



Cadi Ayyad University
FACULTY OF SCIENCES SEMLALIA

Orchestrating Continuous NAS with Early Stopping

A Comprehensive Tutorial on Kubeflow Katib, PVC Storage, and FastAPI

Prepared by: ELAKIL Hakima

Academic Year: 2025–2026

Contents

1	Project Context and Objectives	4
1.1	Core Objectives	4
2	Theoretical Framework	4
2.1	Used Technologies and Their Roles	4
2.2	The Essential Role of Experiments, Suggestions, and Trials	4
2.3	Methodological Framework: The "Student" Model	6
3	Project Structure	6
3.1	File Organization	6
3.2	Workflow	6
3.3	Config of the hyperparams	6
3.4	Tableau des Hyperparamètres et Intervalles pour Iris et MNIST	6
4	Phase 1: Prerequisites & Installation	6
4.1	MicroK8s Setup	7
4.2	Katib Installation (Standalone)	7
4.3	Namespace Configuration	7
4.4	Essential Verification	7
5	Phase 2: Training & Orchestration	8
5.1	Defining the Logic and Environment	8
5.2	Docker Environment Definition	9
5.3	Execution	9
5.4	Verification	9
6	Phase 3: Persistence and Serving	12
6.1	Model Recuperation: The "Bridge" Pod Strategy	12
6.1.1	Defining the Inspection Pod	12
6.1.2	Execution and Extraction	13
6.2	FastAPI Deployment and Testing	13
6.2.1	Testing the API	14
7	Conclusion	15
8	Part 2: MNIST Training and Trial Management	15
8.1	Introduction	15
8.2	Objectives	15
8.3	Python Script: <code>train_mnist.py</code>	15
8.4	Dockerfile for MNIST	17
8.5	Katib Experiment YAML: <code>mnist-nas.yaml</code>	17
8.6	Execution and Trial Tracking	18
8.7	Conclusion	19

9	MNIST Batch Inference Using Katib	20
9.1	Model Preparation	20
9.2	Bridge Pod for Model Access	20
9.3	Batch Inference Job	20
9.4	Inference Script	21
9.5	Execution	21
9.6	Visualization	22
9.7	Notes	22
10	General Conclusion	22
10.1	Key Competencies Acquired	23
11	References	23

1 Project Context and Objectives

This laboratory explores the implementation of an automated Machine Learning pipeline. The transition from static model development to a dynamic MLOps lifecycle is achieved by integrating automated search, resource management, and containerized serving.

1.1 Core Objectives

- **Automated Optimization:** Using Katib to explore the search space of a Multi-Layer Perceptron (MLP) without manual intervention.
- **Resource Efficiency:** Minimizing compute waste via **Early Stopping** based on real-time validation metrics.
- **Data Persistence:** Implementing **Persistent Volume Claims (PVC)** to ensure model artifacts survive the ephemeral nature of Kubernetes pods.
- **Inference Readiness:** Deploying the optimized model through a **FastAPI** REST interface.

2 Theoretical Framework

2.1 Used Technologies and Their Roles

- **MicroK8s:** The lightweight Kubernetes orchestration engine that manages the cluster, networking, and storage.
- **Docker:** Containerization technology ensuring the training environment (Python, Scikit-Learn) is reproducible and portable.
- **Kubeflow Katib:** The AutoML component responsible for managing Experiments, generating Trials, and suggesting Hyperparameters.
- **Jobs and Trials:** Each training run is encapsulated in a Kubernetes Job launched by Katib as a Trial. Trials can run in parallel, and Katib automatically tracks metrics (e.g., accuracy) to determine the best performing model.
- **Persistent Volume Claim (PVC):** A storage abstraction that acts as a "Virtual Hard Drive," ensuring data persists even after training pods are deleted.
- **FastAPI:** A modern, high-performance web framework used to serve the trained model as a REST API.

2.2 The Essential Role of Experiments, Suggestions, and Trials

In the context of Neural Architecture Search (NAS), the concepts of **Experiments**, **Suggestions**, and **Trials** are fundamental:

- **The Experiment:** This is the high-level control unit. It defines the *Objective* (e.g., maximize accuracy), the *Search Space* (e.g., number of layers, learning rate), and the *Search Algorithm* (e.g., Random Search, Bayesian Optimization).
- **The Suggestion:** This represents the decision-making layer. Based on the selected search algorithm, Katib generates suggestions that propose new hyperparameter configurations to explore, using feedback from previous trials to guide the search process.
- **The Trial:** This is the actual execution unit. Each Trial corresponds to one specific configuration of hyperparameters proposed by a Suggestion. Trials run in parallel as Kubernetes pods, enabling efficient exploration of multiple architectures simultaneously.

