

Construyendo software resiliente y escalable con MLOps

PyCON ES 2024 (Vigo)





<https://github.com/alice-biometrics/pycones24>

¿Quién Soy?

Daniel Pérez Cabo

- Ingeniero de telecomunicación (UVigo)
- Research ML Engineer



[LinkedIn](#)

¿Qué hago?

Tech Lead @ Core-Tech en  **Alice**

Solución de verificación de identidad ágil y segura



¿Qué hacemos?

- Reconocimiento Facial
- Prueba de vida
- Reconocimiento de texto
- Análisis forense de documentos

¿Cual es genuino?

¿Cuál es genuino?



¿Cual es genuino?

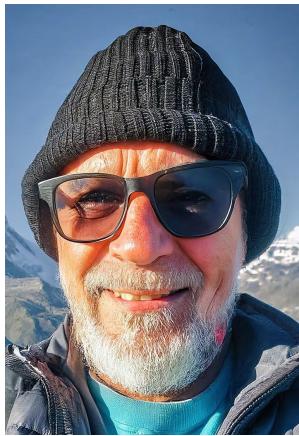


Fake



Genuino

¿Cual es genuino?



¿Cual es genuino?



Fake



Genuino

¿Cual es genuino?



¿Cual es genuino?



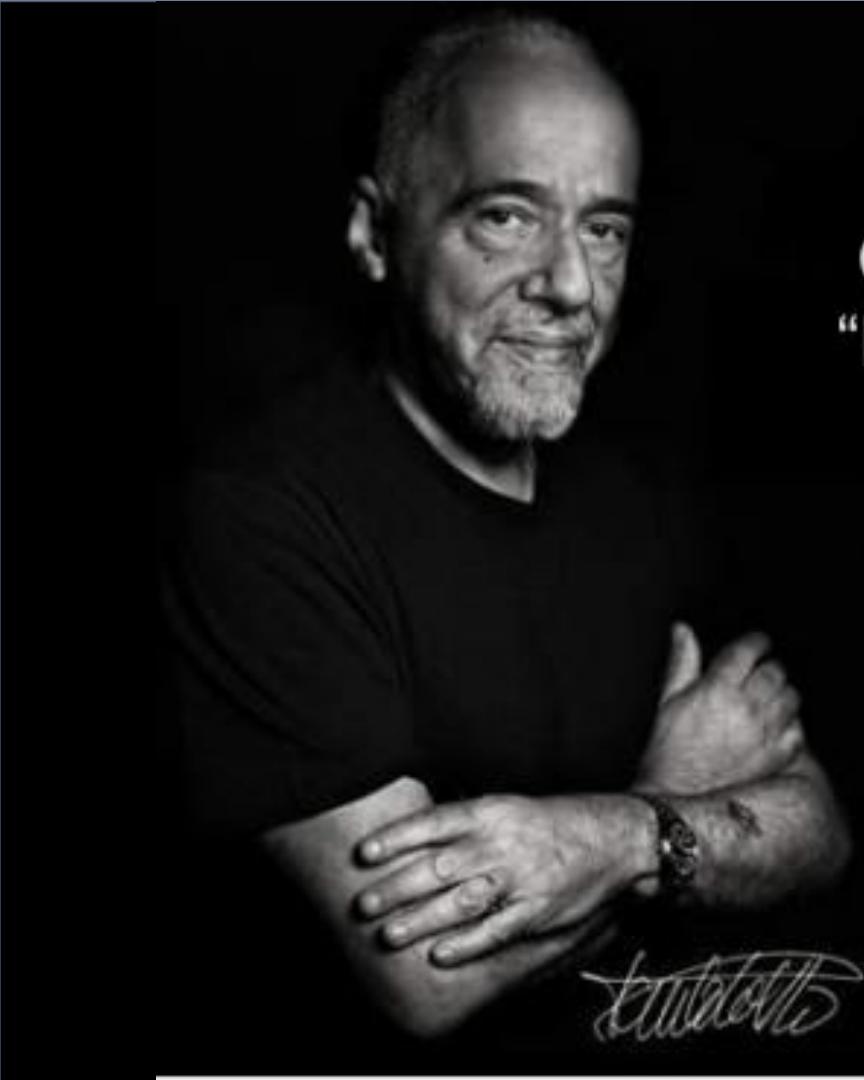
Fake



Fake

Pero ...

