

**Accélérer le déploiement de  
modèles de machine learning**

**1 ML chez DiliTrust**

**2 Maintien en conditions  
opérationnelles**

**3 Impact de l'open source**

**4 Le DevOps appliqué au ML**

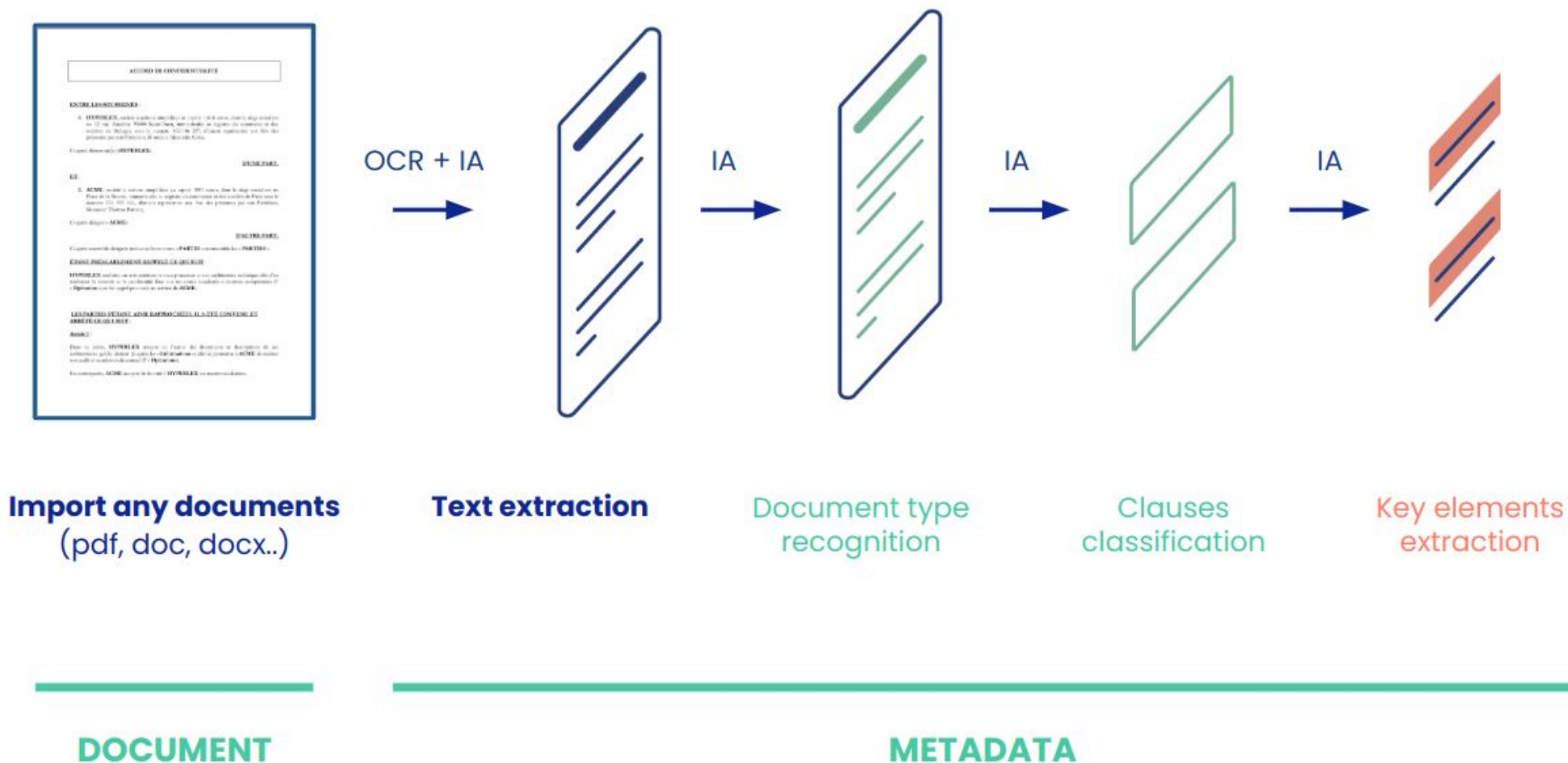
**5 Composants frontaux**

**6 Monitoring**

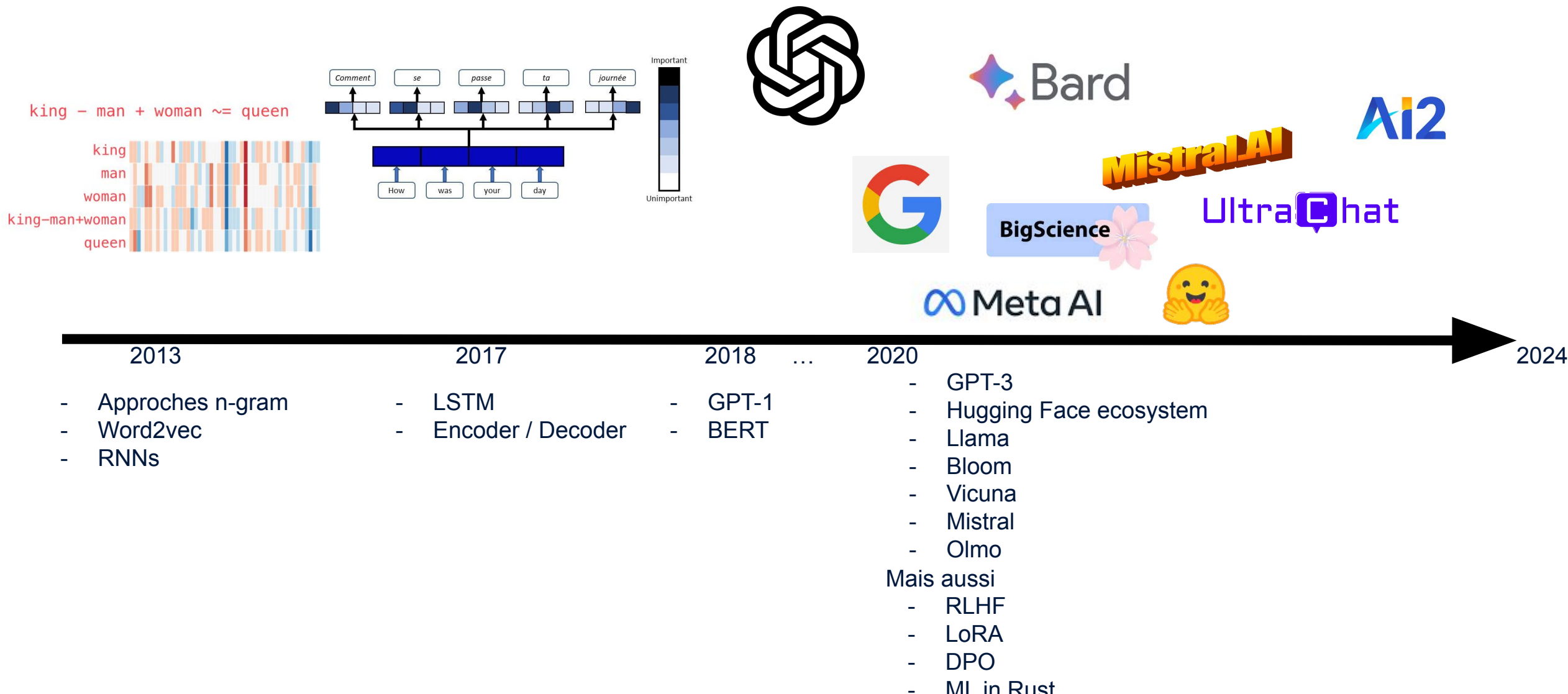
**7 Torchserve**

**8 TGI**

# ML chez DiliTrust



# Evolution des architectures de modèles de NLP



# Evolution des architectures de modèles de NLP



Mixture of Expert

