# Hot-Swapping MLflow models

A pythonic approach on AWS lambda

# Stavros Niafas

ML Engineer | MLflow Ambassador

 https://github.com/sniafas/PyconGR

# Table of contents

PyCon Greece
29-30 August 2025

# 01  Keeping models fresh w/out downtime

# Keeping model fresh w/out downtime

## Updating Models in production 🔄

- Model Drift

- Changes in data / Missing data trends

- Updating a model can cause downtime

## Dynamic Updates 🎯

- Zero Downtime: Important for high availability services

- Speeding up the ML lifecycle

- Avoids the need for a full re-deployment of application.

## Hot-Swapping 🔥

- Event driven architecture

- Replace the current model with the updated, while the service is still running

- Pythonic way to achieve that

# 02  ML Inference w/ Serverless

# ML Inference w/ Serverless

### Serverless over simplicity 🙌🏻

– Confusing term, servers are there, you just not manage them.

- Zero infra management, self-managed provisioning & auto-scaling

- Ideal for simple inference environments

- Straightforward to implement. Just like a regular Python script
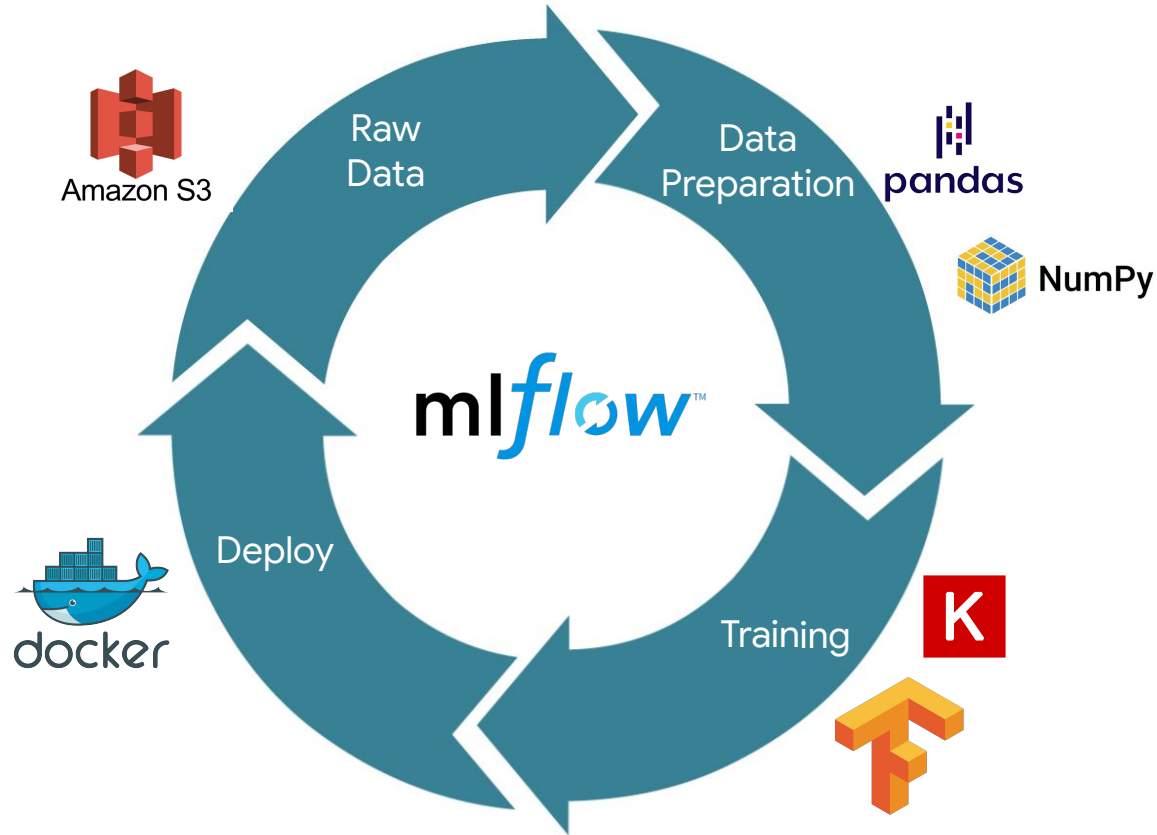
### Reduced time-to-market 📉