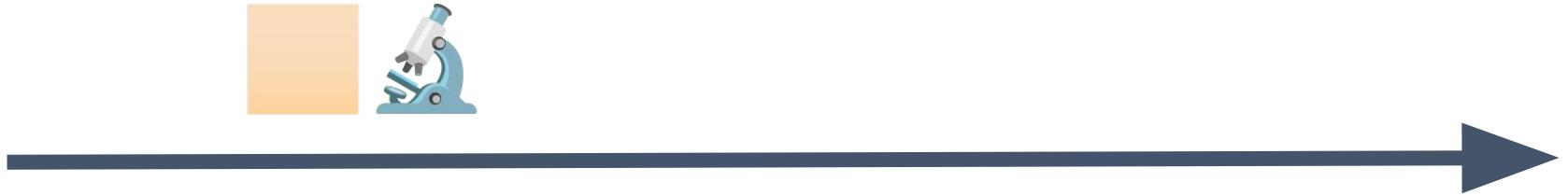
... en qué nos afecta la escalabilidad y la resiliencia del software

A black and white photograph of Paulo Coelho. He is an elderly man with a beard and mustache, wearing a dark t-shirt. He is resting his chin on his hand, looking slightly to the side with a thoughtful expression. His other hand is visible, showing a tattooed forearm and a ring on his finger.

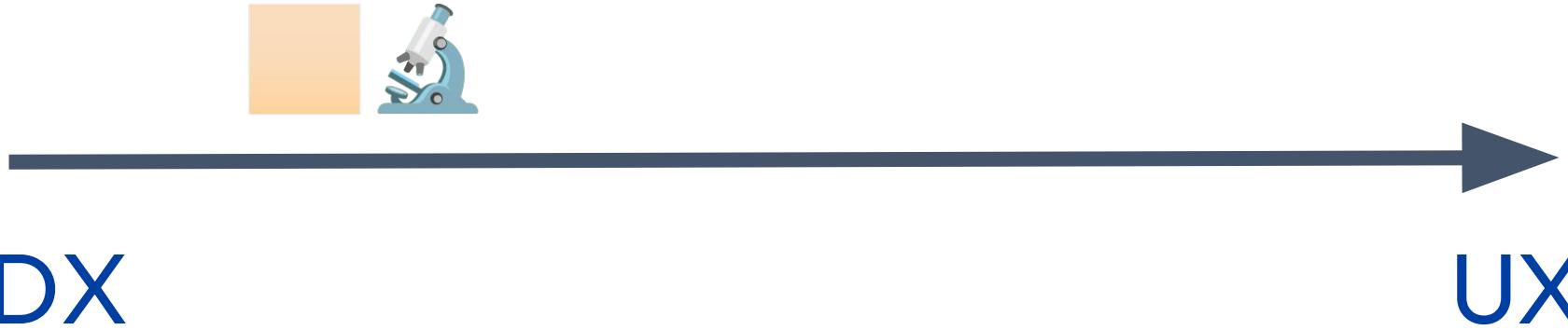
Complejidad Accidental:
“No tenemos tiempo para
perderlo con mierdas”

~Paulo Coelho

A handwritten signature in cursive script, appearing to read "PAULO COELHO".