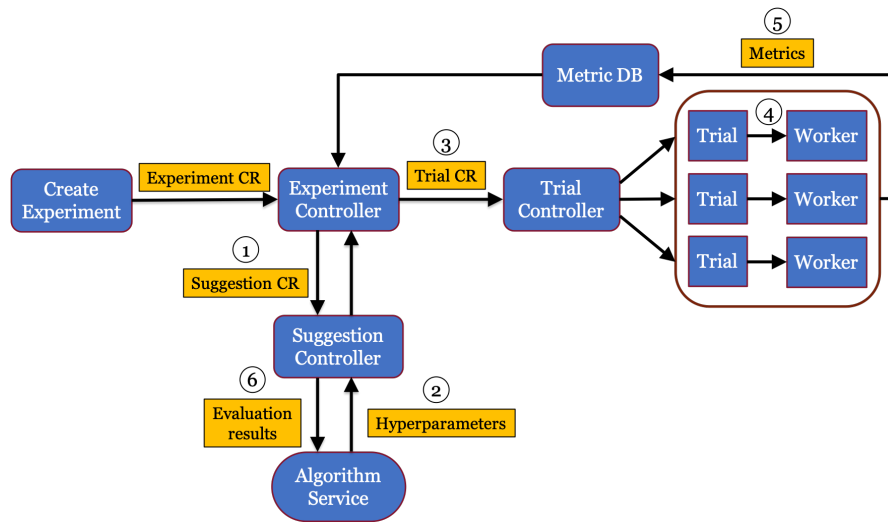


Figure 1: KubeFlow Katib: System Architecture

Katib follows the iterative workflow illustrated in the diagram. Once an **Experiment** is submitted, the Experiment Controller creates a **Suggestion** object. The Suggestion Controller then launches the AutoML algorithm service, which generates new hyperparameter configurations. Based on these suggestions, the Experiment Controller creates corresponding **Trial** objects. Each Trial is executed as a Kubernetes pod using the Trial template, while its execution is monitored by the Trial Controller. After completion, the Metrics Collector retrieves and stores the reported metrics. These results are then fed back to the algorithm service, enabling the generation of improved suggestions and forming a closed optimization loop.

2.3 Methodological Framework: The "Student" Model

To define our training strategy, we adopt a "Student" analogy within a supervised learning curriculum:

- **Continuous Training:** Unlike standard models that restart from scratch, our model utilizes `warm_start`, allowing the "Student" to retain weights across incremental iterations.
- **Early Stopping (Dynamic Termination):** This acts as a graduation criteria. If the student's performance (Accuracy) plateaus for 5 consecutive rounds (Epochs), we terminate the session to prevent overfitting and save resources.

3 Project Structure

The project follows a modular structure where each file represents a distinct stage in the MLOps lifecycle.

3.1 File Organization

```
1 katib-nas-atelier/  
2 |-- train_continuous_nas.py # Core training logic with Early Stopping  
3 |-- Dockerfile             # Environment blueprint for the training pod  
4 |-- continuous_nas.yaml    # Katib Experiment definition (Search Space)  
5 |-- tmp-pvc-pod.yaml       # Helper pod to extract models from PVC  
6 |-- main.py                # FastAPI inference service  
7 |-- model/                 # Local mount point for PVC artifacts
```

3.2 Workflow

1. **Environment Packaging:** The training script is containerized using Docker.
2. **NAS Orchestration:** Katib interprets the YAML blueprint to launch multiple Trials.
3. **Metric Tracking:** A sidecar container captures accuracy logs from the pods.
4. **Artifact Serialization:** The best model is saved to the Persistent Volume.
5. **Model Serving:** FastAPI mounts the volume and exposes the model to the network.

3.3 Config of the hyperparams

3.4 Tableau des Hyperparamètres et Intervalles pour Iris et MNIST

4 Phase 1: Prerequisites & Installation

Before running the experiment, the infrastructure must be correctly provisioned.

Dataset	Hyperparamètre	Type	Intervalle / Valeurs	Description
Iris	layers	int	1 – 3	Nombre de couches cachées du MLP
	units	int	8 – 32	Nombre de neurones par couche
	lr	double	0.001 – 0.05	Taux d'apprentissage initial
MNIST	layers	int	2 – 4	Nombre de couches cachées du MLP
	units	int	64 – 256	Nombre de neurones par couche
	lr	double	0.0005 – 0.01	Taux d'apprentissage initial

Table 1: Configuration des hyperparamètres pour les expériences Iris et MNIST

4.1 MicroK8s Setup

```

1 sudo snap install microk8s --classic
2 microk8s status --wait-ready
3 microk8s enable dns storage

```

Role: *DNS* is critical for service discovery; *Storage* enables the Host-Path provisioner to create PVCs on our local disk.

4.2 Katib Installation (Standalone)

We install Katib by cloning the official repository and applying the standalone manifests.

```

1 # 1. Clone the repository
2 git clone https://github.com/kubeflow/katib.git
3
4 # 2. Navigate to the standalone installation directory
5 cd katib/manifests/v1beta1/installs/katib-standalone
6
7 # 3. Apply the manifests
8 microk8s kubectl apply -k .
9
10 # 4. Verify Katib Pods are Running
11 microk8s kubectl get pods -n kubeflow

```

4.3 Namespace Configuration

```

1 # Create Namespace
2 microk8s kubectl create namespace kubeflow
3
4 # Enable Metrics Injection
5 microk8s kubectl label namespace kubeflow katib.kubeflow.org/metrics-
  collector-injection=enabled

```

Role: This labels the namespace, signaling Katib to inject a "Metrics Collector" container into every trial pod to read the "accuracy=" string.

4.4 Essential Verification

After installation, verify that the core Katib components are running:

```
1 microk8s kubectl get pods -n kubeflow
```

You should see:

- katib-controller: Ready 1/1 & status Running
- katib-db-manager: Ready 1/1 & status Running
- katib-mysql: Ready 1/1 & status Running
- katib-ui: Ready 1/1 & status Running

```
hp@DESKTOP-1FTUR50:~/katib-iris-v3/mnist$ microk8s kubectl get pods -n kubeflow
NAME                                READY   STATUS    RESTARTS   AGE
katib-controller-85d6f4497d-mvnw9   1/1     Running   2          4d7h
katib-db-manager-6dbf658d74-qhz6l   1/1     Running   3          4d7h
katib-mysql-6d756b5657-vvl6v       1/1     Running   2          4d7h
katib-ui-5954947794-vmvhj          1/1     Running   2          4d7h
```

This ensures the environment is fully operational before launching any experiments.