MistralAI

OpenChat update

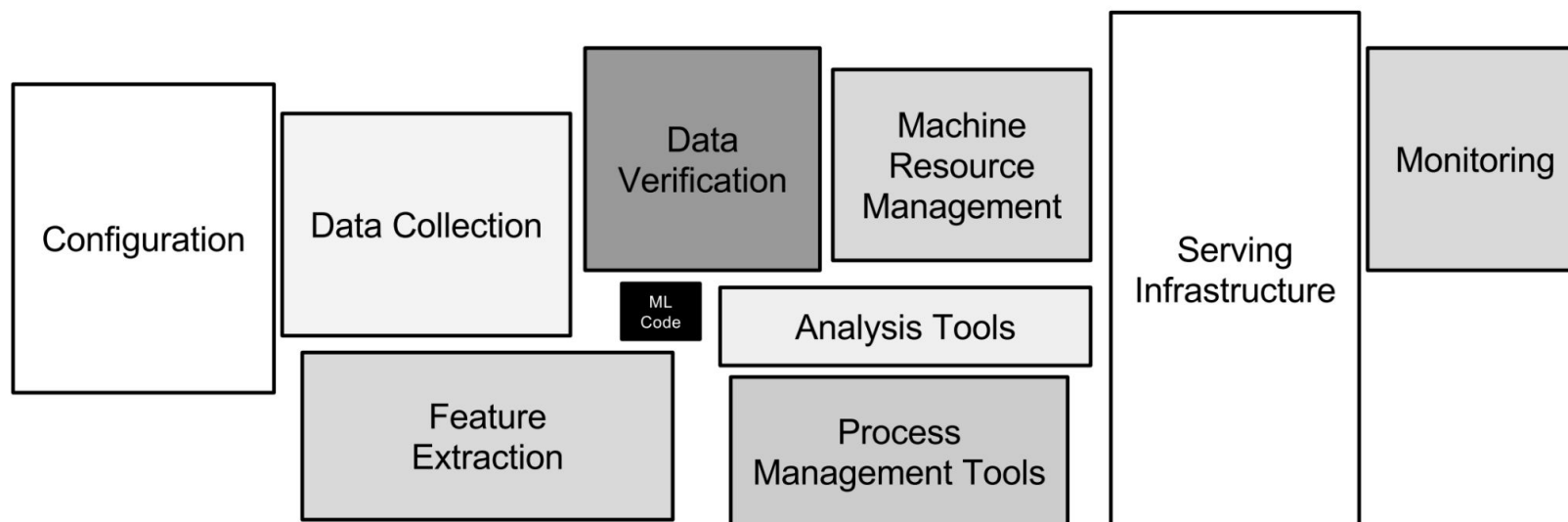
MLX,  
efficient inference on iPhones



Efficient inference on  
distributed serving

Alibaba.com

# Maintenir un modèle de NLP embarqué dans une application



Recherche

Production

- Code volatile
- Scripts, jobs automatisés
- Latences fortes

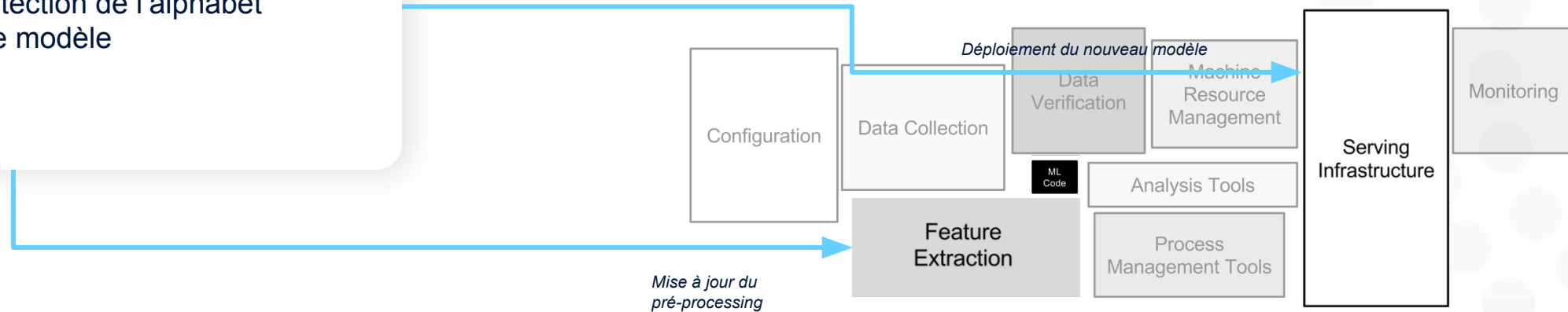
- Stabilité du code
- Applications scalables, relativement urbanisées dans le SI
- Latences faibles