DX Dev eXperience

UX



DX Dev eXperience
REX

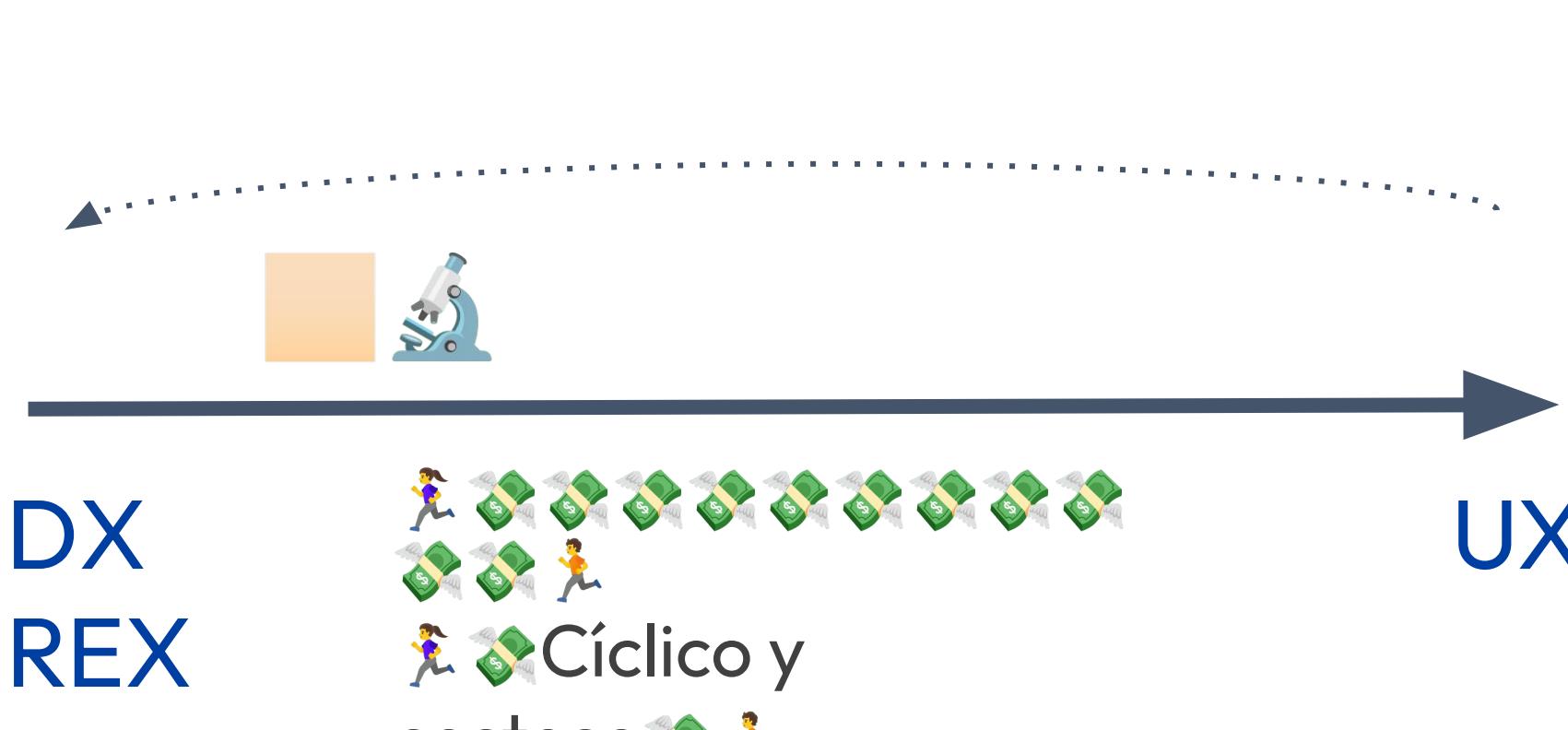
UX

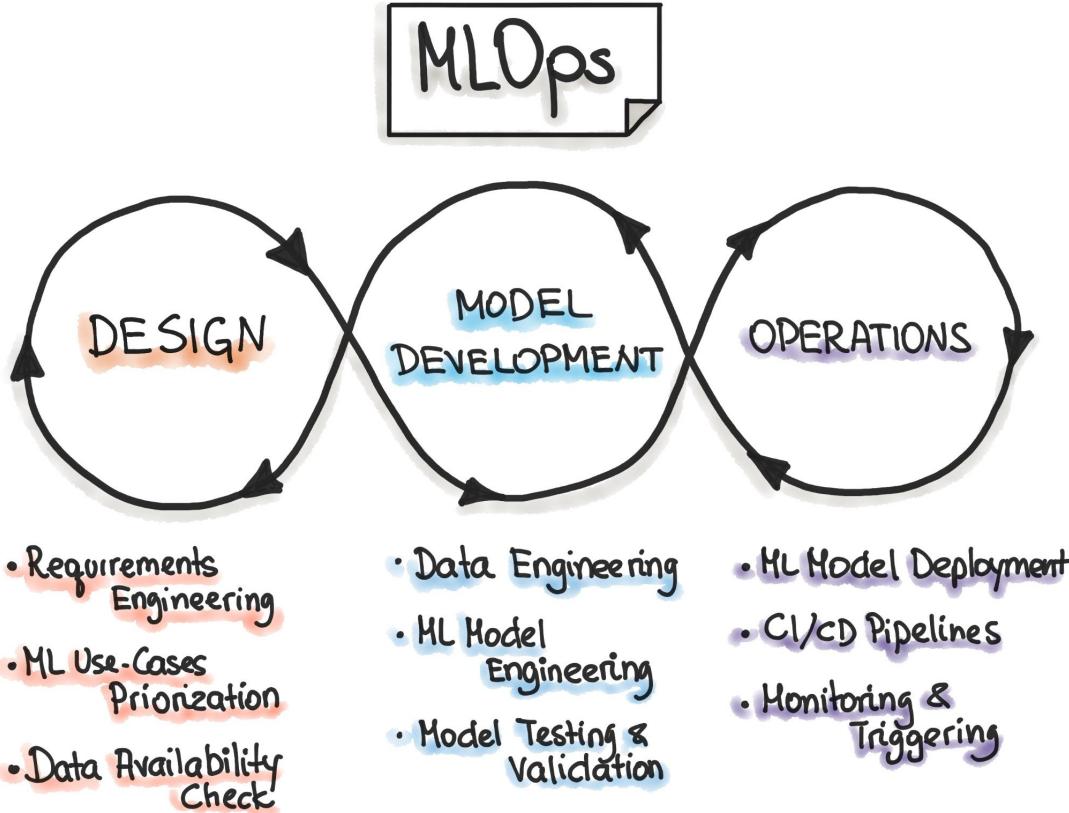


DX Dev eXperience

REX Research & Experimentation eXperience

UX





Objetivos

⬆ Escalabilidad ⬆ Resiliencia ⬇ Coste



Objetivos



Escalabilidad



Resiliencia



Coste

La **evaluación incorrecta** de nuestros modelos puede generar un **alto coste económico** .

Automatizar y versionar el despliegue de modelos acelera el flujo de desarrollo, evitando contratiempos y permitiendo concentrarse en la complejidad intrínseca del problema a resolver. 

Un buen diseño y arquitectura hace que tu software sea mantenible, flexible y confiable ( también en experimentación y MLOps )

