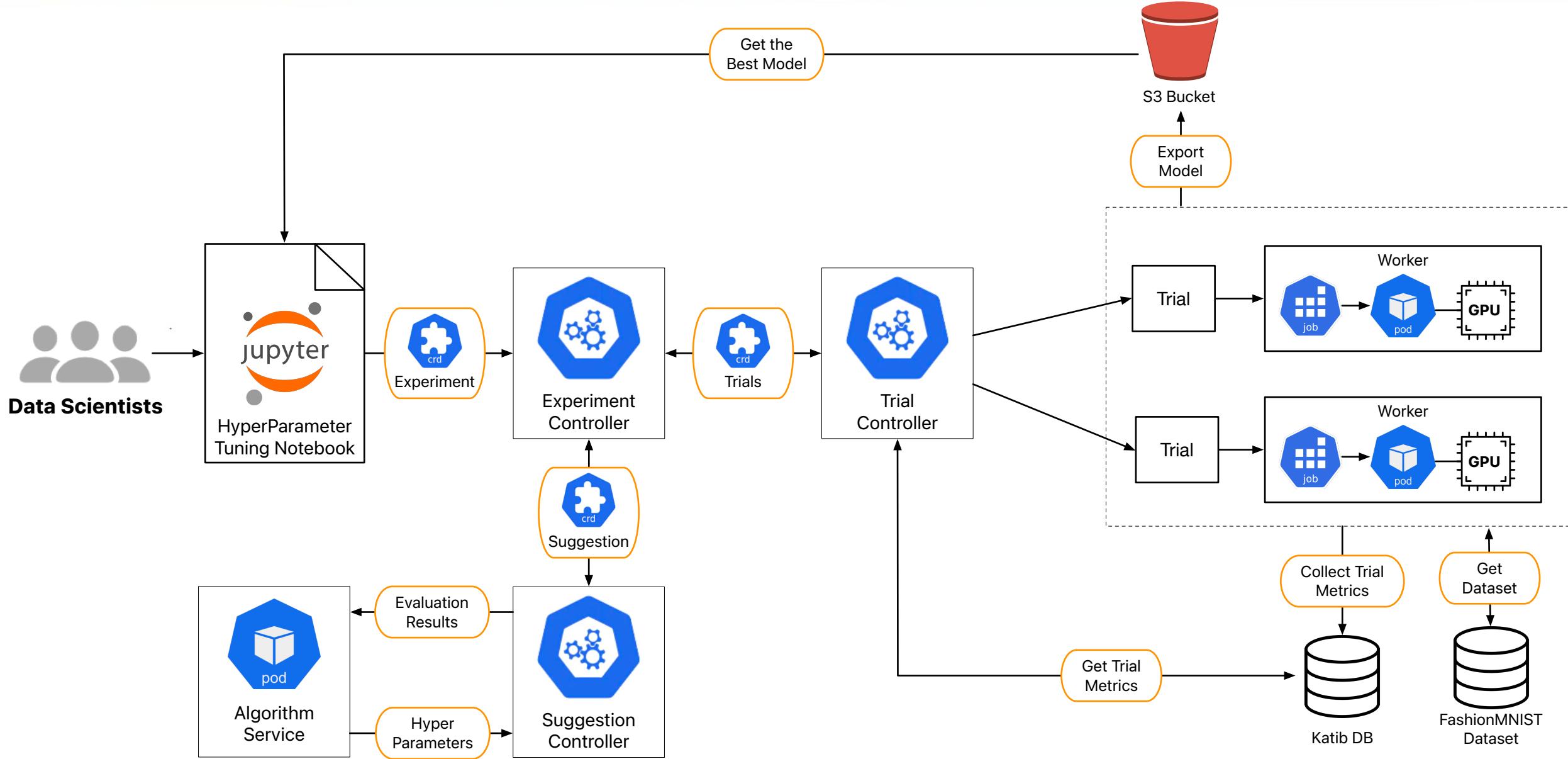
BEST_TRIAL.parameter_assignments[0]
Last executed at 2023-09-28 12:53:21 in 29ms

{'name': 'num_epochs', 'value': '7'}

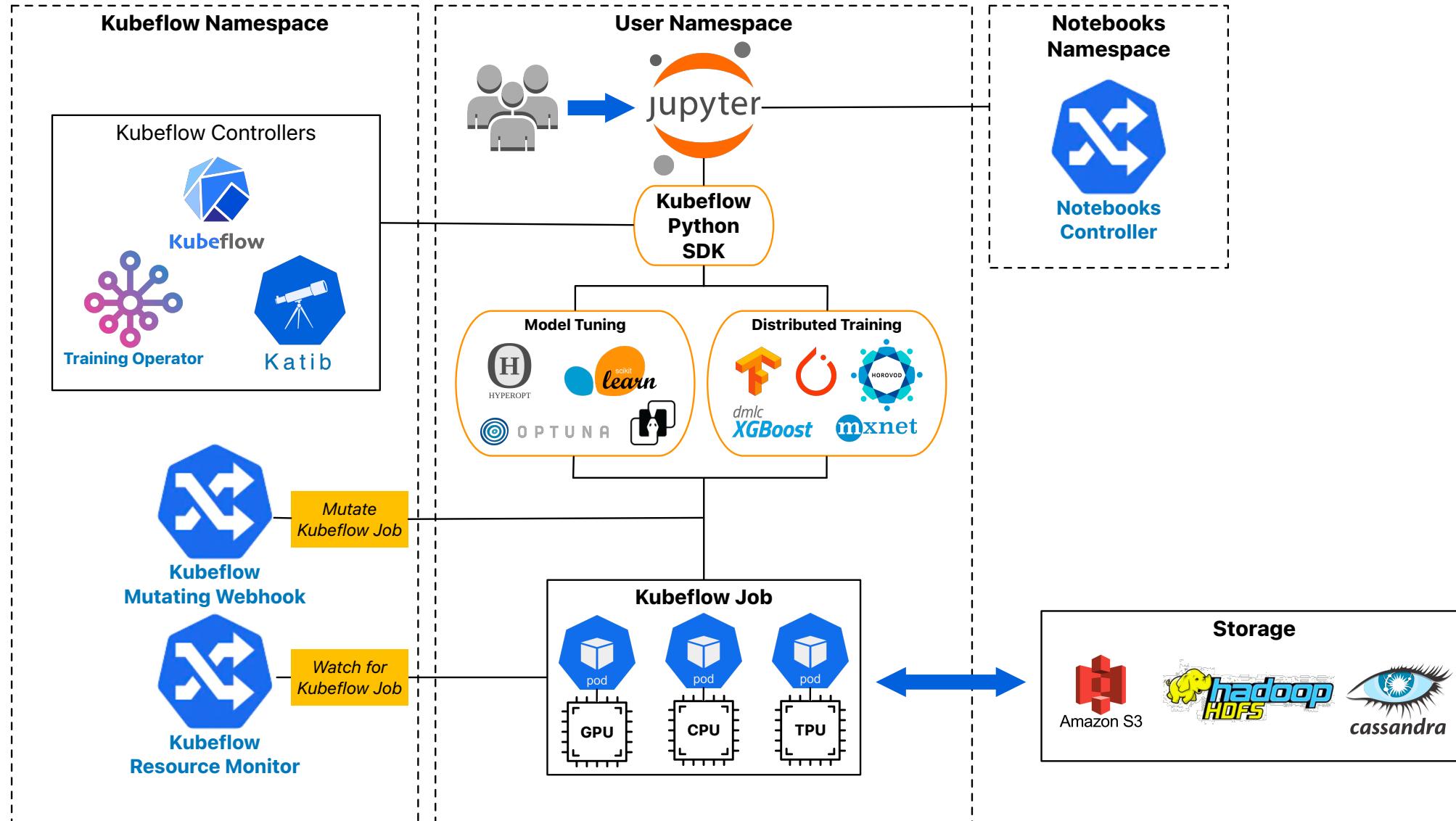
[4]: BEST_TRIAL.observation.metrics[0]
Last executed at 2023-09-28 12:53:21 in 3ms

{'latest': '0.8734',
 'max': '0.8734',
 'min': '0.6583333333333333',
 'name': 'Accuracy'}
```