# En pratique, qu'est ce que cela implique ?

Faire évoluer le périmètre fonctionnel

Ajouter une langue au modèle nécessite une nouvelle méthode de pré-processing : détecter l'alphabet avant de passer à l'inférence :

- Déployer la détection de l'alphabet
- Mettre à jour le modèle

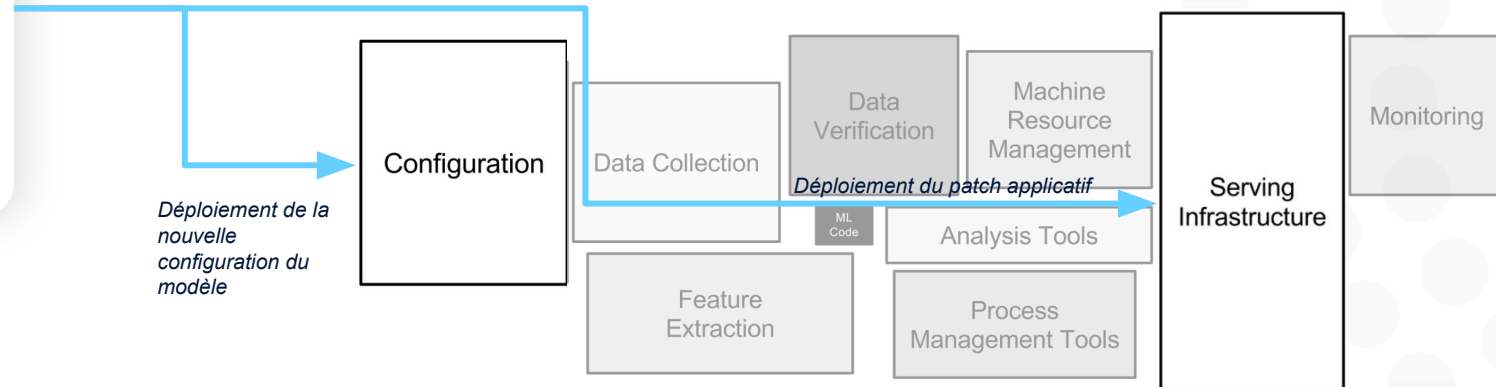


# En pratique, qu'est ce que cela implique ?

## Patcher un bug

Bug détecté en production, nécessitant la mise à jour de la configuration du modèle ainsi que du code servant à préparer le document pour l'inférence