5 Phase 2: Training & Orchestration

5.1 Defining the Logic and Environment

First, we define the training script that implements the Early Stopping logic.

```
1 # ... imports ...
2 model = MLPClassifier(hidden_layer_sizes=tuple([args.units]*args.layers),
3                       warm_start=True, max_iter=1)
4
5 # Continuous Training Loop
6 for epoch in range(50):
7     model.fit(X_train, y_train)
8     acc = accuracy_score(y_test, model.predict(X_test))
9     print(f"accuracy={acc}") # Metric for Katib
10
11 # Early Stopping Check
12 if acc > best_acc:
13     best_acc, wait = acc, 0
14 else:
15     wait += 1
16     if wait >= patience: break
17
18 # Persist to PVC
19 with open(f"{args.export_path}/model.pkl", "wb") as f:
20     pickle.dump(model, f)
```

Listing 1: train_continuous_nas.py

5.2 Docker Environment Definition

We define a lean container image including `scipy` and our training logic.

```
1 FROM python:3.10-slim
2
3 WORKDIR /app
4 COPY train_continuous_nas.py /app/train_continuous_nas.py
5
6 RUN pip install --no-cache-dir numpy scipy scikit-learn
7
8 # Corrected Entrypoint to match the copied file
9 ENTRYPOINT ["python3", "/app/train_continuous_nas.py"]
```

Listing 2: Dockerfile

5.3 Execution

We build the image, import it into the cluster registry, and submit the experiment manifest.

```
1 # 1. Build the Docker image
2 docker build -t katib-iris:v3 .
3
4 # 2. Import image to MicroK8s registry
5 docker save katib-iris:v3 | microk8s ctr image import -
6
7 # 3. Submit the Experiment to the Cluster
8 microk8s kubectl apply -f continuous_nas.yaml
```

Role: *Build* creates a portable environment; *Import* pushes that image into the MicroK8s internal registry; *Apply* triggers the Katib controller to start the "Bi-level optimization" process.

```
---> Removed intermediate container f4938f855724
---> 412a0d387866
Step 5/5 : ENTRYPOINT ["python3", "/app/train.py"]
---> Running in f1cf6f8527a4
---> Removed intermediate container f1cf6f8527a4
---> d9e0c4ff6d5e
Successfully built d9e0c4ff6d5e
Successfully tagged katib-iris:v3
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ docker save katib-iris:v3 | microk8s ctr image import -
[sudo] password for hp:
unpacking docker.io/library/katib-iris:v3 (sha256:b5c878ce89d7ab56d311115c2ea6c48cc7080c07b25f1f800d504231d126119d)...done
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl apply -f continuous_nas.yaml
experiment.kubeflow.org/iris-continuous-nas created
```

Figure 2: image building and submission of the experiment

5.4 Verification

```
1 # Check Pods
2 microk8s kubectl get pods -n kubeflow
3 # Check Trials
```

```

4 microk8s kubectl get trials -n kubeflow
5 # Access UI
6 microk8s kubectl port-forward -n kubeflow svc/katib-ui 8080:80

```

```

hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl get pods -n kubeflow
NAME                                READY   STATUS    RESTARTS   AGE
iris-continuous-nas-7988zrq2-hd5br  2/2     Running   0           8s
iris-continuous-nas-dw7gw4sl-f492d  2/2     Running   0           8s
iris-continuous-nas-random-66578c6b66-24bcz  1/1     Running   0          26s
katib-controller-85d6f4497d-mvnmw9  1/1     Running   0          160m
katib-db-manager-6dbf658d74-qhz6l   1/1     Running   0          160m
katib-mysql-6d756b5657-vvl6v        1/1     Running   0          160m
katib-ui-5954947794-vmvhj           1/1     Running   0          160m

```

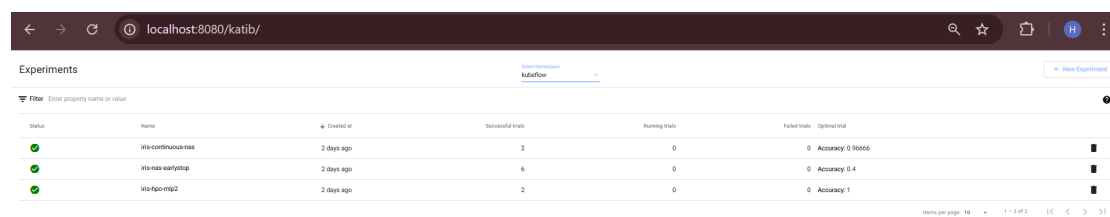
Figure 3: Pods status verification

```

hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl get trials -n kubeflow
NAME                                TYPE          STATUS    AGE
iris-continuous-nas-7988zrq2        Succeeded     True      23s
iris-continuous-nas-dw7gw4sl        Succeeded     True      23s

```

Figure 4: Trials List and Status



Status	Name	Created at	Successful trials	Running trials	Failed trials	Optimal trial
✔	iris-continuous-nas	2 days ago	2	0	0	Accuracy: 0.9666
✔	iris-nas-earlystop	2 days ago	6	0	0	Accuracy: 0.4
✔	iris-tpo-mp2	2 days ago	2	0	0	Accuracy: 1

