

28/01/2026

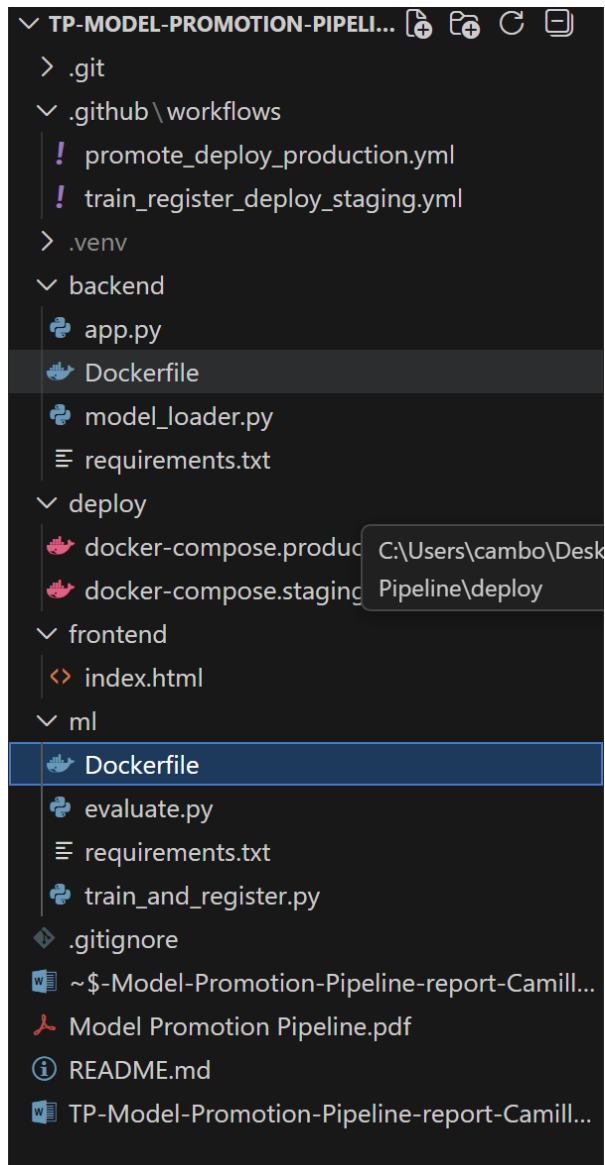
# TP-Model- Promotion- Pipeline

Machine Learning in Production



Camille BORDES  
GROUP DAI

## Configuration :



### Repository variables

Name	↑
DAGSHUB_MLFLOW_TRACKING_URI	
DAGSHUB_TOKEN	
DOCKERHUB_TOKEN	
DOCKERHUB_USERNAME	

## Task 1: Run Candidate -> Staging workflow

```

1  Run echo '{"run_id": "973906543027475ca08b3f182da8ddd7", "accuracy": 0.8975, "model_version": 1}' | python ml/evaluate.py
2  echo '{"run_id": "973906543027475ca08b3f182da8ddd7", "accuracy": 0.8975, "model_version": 1}' | python ml/evaluate.py
3  shell: /usr/bin/bash -e {0}
4  env:
5    pythonLocation: /opt/hostedtoolcache/Python/3.11.14/x64
6    PKG_CONFIG_PATH: /opt/hostedtoolcache/Python/3.11.14/x64/lib/pkgconfig
7    Python_ROOT_DIR: /opt/hostedtoolcache/Python/3.11.14/x64
8    Python2_ROOT_DIR: /opt/hostedtoolcache/Python/3.11.14/x64
9    Python3_ROOT_DIR: /opt/hostedtoolcache/Python/3.11.14/x64
10   LD_LIBRARY_PATH: /opt/hostedtoolcache/Python/3.11.14/x64/lib
11   ACCURACY_THRESHOLD: 0.90
12   {"passed": false, "accuracy": 0.8975, "threshold": 0.9, "model_version": 1, "run_id": "973906543027475ca08b3f182da8ddd7"}
13  Error: Process completed with exit code 2.
```

? **Model version number : 1.**

? **Accuracy : 0.8975.**

? **Did the gate pass? : No**

## Task 2: Explain what "staging" proves

**What staging tests that offline evaluation does not ?**

In a production pipeline, staging serves to prove operational readiness by validating factors that offline evaluation cannot, such as the successful integration of the model within a Flask API, the correctness of Docker environment configurations including all dependencies, and the model's ability to handle real-time API requests and schema compatibility through the /predict endpoint. While offline evaluation focuses solely on mathematical metrics like accuracy, staging ensures the entire system, from model loading in MLflow to serving predictions, works as a cohesive unit in a production-like setting.

## Task 3: Promote to production

I changed the ACCURACY\_THRESHOLD to 0.85 to the model with the Candidate to Staging workflow :

change accuracy #3

Re-run all jobs Latest #2

Summary

All jobs

train\_register  
deploy\_staging  
staging\_smoke\_test

Run details  
Usage  
Workflow file

Re-run triggered 21 minutes ago Status Total duration Artifacts  
camm22 6c684ac main Success 2m 6s -

train\_register\_deploy\_staging.yml  
on: push

train\_register 42s → deploy\_staging 1m 10s → staging\_smoke\_test 3s

Apply Quality Gate 0s

```

1 ▶ Run echo '{"run_id": "f3d558ab8a454c5da714f3f0f642d9c3", "accuracy": 0.8975, "model_version": "2"}' | python ml/evaluate.py
12 {"passed": true, "accuracy": 0.8975, "threshold": 0.85, "model_version": "2", "run_id": "f3d558ab8a454c5da714f3f0f642d9c3"}