- Déploiement de la configuration du modèle
- Livraison du patch applicatif

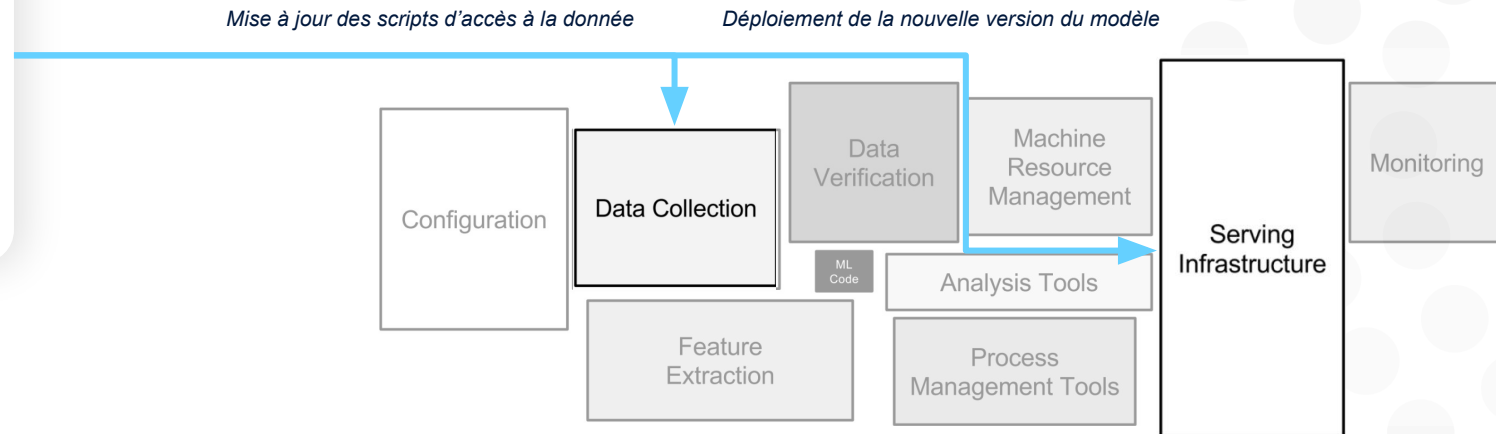




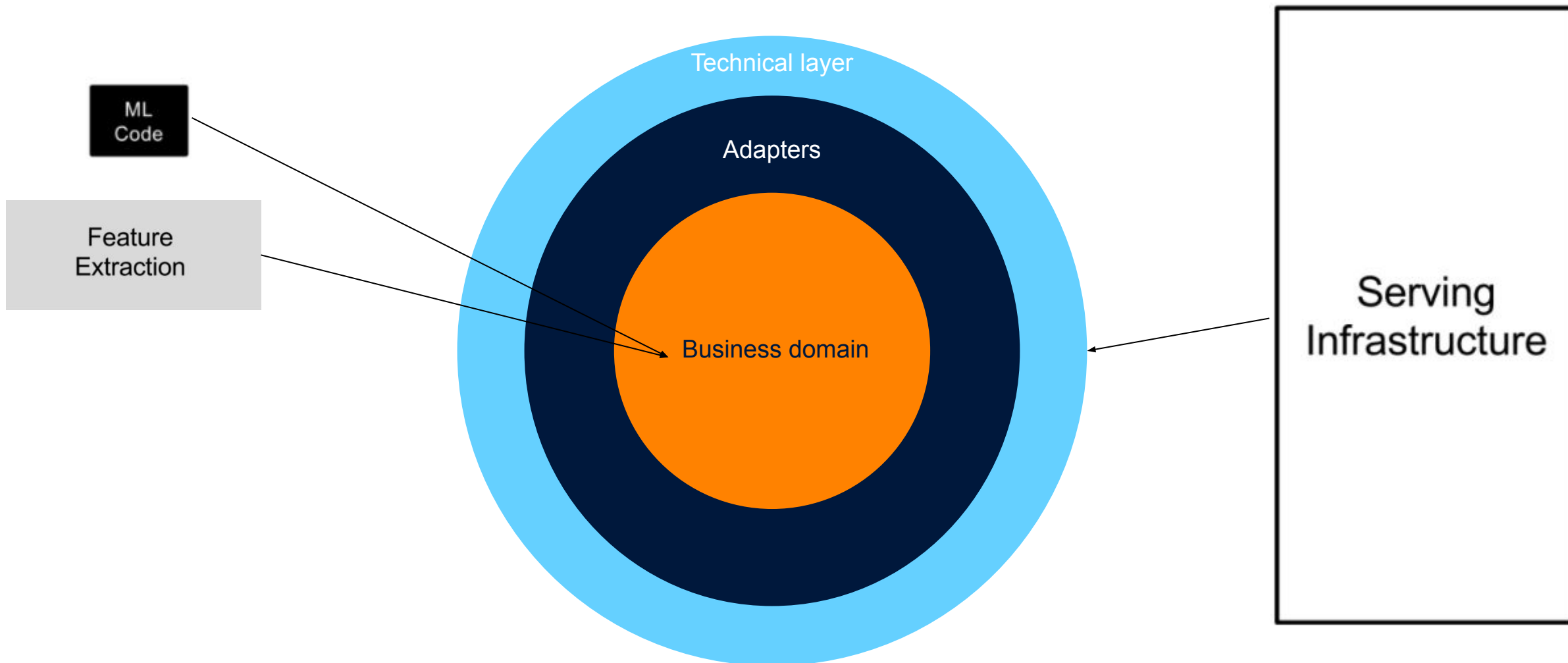
# En pratique, qu'est ce que cela implique ?