Figure 5: Katib UI

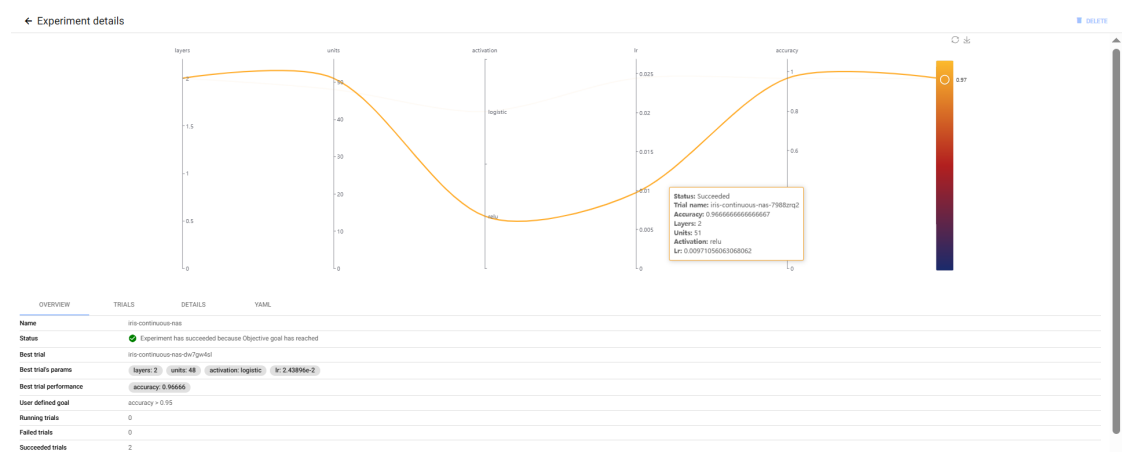


Figure 6: Expiremnet details

OVERVIEW		TRIALS	DETAILS	YAML		
Filter <input type="text" value="Enter property name or value"/>						
Status	Trial name	Accuracy	Layers	Units	Activation	Lr
✓	iris-continuous-nao-7988mq2	0.96666	2	51	relu	9.71056e-3
✓	iris-continuous-nao-dw7gwfdl	0.96666	2	48	logistic	2.43896e-2
						Items per page: 10 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20

Figure 7: The complete set of trials, where the best performing iteration is marked in yellow

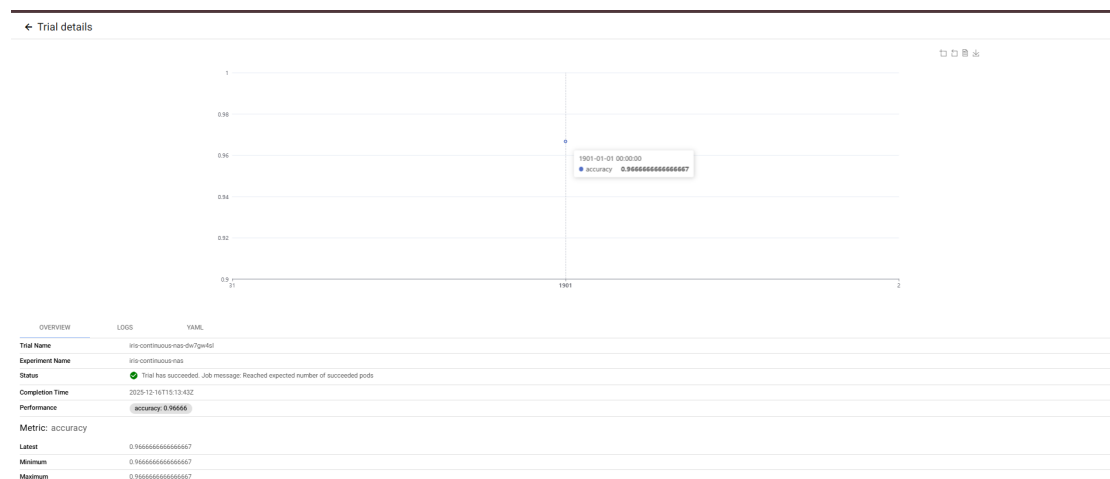


Figure 8: Best Trial details

6 Phase 3: Persistence and Serving

```
1 microk8s kubectl get pvc -n kubeflow
```

Role: Confirms that our "Virtual Disk" is **Bound**. Without this, the training script has no place to save the .pkl file.

```
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl get storageclass
NAME                                PROVISIONER            RECLAIMPOLICY    VOLUMEBINDINGMODE    ALLOWVOLUMEEXPANSION    AGE
microk8s-hostpath (default)        microk8s.io/hostpath   Delete           WaitForFirstConsumer  false                   10h
```

Figure 9: Verification of Persistent Volume Claim status.

```
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl get pvc -n kubeflow
NAME                                STATUS    VOLUME                                     CAPACITY   ACCESS MODES   STORAGECLASS    VOLUMEATTRIBUTESCLASS    AGE
katib-iris-models-pvc              Bound     pvc-70b2678a-bb16-4f4a-9623-8df365f217ed  1Gi        RWO            microk8s-hostpath <unset>                  137m
katib-mysql                        Bound     pvc-a2e0d183-ffaf-4a4e-87e8-ac4cd1abcb74  10Gi       RWO            microk8s-hostpath <unset>                  161m
```

Figure 10: PVC verification in "Bound" state

6.1 Model Recuperation: The "Bridge" Pod Strategy

Since Katib trial pods are ephemeral (they terminate upon completion), direct access to the generated artifacts is lost once the trial finishes. To recuperate the trained model from the **Persistent Volume Claim (PVC)**, we must employ a temporary "Bridge" Pod.