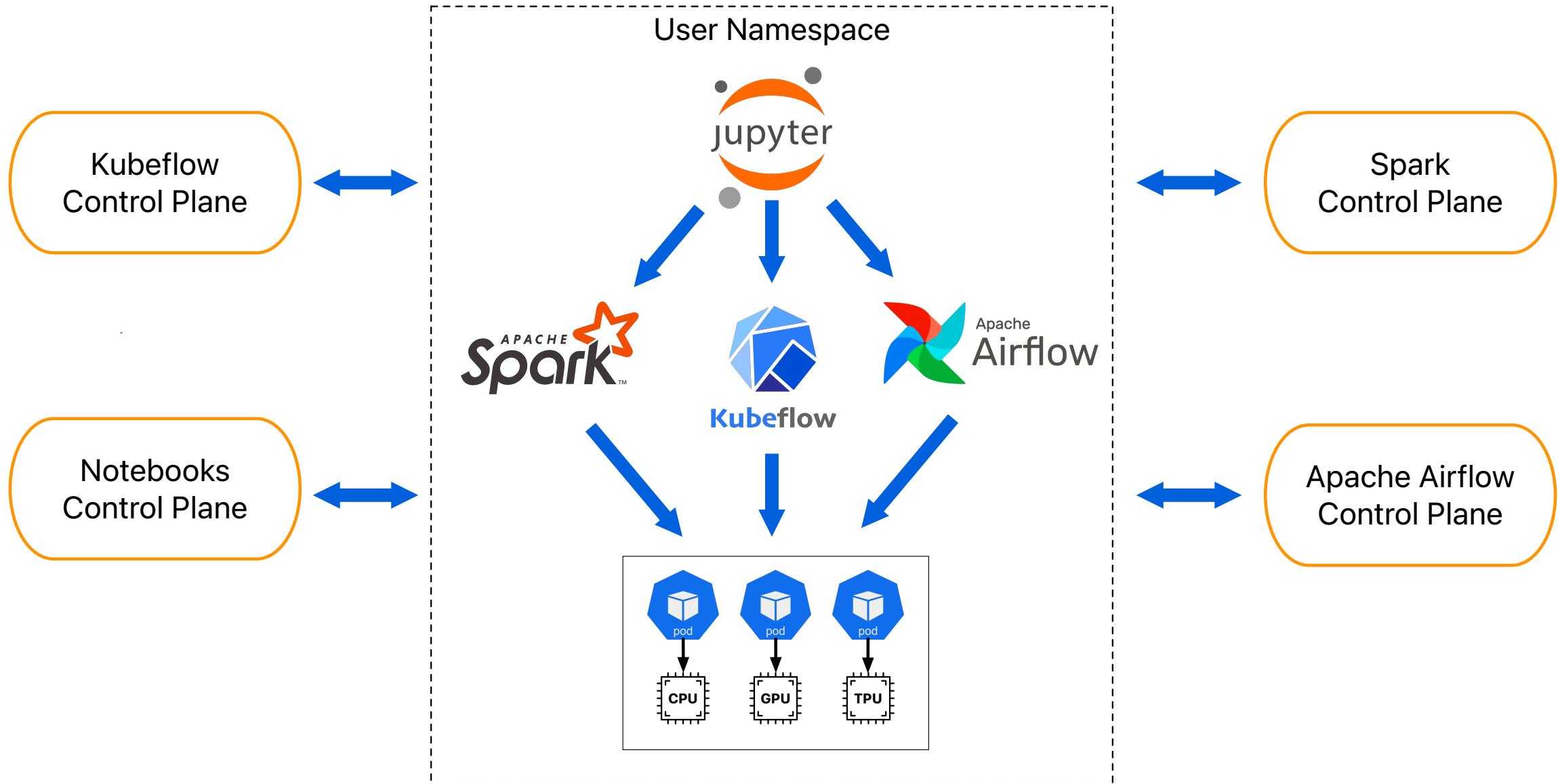
# The complexity you will never know about



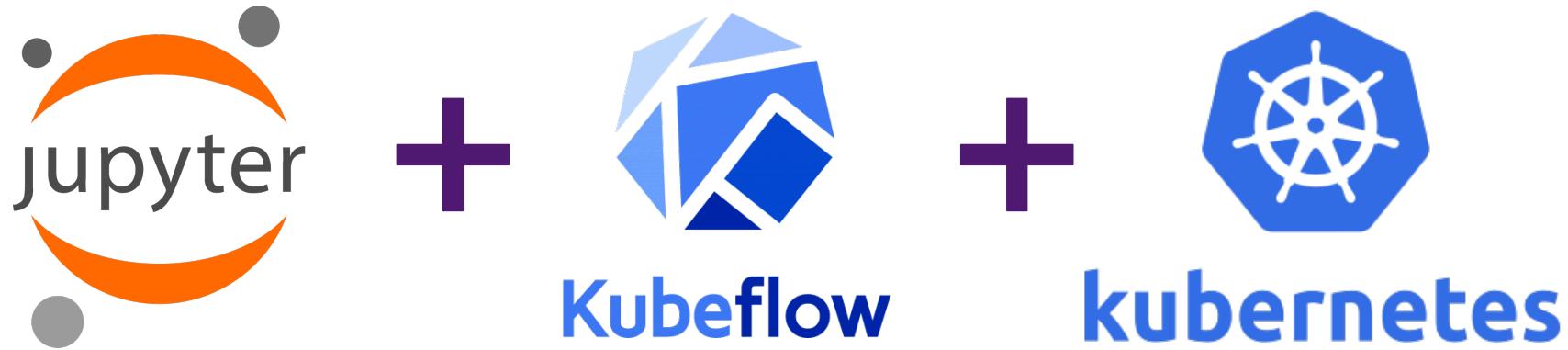
# Notebooks & Kubeflow Architecture



# Isolated User Environment



# Jupyter + Kubeflow

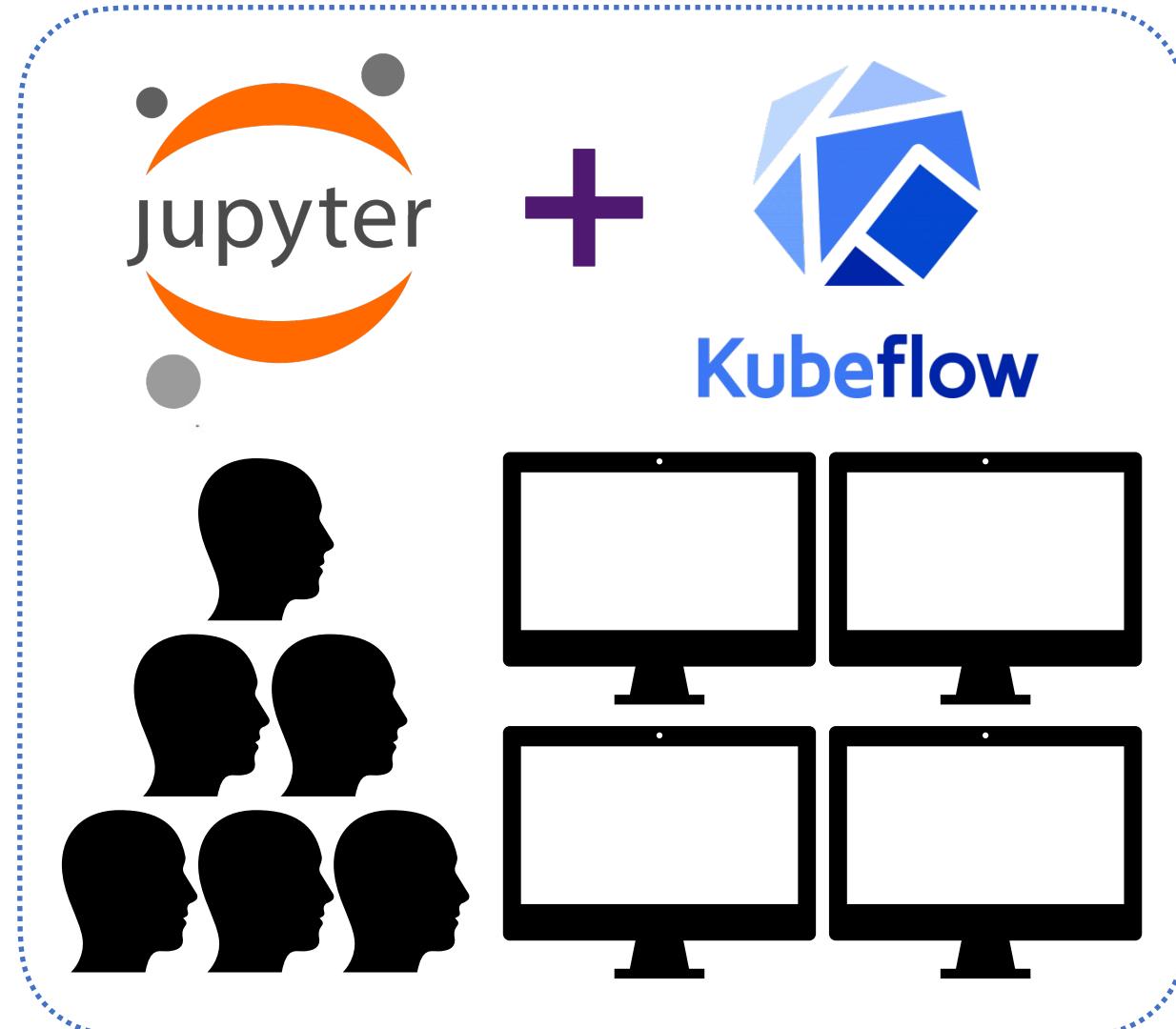


Interactive • Simple • Portable • Scalable

---

**ML Platform on Kubernetes**

# Interactive Data Science at Scale



- ✓ Computational
- ✓ Exploratory
- ✓ Collaborative
- ✓ Publication ready
- ✓ Built for AI/ML
- ✓ Simple to scale

# Jupyter and Kubeflow Communities

## Kubeflow

- Attend the [Kubeflow AutoML and Training community meetings](#)
- Join [Kubeflow Slack Channels](#)
- How to contribute ?
  - Check the Kubeflow [Training Operator](#) and [Katib](#) GitHub
  - Issues with the ["Help Wanted" Label](#)
  - Submit [a new Enhancement Proposal](#) to Kubeflow

## Jupyter

- Attend *any* public [Jupyter meeting](#)
- How to contribute ?
  - Join a ["Contributing Hour"](#)
  - [JupyterLab](#) and [Jupyter Collaboration](#) Projects

# Kubeflow & CNCF

**Kubeflow is now  
a CNCF  
incubating  
project!**

BY CNCF





**Please scan the QR Code above  
to leave feedback on this session**