Ré-entraîner un modèle

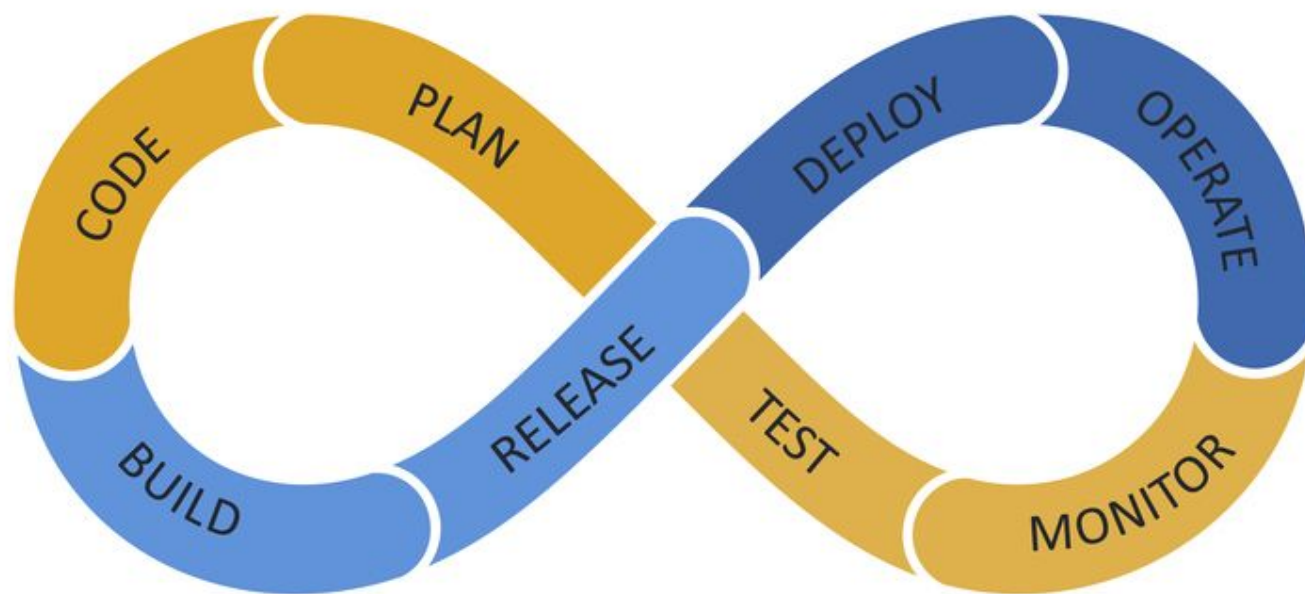
Les annotateurs ont fini une grosse campagne d'analyse de documents, il est temps de réentraîner et déployer une nouvelle version du modèle



## Comment l'architecture hexagonale sert la maintenabilité d'un service basé sur du ML ?



# Comment s'inspirer des pratiques de déploiement du DevOps ?



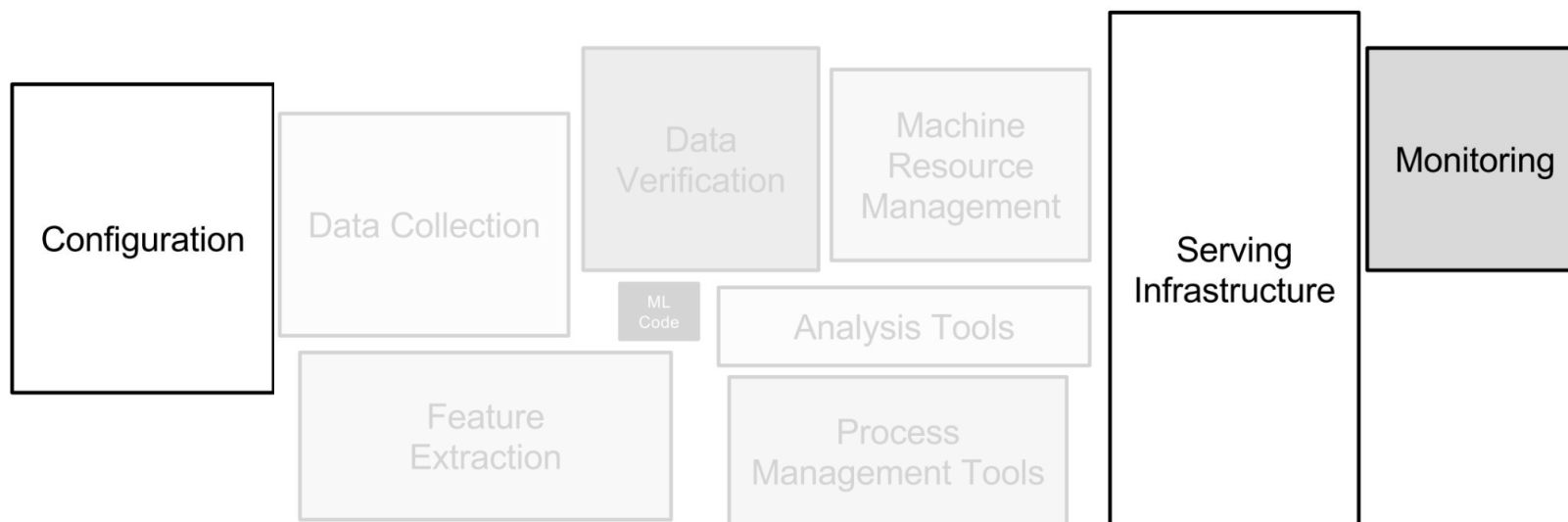
Designons notre pipeline de déploiement en lui assurant ces attributs :

- Idempotente
- Automatisé, pilotable manuellement
- Versionning général, GitOps
- Générique, simple à étendre

Un pilotage plus fin du déploiement permet également :

- Optimization des coûts en diminuant l'empreinte mémoire GPU d'un modèle, un GPU peut exposer plus de modèles
- Itération plus rapide sur les optimisations à l'inférence (quantization, re-write de tooling en Rust)