```

So now we can launch the Promote to Production :

Promote to Production

Promote to Production #8

Summary

All jobs

promote  
deploy\_production

Run details  
Usage  
Workflow file

Manually triggered 13 minutes ago Status Total duration Artifacts  
camm22 6c684ac main Success 1m 50s -

promote\_deploy\_production.yml  
on: workflow\_dispatch

promote 35s → deploy\_production 1m 9s

Promote model version in MLflow Registry 1s

```

1 ▶ Run python - << 'PV'
33 <stdin>:10: FutureWarning: `mlflow.tracking.client.MlflowClient.transition_model_version_stage` is deprecated since 2.9.0. Model registry
  stages will be removed in a future major release. To learn more about the deprecation of model registry stages, see our migration guide here:
  https://mlflow.org/docs/latest/model-registry.html#migrating-from-stages
34 Promoted model version 1 to Production

```

Task 4: Prove production uses registry stage, not "latest code"

## Staging:

```
PS C:\Users\cambo> curl http://localhost:8000/health

Avertissement de sécurité : risque d'exécution de script
Invoke-WebRequest analyse le contenu de la page web. Il se peut que le code de script de la page web s'exécute lors de l'analyse de la page.
ACTION RECOMMANDÉE :
    Utilisez le commutateur -UseBasicParsing pour éviter l'exécution du code de script.

Voulez-vous continuer ?

[O] Oui [I] Oui pour tout [N] Non [U] Non pour tout [S] Suspendre [?] Aide (la valeur par défaut est « N ») : t

StatusCode      : 200
StatusDescription : OK
Content          : {"stage":"Staging","status":"ok"}

RawContent       : HTTP/1.1 200 OK
                  Connection: close
                  Content-Length: 34
                  Content-Type: application/json
                  Date: Wed, 28 Jan 2026 23:37:01 GMT
                  Server: Werkzeug/3.1.5 Python/3.11.14

                  {"stage":"Staging","status":"ok"}
Forms            : {}
Headers          : {[Connection, close], [Content-Length, 34], [Content-Type, application/json], [Date, Wed, 28 Jan 2026 23:37:01 GMT]...}
Images           : {}
InputFields      : {}
Links            : {}
ParsedHtml       : mshtml.HTMLDocumentClass
RawContentLength : 34
```

## Production:

```
PS C:\Users\cambo> curl http://localhost:8001/health

StatusCode      : 200
StatusDescription : OK
Content          : {"stage":"Production","status":"ok"}

RawContent       : HTTP/1.1 200 OK
                  Connection: close
                  Content-Length: 37
                  Content-Type: application/json
                  Date: Wed, 28 Jan 2026 23:42:59 GMT
                  Server: Werkzeug/3.1.5 Python/3.11.14

                  {"stage":"Production","status":"ok...
Forms            : {}
Headers          : {[Connection, close], [Content-Length, 37], [Content-Type, application/json], [Date, Wed, 28 Jan 2026 23:42:59 GMT]...}
Images           : {}
InputFields      : {}
Links            : {}
ParsedHtml       : mshtml.HTMLDocumentClass
RawContentLength : 37
```

## Discussion questions

### 1. Why is it dangerous to deploy "whatever just merged to main" as the model?

- **Performance Uncertainty:** A code merge only guarantees that the scripts are syntactically correct and pass unit tests, but it does not guarantee that the resulting model meets performance requirements.
- **Quality Control:** Without a "Quality Gate," a model with poor accuracy or high loss could be deployed to production, potentially causing business failures.

- **Stochastic Nature:** Model training is often stochastic; the same code can produce different results across runs, making it essential to validate the specific artifact (the model file) rather than just the code.

## 2. What does the registry stage give you that a Git tag does not?

- **Dynamic Decoupling:** Registry stages (Staging/Production) allow you to update which model is "live" without changing the code or rebuilding Docker images.
- **State Management:** It provides a centralized source of truth for the model's lifecycle that is independent of the Git branch, allowing different environments (ports 8000 vs 8001) to pull specific versions based on their status.
- **Auditability:** Registries maintain metadata about who promoted a model and when, which is more difficult to track and restrict via standard Git tags in a collaborative ML environment.

## 3. If staging passes but production fails, what could be the causes?

- **Data Drift:** The real-world data in production may differ significantly from the test/validation data used in staging, leading to poor model performance.
- **Scale and Load:** Production often handles much higher traffic, which might expose latency issues or memory leaks that weren't visible in the staging "smoke test."
- **Configuration Drift:** Discrepancies in environment variables, hardware (CPU vs GPU), or software dependencies between the two environments.

## 4. Where should DVC fit in a serious pipeline?

- **Training Data Snapshot:** DVC should version the exact dataset used for training to ensure reproducibility of the model version stored in MLflow.
- **Evaluation Dataset Snapshot:** It must version the "Golden Dataset" used by the Quality Gate to ensure that model comparisons over time are consistent.
- **Drift Reference Dataset:** DVC should store the baseline data distribution used to detect if production data has started to deviate from training data.

## 5. What should be added to the gate beyond accuracy?

- **Latency:** Ensure the model predicts within a specific time limit (e.g., < 100ms) to meet User Experience requirements.
- **Schema Checks:** Verify that the input data format and feature types haven't changed, preventing the API from crashing.

- **Fairness Constraints:** Check for biases against specific subgroups to ensure ethical and legal compliance.
- **Adversarial/Robustness Tests:** Test the model against edge cases or intentionally noisy data to ensure it doesn't fail catastrophically in unexpected scenarios