6.1.1 Defining the Inspection Pod

We define a lightweight Pod configuration (`tmp-pvc-pod.yaml`) whose sole purpose is to mount the shared storage volume and sleep, allowing us time to access the file system.

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: tmp-pvc-pod
5   namespace: kubeflow
6 spec:
7   containers:
8     - name: pvc-container
9       image: python:3.10-slim
10      # Sleep keeps the pod alive so we can enter it
11      command: ["sleep", "3600"]
12      volumeMounts:
13        - name: model-volume
14          mountPath: /model
15   volumes:
16     - name: model-volume
17       persistentVolumeClaim:
```

```

18 # Must match the PVC name used in the Experiment
19 claimName: katib-iris-models-pvc

```

Listing 3: tmp-pvc-pod.yaml

6.1.2 Execution and Extraction

Once the manifest is defined, we deploy the pod and use `kubectl cp` to transfer the artifact from the Kubernetes cluster to the local development environment.

```

1 # 1. Deploy the bridge pod
2 microk8s kubectl apply -f tmp-pvc-pod.yaml
3
4 # 2. Wait for the Pod to be 'Running', then verify artifact existence
5 microk8s kubectl exec -n kubeflow -it tmp-pvc-pod -- ls -l /model/
6
7 # 3. Recuperate the model (Copy from Cluster to Local Host)
8 microk8s kubectl cp kubeflow/tmp-pvc-pod:/model/model.pkl ./model.pkl

```

```

hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl apply -f tmp-pvc-pod.yaml
pod/tmp-pvc-pod created
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl get pods -n kubeflow
NAME                                READY   STATUS    RESTARTS   AGE
katib-controller-85d6f4497d-mvnw9   1/1     Running   0           171m
katib-db-manager-6dbf658d74-qhz6l   1/1     Running   0           171m
katib-mysql-6d756b5657-vv16v        1/1     Running   0           171m
katib-ui-5954947794-vmvhj           1/1     Running   0           171m
tmp-pvc-pod                         1/1     Running   0           8s
hp@DESKTOP-1FTUR50:~/katib-iris-v3/v3$ microk8s kubectl exec -n kubeflow -it tmp-pvc-pod -- /bin/sh
/model/
# ls /model/
best_model.pkl  model.pkl

```

Figure 11: tmp-pvc-pod creation and model artifact existence verification

```

hp@DESKTOP-1FTUR50:~/mnt/c/Users/hp$ microk8s kubectl cp kubeflow/tmp-pvc-pod:/model/best_model.pkl best_model2.pkl
hp@DESKTOP-1FTUR50:~/mnt/c/Users/hp$ microk8s kubectl cp kubeflow/tmp-pvc-pod:/model/model.pkl model2.pkl

```

Figure 12: Model recuperation from cluster to local host

Note: After the file is successfully copied, the temporary pod can be deleted to free up resources.

6.2 FastAPI Deployment and Testing

The following logic ensures the model is available for external requests.

```

1 from fastapi import FastAPI
2 import pickle, numpy as np
3
4 app = FastAPI()
5 model = pickle.load(open("model.pkl", "rb"))
6

```

```

7 @app.post("/predict/")
8 def predict(data: list):
9     # Converts input JSON into NumPy array for inference
10    prediction = model.predict([np.array(data)])
11    return {"prediction": prediction.tolist()}

```

Listing 4: main.py Deployment Logic

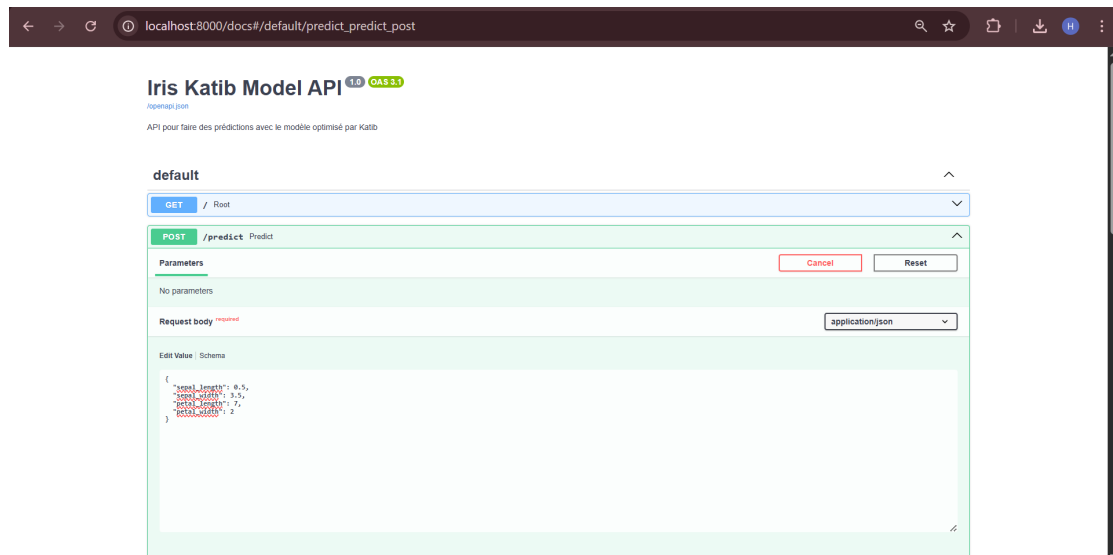


Figure 13: FastAPI

6.2.1 Testing the API

Once the FastAPI server is running (`uvicorn main:app -reload`), we verify it using `curl`:

```

1 curl -X POST "http://localhost:8000/predict/" \
2     -H "Content-Type: application/json" \
3     -d "[5.1, 3.5, 1.4, 0.2]"

```

Expected Output: A JSON response containing the class prediction.