# Composants frontaux



## Infrastructure :

Code : terraform, Kubernetes

Configuration : GitOps

Runtime: Container docker standardisé

## Serving layer (inférence) :

Configuration : artefact versionné

Modèle : artefact(s) versionné

Code : code applicatif versionné

## Monitoring :

Collecte de métriques

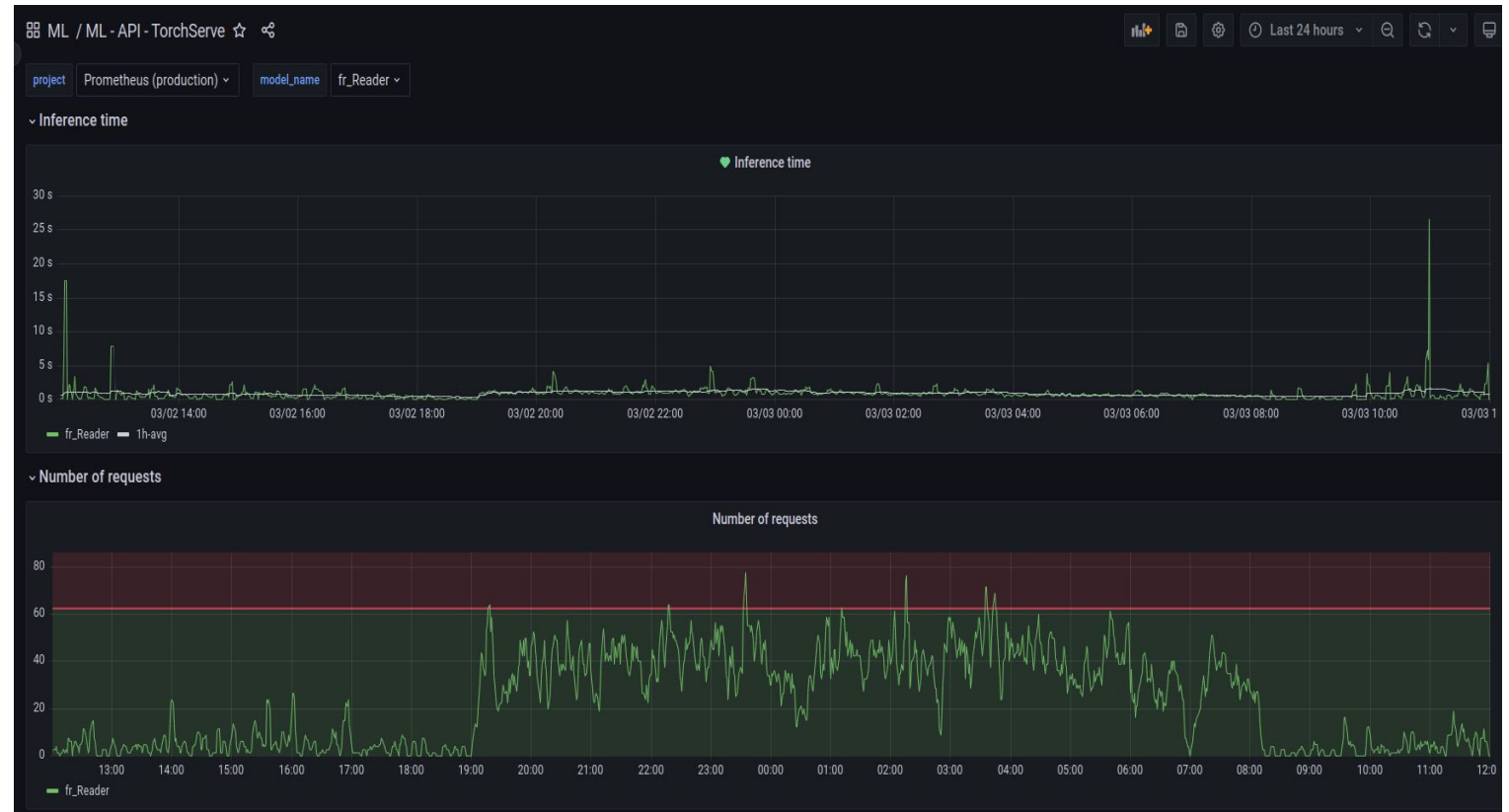
Alerting

Autoscaling

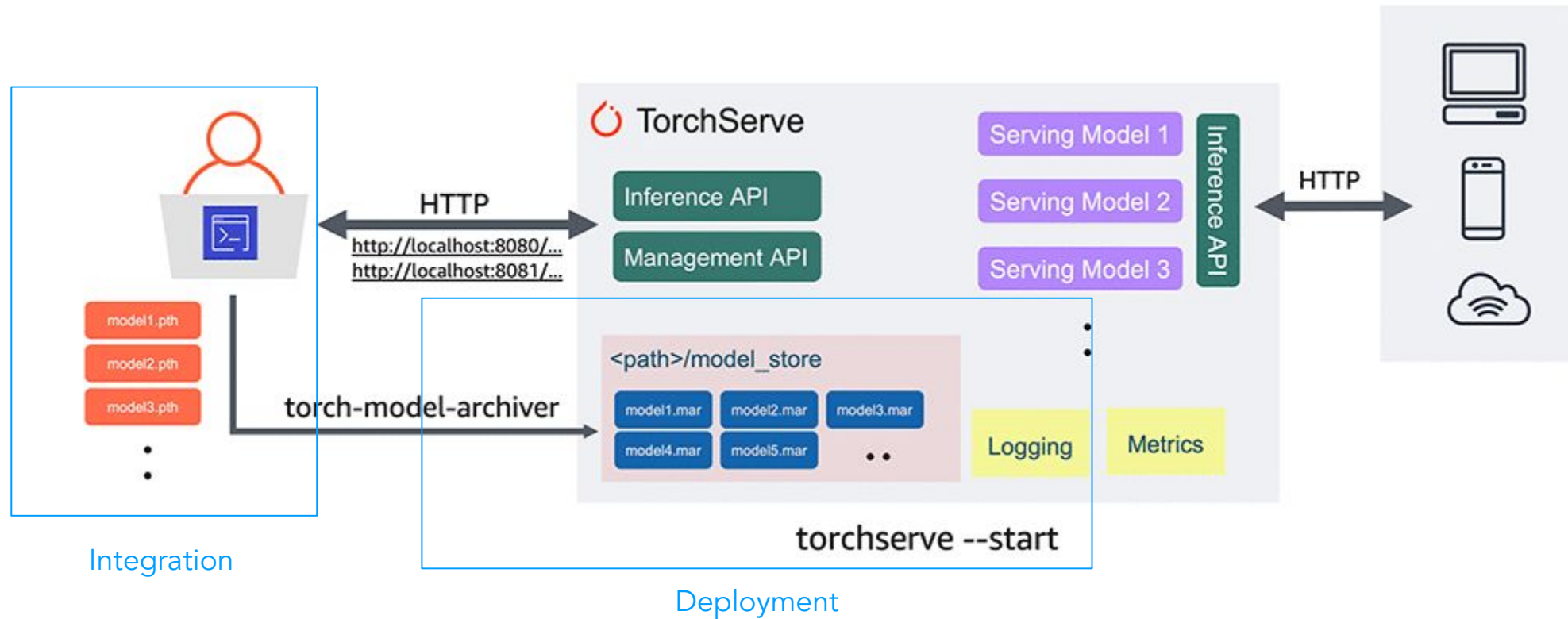
# Monitoring

## Métriques :

- Technique : État de l'infrastructure (empreinte mémoire, latence des requêtes, etc.)
- Fonctionnelle : Pertinence du modèle face aux cas d'usages utilisateurs :
  - Drifts données, drift de concepts (trends des réseaux sociaux, contenu généré par LLM)



# Torchserve



# Torchserve: model & configuration artefacts

```
{
  "fr_reader": {
    "1.0": {
      "marName": "fr_reader.mar",
      "defaultVersion": false,
      "minWorkers": 1,
      "maxWorkers": 1,
      "batchSize": 1,
      "maxBatchDelay": 100,
      "responseTimeout": 120
    }
  }
}
```

```
{
  "QuestionAnswering": {
    "Reader": {
      "fr": {
        "version": 20230207180600,
        "description": "hlx-tools version: ",
        "model_config": {
          "source_bucket_name": "production-models.hyperlex.fr",
          "source_path": "RemoteModels/fr/QuestionAnswering/ReaderWPS",
          "target_bucket_name": "develop-models.hyperlex.fr",
          "target_path": "Deployments/RemoteModels/fr/QuestionAnswering/Reader/"
        }
      }
    }
  }
}
```

Configuration de service cible :

- Aspects techniques du serveur d'inférence
- Aspects applicatifs, sources des modèles



```
{
  "fr_reader": {
    "1.0": {
      "marName": "fr_reader.mar",
      "defaultVersion": false,
      "minWorkers": 1,
      "maxWorkers": 1,
      "batchSize": 1,
      "maxBatchDelay": 100,
      "responseTimeout": 120
    },
    "2": {
      "defaultVersion": true,
      "marName": "fr_reader_2.mar",
      "minWorkers": 1,
      "responseTimeout": 120
    }
  }
}
```