- No overhead for API specifics

- Seamless integration, easy to consume

- Unlimited space for automations and parametrization

- Faster experimentation: Commit, Rebuild, Push

### What to look at ⚠️

- Docker packaging & dependencies

- Cold starts 🥶

# 03 ML Lifecycle & MLflow

# ML Lifecycle

# MLflow: Managing the ML lifecycle

1. Implement ML lifecycle

3. Governance to track experiments and models efficiently

2. Fluent API in Python

4. Reduce friction between the teams

# 04  AWS λ

# AWS λ

Cold start initialization

Invocation

Execution Environment

```python
1  # lambda_function.py
2  import ...
3
4  def get_model():
5
6      model = load_from_s3()
7      return model
8
9  def handler(event, context):
10
11     input_data = json.loads(event["body"])
12
13     model = get_model()
14
15     prediction = model.predict(input_data["data"])
16
17     return {
18         "statusCode": 200,
19         "body": {
20             "prediction": prediction.tolist(),
21         }
22     }
```
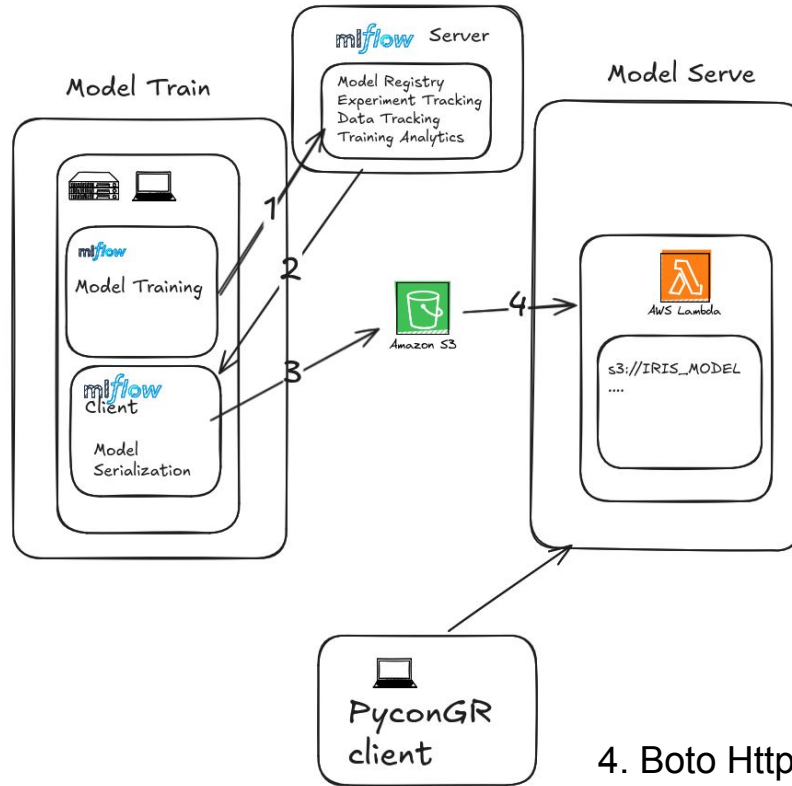
# 05  Hot-Swapping!

# Hot-Swapping



1. Model Traning (MLflow)

2. Model Serialization (Boto)

3. Model Server (PyLambda + Boto)

4. Boto Http Client

# Hot-Swapping
# What's behind zero downtime

Globals →

Time check →

```python
1 model_list, local_timestamp = None, None
2
3 def refresh_model():
4     global model_list, local_timestamp
5
6     # get the current time
7     now = datetime.now()
8
9     # round the time to the minute
10    current_time = now.minute
11
12    # check for updates every N minutes
13    if current_time - local_timestamp > N:
14        new_model_list = load_object_from_s3()
15        local_timestamp = current_time
16        if model_list["version"] != new_model_list["version"]:
17            logger.info("Refreshing models..")
18            load_model_artifacts()
19            model_list = new_model_list
20
21    return model_list
```

PyCon Greece
29-30 August 2025
Technopolis City of Athens

# 05  Live Demo 🤞🚀

# Thank you