```

PS C:\Users\hp> Invoke-RestMethod `
>> -Uri http://localhost:8000/predict `
>> -Method POST `
>> -ContentType "application/json" `
>> -Body '{ "sepal_length": 5.1, "sepal_width": 3.5, "petal_length": 1.4, "petal_width": 0.2 }'

prediction
-----
1

```

Figure 14: API testing using PowerShell

7 Conclusion

This atelier demonstrated the complete lifecycle of an MLOps project. We successfully:

1. **Installed and Configured** a Kubernetes environment with Kubeflow Katib.
2. **Implemented Continuous NAS**, automating the search for the optimal neural network architecture using Early Stopping.
3. **Orchestrated Persistence** using PVCs, ensuring valuable model artifacts were not lost.
4. **Deployed the Solution** using FastAPI, bridging the gap between research (training) and production (serving).

This workflow ensures that our ML solutions are reproducible, scalable, and ready for real-world deployment.

8 Part 2: MNIST Training and Trial Management

8.1 Introduction

In this phase, we train a Fully Connected Neural Network (FCNN) on the MNIST dataset. The goal is to demonstrate the use of Katib for automatic hyperparameter optimization and trial tracking. Each trial corresponds to a different combination of parameters, allowing us to identify the model with the best performance.

8.2 Objectives

- Create a parameterized Python script for training the MNIST model.
- Define a Dockerfile to execute this script on Kubernetes.
- Configure a Katib Experiment to launch multiple trials and track performance.
- Identify the best model and save its weights to the persistent volume.

8.3 Python Script: train_mnist.py

This script utilizes PyTorch to train the model, dynamically accepting hyperparameters via command-line arguments.

```
1 import argparse
2 import torch
3 from torch import nn, optim
4 from torchvision import datasets, transforms
5 from torch.utils.data import DataLoader
6 import pickle
7
8 # Argument Parsing
```

```

9  parser = argparse.ArgumentParser()
10 parser.add_argument("--layers", type=int, required=True)
11 parser.add_argument("--units", type=int, required=True)
12 parser.add_argument("--lr", type=float, required=True)
13 parser.add_argument("--export_path", type=str, required=True)
14 args = parser.parse_args()
15
16 # Data Loading
17 transform = transforms.Compose([transforms.ToTensor(), transforms.
    Normalize((0.5,), (0.5,))])
18 train_dataset = datasets.MNIST(root="./data", train=True, download=True,
    transform=transform)
19 test_dataset = datasets.MNIST(root="./data", train=False, download=True,
    transform=transform)
20 train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
21 test_loader = DataLoader(test_dataset, batch_size=1000)
22
23 # Model Definition
24 layers = [784] + [args.units]*args.layers + [10]
25 model = nn.Sequential(*[nn.Linear(layers[i], layers[i+1]) for i in range(
    len(layers)-1)])
26 loss_fn = nn.CrossEntropyLoss()
27 optimizer = optim.Adam(model.parameters(), lr=args.lr)
28
29 # Training with Early Stopping Logic
30 best_acc = 0
31 patience = 5
32 wait = 0
33
34 for epoch in range(20):
35     model.train()
36     for X, y in train_loader:
37         X = X.view(X.size(0), -1)
38         optimizer.zero_grad()
39         output = model(X)
40         loss = loss_fn(output, y)
41         loss.backward()
42         optimizer.step()
43
44     # Validation
45     model.eval()
46     correct, total = 0, 0
47     with torch.no_grad():
48         for X, y in test_loader:
49             X = X.view(X.size(0), -1)
50             pred = model(X).argmax(dim=1)
51             correct += (pred == y).sum().item()
52             total += y.size(0)
53     acc = correct / total
54     print(f"accuracy={acc}") # Metric for Katib
55
56     if acc > best_acc:
57         best_acc, wait = acc, 0
58     else:

```



```

59         wait += 1
60         if wait >= patience:
61             break
62
63 # Save Model
64 with open(f"{args.export_path}/mnist_model.pkl", "wb") as f:
65     pickle.dump(model, f)

```

Listing 5: train_mnist.py

8.4 Dockerfile for MNIST

```

1 FROM python:3.10-slim
2
3 WORKDIR /app
4 COPY train_mnist.py /app/train_mnist.py
5
6 RUN pip install --no-cache-dir torch torchvision numpy
7
8 ENTRYPOINT ["python3", "/app/train_mnist.py"]

```

Listing 6: Dockerfile

8.5 Katib Experiment YAML: mnist-nas.yaml

```

1 apiVersion: kubeflow.org/v1beta1
2 kind: Experiment
3 metadata:
4   name: mnist-nas
5   namespace: kubeflow
6 spec:
7   maxTrialCount: 5
8   parallelTrialCount: 2
9   maxFailedTrialCount: 3
10  objective:
11    type: maximize
12    goal: 0.98
13    objectiveMetricName: accuracy
14  algorithm:
15    algorithmName: random
16  parameters:
17    - name: layers
18      parameterType: int
19      feasibleSpace:
20        min: "1"
21        max: "3"
22    - name: units
23      parameterType: int
24      feasibleSpace:
25        min: "64"
26        max: "256"
27    - name: lr

```