Configuration à déployer

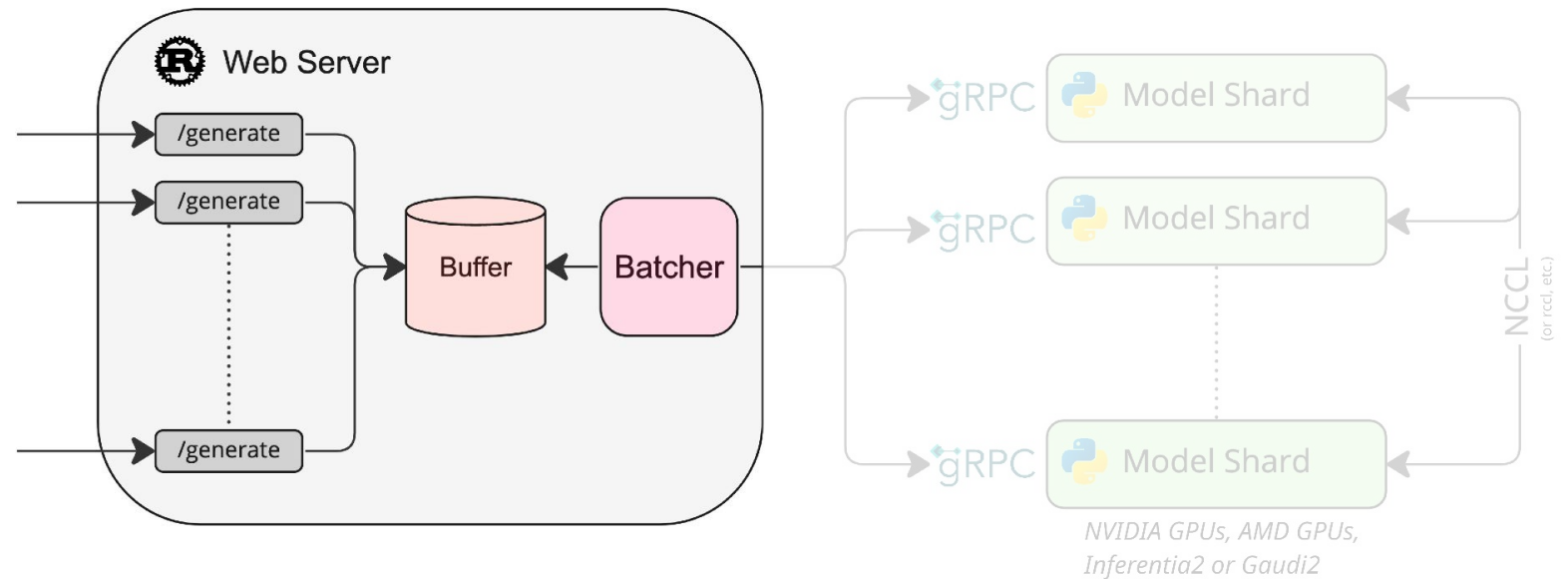
Model binary &  
configuration

Model archive

# Text Generation Inference

Passer l'infrastructure d'inférence à l'échelle, pour exposer des modèles extrêmement performants.

- Back end d'exécution en Rust
- Interfaces en Python
- Intégration des features d'optimisation des performances au fil de l'état de l'art

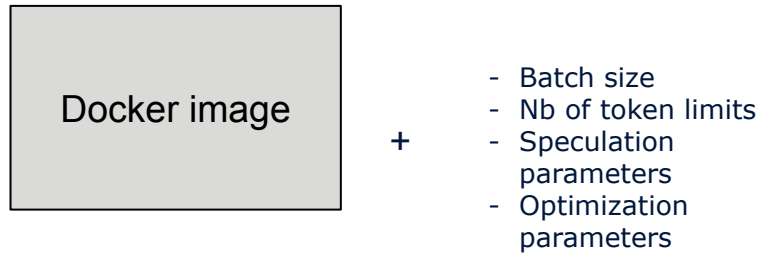




## Quelques techniques d'optimisation

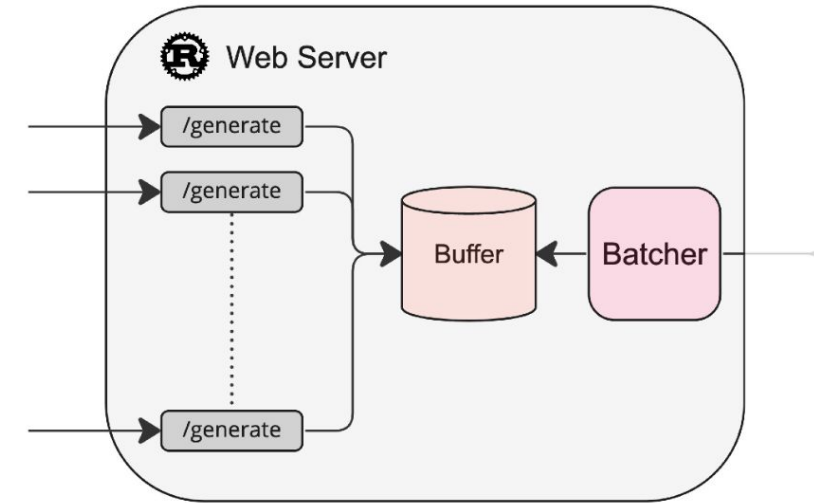
Feature	TGI - Text Generation Interface	vLLM - Versatile Large Language Model
<b>Tensor Parallelism</b>	✓ Utilizes Tensor Parallelism for faster inference	✓ Leverages Tensor Parallelism for high throughput
<b>Optimized Transformers Code</b>	✓ Optimizes transformer code for inference	✓ Employs optimized CUDA kernels for speed
<b>Quantization</b>	✓ Supports quantization methods	✗ Does not support quantization
<b>Continuous Batching</b>	✓ Batches incoming requests continuously	✓ Maximizes GPU utilization with continuous batching
<b>Accelerated Weight Loading</b>	✓ Accelerated weight loading for quick startup	✗ Does not include accelerated weight loading
<b>Logits Warping</b>	✓ Offers logits warping options	✗ Does not provide logits warping features
<b>Custom Prompt Generation</b>	✓ Allows custom prompt generation	✗ Does not support custom prompt generation
<b>Fine-tuning Support</b>	✓ Supports fine-tuned models for higher accuracy	✗ Lacks fine-tuning support
<b>High Throughput</b>	Moderate throughput	Excellent serving throughput
<b>Paged Attention</b>	✗ Does not include Paged Attention	✓ Efficiently manages memory with Paged Attention
<b>Streaming Outputs</b>	✓ Supports streaming outputs	✓ Enables streaming outputs for improved performance
<b>Support for Various Models</b>	✓ Compatible with various LLMs	✓ Seamlessly supports a wide range of Hugging Face models

# Text Generation Inference : un déploiement sans aspérités



Build du runtime cible

Model binary with  
technical configuration



Déploiement du conteneur autonome

Model Instance

# Questions

Amine Saboni  
@Tundji

[dilitrust.com](https://dilitrust.com)