```

28     parameterType: double
29     feasibleSpace:
30       min: "0.0005"
31       max: "0.01"
32   trialTemplate:
33     primaryContainerName: training-container
34     trialParameters:
35       - name: layers
36         description: Number of hidden layers
37         reference: "{{.TrialParameters.layers}}"
38       - name: units
39         description: Units per layer
40         reference: "{{.TrialParameters.units}}"
41       - name: lr
42         description: Learning rate
43         reference: "{{.TrialParameters.lr}}"
44   trialSpec:
45     apiVersion: batch/v1
46     kind: Job
47     spec:
48       backoffLimit: 0
49       template:
50         spec:
51           restartPolicy: Never
52           containers:
53             - name: training-container
54               image: katib-mnist:latest
55               command:
56                 - "python3"
57                 - "/app/train_mnist.py"
58               args:
59                 - "--layers={{.TrialParameters.layers}}"
60                 - "--units={{.TrialParameters.units}}"
61                 - "--lr={{.TrialParameters.lr}}"
62                 - "--export_path=/model"

```

Listing 7: mnist-nas.yaml

8.6 Execution and Trial Tracking

We build the specific MNIST image and launch the experiment.

```

1 # 1. Build Docker image
2 docker build -t katib-mnist:latest .
3
4 # 2. Import image into MicroK8s registry
5 docker save katib-mnist:latest | microk8s ctr image import -
6
7 # 3. Launch the Katib Experiment
8 microk8s kubectl apply -f mnist-nas.yaml
9
10 # 4. Track trials in real-time
11 microk8s kubectl get trials -n kubeflow -w

```

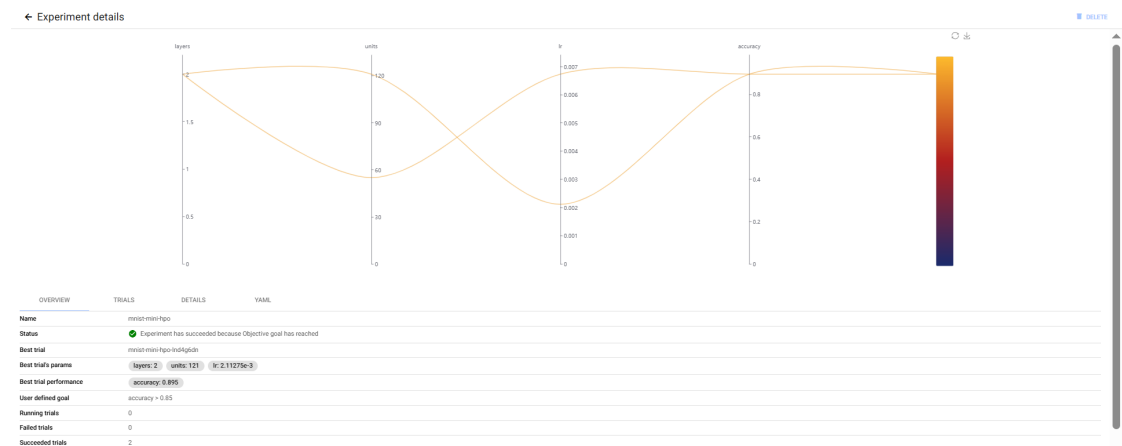


Figure 15: experiment details

OVERVIEW

TRIALS

DETAILS

YAML

Filter

Enter property name or value

Status	Trial name	Accuracy	Layers	Units	Lr
✓	mnist-mnist-ipo-b01b0f5z	0.895	2	55	6.72864e-3
✓	mnist-mnist-ipo-b01d64de	0.895	2	121	2.11275e-3

Items per page: 10

1 - 2 of 2

<

>

>>

Figure 16: Trials

OVERVIEW	TRIALS	DETAILS	YAML
Objective			
Name	accuracy		
Type	maximize		
Goal	0.85		
Additional metrics	No additional metrics		
Trials			
Max failed trials			
Max trials	4		
Parallel trials	2		
Parameters			
layers	Parameter type: int (Min: 1) (Max: 3)		
units	Parameter type: int (Min: 32) (Max: 128)		
lr	Parameter type: double (Min: 0.0005) (Max: 0.01)		
Algorithm			
Name	random		
Metrics collector			
Collector type	StdOut		

Figure 17: Config details

8.7 Conclusion

Through this phase, we have:

- Developed a flexible Python script for training MNIST using PyTorch.

- Created a Dockerfile to execute the model within Kubernetes.
- Utilized Katib to launch multiple trials and optimize hyperparameters.
- Saved the best-performing model for future use.

Cette approche illustre le workflow complet d'optimisation de modèles dans un environnement Kubernetes avec Katib.

9 MNIST Batch Inference Using Katib

After exploring the Iris dataset with continuous NAS, we extend the workflow to the MNIST dataset for batch inference.

9.1 Model Preparation

The MNIST model was pre-trained and stored as `model.pkl` in the PVC:

```
1 -rw-r--r-- 1 root root 3558862 Dec 20 18:26 model.pkl
```

Role: This model will be loaded by a Job for batch inference.

9.2 Bridge Pod for Model Access

Since pods are ephemeral, we use a temporary pod to access the model:

```
1 apiVersion: v1
2 kind: Pod
3 metadata:
4   name: mnist-bridge-pod
5   namespace: kubeflow
6 spec:
7   containers:
8     - name: mnist-container
9       image: python:3.10-slim
10      command: ["sleep", "3600"]
11      volumeMounts:
12        - name: model-volume
13          mountPath: /model
14   volumes:
15     - name: model-volume
16       persistentVolumeClaim:
17         claimName: katib-iris-models-pvc
```

Listing 8: tmp-pvc-pod.yaml for MNIST

9.3 Batch Inference Job

We define a Kubernetes Job to run predictions on MNIST samples:

```

1  apiVersion: batch/v1
2  kind: Job
3  metadata:
4    name: mnist-batch-inference
5    namespace: kubeflow
6  spec:
7    template:
8      spec:
9        restartPolicy: Never
10       containers:
11         - name: mnist-inference
12           image: python:3.10-slim
13           command: ["python3", "/app/mnist_inference.py"]
14           volumeMounts:
15             - name: model-volume
16               mountPath: /model
17       volumes:
18         - name: model-volume
19           persistentVolumeClaim:
20             claimName: katib-iris-models-pvc

```

Listing 9: mnist-batch-job.yaml

9.4 Inference Script

The Python script `mnist_inference.py` loads the model and runs batch predictions:

```

1  import pickle
2  import numpy as np
3  from tensorflow.keras.datasets import mnist
4
5  # Load MNIST data
6  (_, _), (X_test, y_test) = mnist.load_data()
7  X_test = X_test.reshape(len(X_test), -1) / 255.0
8
9  # Load trained model
10 with open("/model/model.pkl", "rb") as f:
11     model = pickle.load(f)
12
13 # Run batch predictions
14 preds = model.predict(X_test[:100])
15 print("Predictions:", preds)

```

9.5 Execution

Deploy the Job and check logs:

```

1  microk8s kubectl apply -f mnist-batch-job.yaml
2  microk8s kubectl get pods -n kubeflow
3  microk8s kubectl logs -n kubeflow <mnist-pod-name>

```

9.6 Visualization

Optionally, the Kubernetes Dashboard can be used to monitor pod status, logs, and PVC mounting:

```
py@DESKTOP-1FTURSO2:/k8s-iris-03/mnist$ microk8s kubectl get jobs -n kubeflow
NAME                                STATUS    COMPLECTIONS   DURATION    AGE
mnist-batch-inference              Running   0/1             24s         24s
py@DESKTOP-1FTURSO2:/k8s-iris-03/mnist$ microk8s kubectl logs -n kubeflow job/mnist-batch-inference
2025-12-20 19:33:11.323260: I external/local_xla/xla/tsl/cuda/cudart_stub.cc:31] Could not find cuda drivers on your machine, GPU will not be used.
2025-12-20 19:33:13.061489: I tensorflow/core/platform/cpu_feature_guard.cc:210] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2025-12-20 19:33:20.395646: I external/local_xla/xla/tsl/cuda/cudart_stub.cc:31] Could not find cuda drivers on your machine, GPU will not be used.
Model loaded
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/mnist.npz
11490434/11490434 8s 1us/step
Sample 0: predicted class = 6
Sample 1: predicted class = 5
Sample 2: predicted class = 6
Sample 3: predicted class = 5
Sample 4: predicted class = 7
Sample 5: predicted class = 2
Sample 6: predicted class = 6
Sample 7: predicted class = 4
Sample 8: predicted class = 3
Sample 9: predicted class = 2
```

Figure 18: MNIST Batch Job running in the Kubernetes Dashboard.

9.7 Notes

- The model artifact is shared across experiments via the PVC.
- Using a bridge pod ensures safe retrieval and batch processing without losing data when pods terminate.
- This workflow can be extended to larger datasets and other models.

10 General Conclusion

This laboratory successfully demonstrated the implementation of a complete, production-ready MLOps pipeline, bridging the gap between static model development and dynamic, automated orchestration. By leveraging MicroK8s and Kubeflow Katib, we transitioned from manual hyperparameter tuning to an automated Neural Architecture Search (NAS).

Throughout this project, we explored two distinct use cases that highlighted the versatility of our infrastructure:

1. **The Iris Experiment (Part 1):** Established the fundamental workflow using Scikit-Learn. We successfully automated the search for an optimal MLP architecture, implemented **Early Stopping** to conserve resources, and deployed the final model as a **FastAPI** microservice.
2. **The MNIST Experiment (Part 2):** Demonstrated the scalability of the pipeline by applying it to a Deep Learning task using PyTorch. This phase confirmed that our containerized architecture is framework-agnostic and capable of handling complex image classification tasks.

10.1 Key Competencies Acquired

This atelier provided practical experience in the essential pillars of modern AI engineering:

- **Containerization (Docker):** Encapsulating dependencies to ensure that "it works on my machine" translates to "it works on the cluster."
- **Orchestration (Kubernetes/Katib):** Managing the lifecycle of multiple parallel trials without manual intervention.
- **Persistence (PVC):** Solving the challenge of ephemeral pods by implementing robust storage strategies for model artifacts.
- **Serving (FastAPI):** Exposing trained models as consumable APIs, completing the loop from research to production.

In conclusion, this architecture provides a robust foundation for scaling machine learning workflows. It ensures that models are not only accurate but also reproducible, deployable, and easier to maintain in a real-world industrial context.

11 References

<https://blog.kubeflow.org/katib/>