Test y Validación del Modelo

- Formato final de puesta en producción

PyTorch

TensorFlow

Keras



Torch Script

SavedModel

ONNX



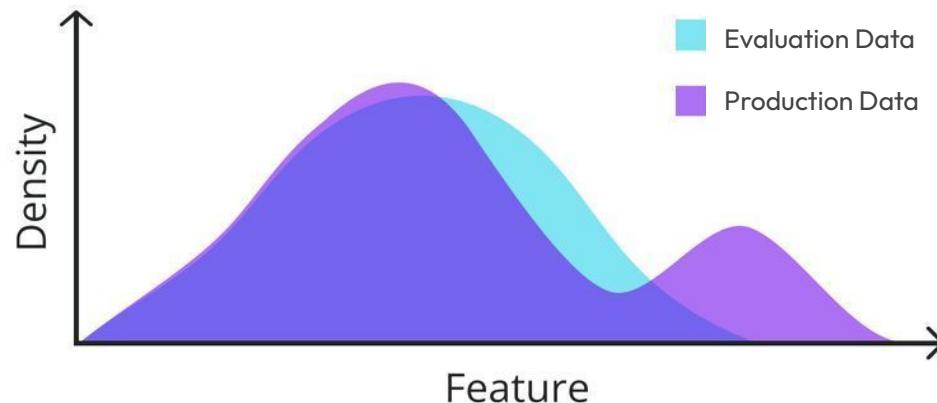
TensorFlow Lite



NVIDIA
TENSORRT

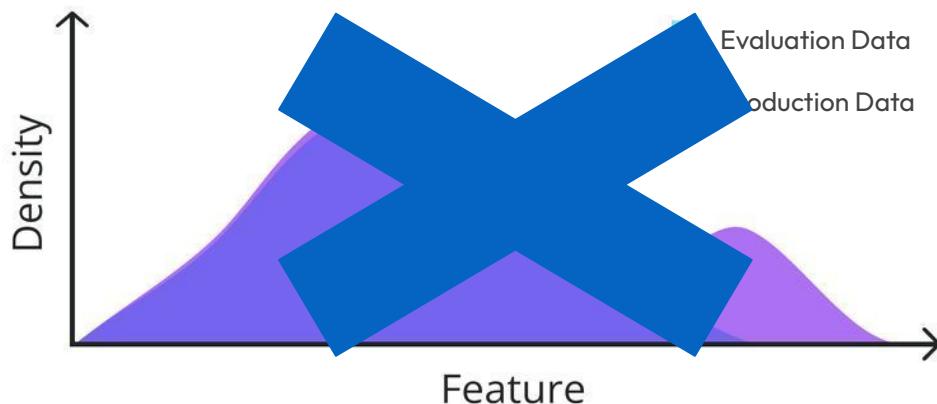
Test y Validación del Modelo

- Formato final de puesta en producción
- Representatividad de producción



Test y Validación del Modelo

- Formato final de puesta en producción
- Representatividad de producción



Evaluación



Producción

Test y Validación del Modelo

- Formato final de puesta en producción
- Representatividad de producción
- Métricas correctas

Problema de Clasificación

Clase	# Muestras Evaluación
	400
	25
	174
	4

Test y Validación del Modelo

- Formato final de puesta en producción
- Representatividad de producción
- Métricas correctas

Accuracy

F1-Score

MSE

Problema de Clasificación

Clase	# Muestras Evaluación
	400
	25
	174
	4

Test y Validación del Modelo

- Formato final de puesta en producción
- Representatividad de producción
- Métricas correctas

Ac~~X~~acy

AI~~X~~

F1-Score

Problema de Clasificación

Clase	# Muestras Evaluación
zebra	400
lion	25
deer	174
puppy	4

Objetivos



Escalabilidad



Resiliencia



Coste

La **evaluación incorrecta** de nuestros modelos puede generar un **alto coste económico** .

Automatizar y versionar el despliegue de modelos acelera el flujo de desarrollo, evitando contratiempos y permitiendo concentrarse en la complejidad intrínseca del problema a resolver. 

Un buen diseño y arquitectura hace que tu software sea mantenible, flexible y confiable ( también en experimentación y MLOps )

Objetivos



Escalabilidad



Resiliencia



Coste

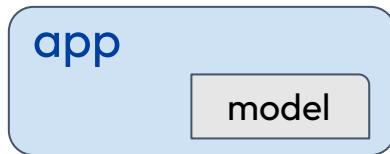
La evaluación incorrecta de nuestros modelos puede generar un alto coste económico 💰.

Automatizar y versionar el despliegue de modelos **acelera el flujo de desarrollo**, evitando contratiempos y permitiendo mantener el **foco en la complejidad intrínseca del problema a resolver**. 

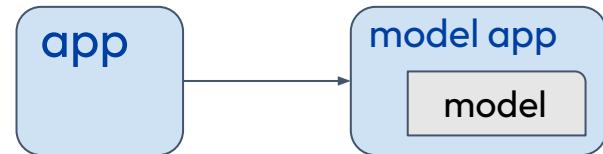
Un buen diseño y arquitectura hace que tus software sea mantenible, flexible y confiable ( también en experimentación y MLOps )

¿Dónde despliego mi modelo?

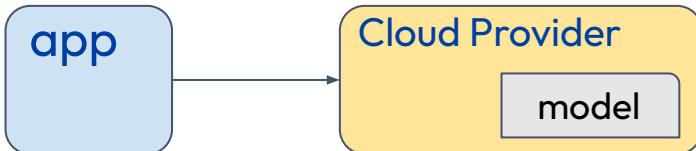
Dentro de mi aplicación



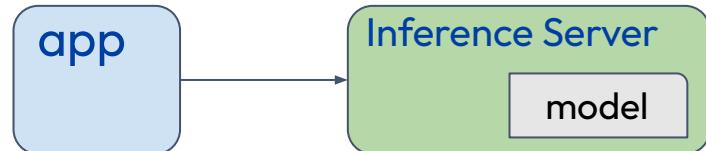
Como una web app independiente



Cloud Provider



Con un Inference Server



¿Dónde despliego mi modelo?

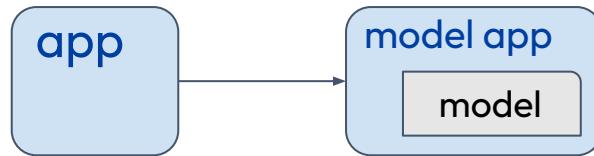
Dentro de mi aplicación



Pros	Cons
<ul style="list-style-type: none">• Facilidad de integración• Complejidad mínima• Versión de modelo asociada a versión de aplicación	<ul style="list-style-type: none">• Tests complejos• Dependencias extra• Versión de modelo asociada a versión de aplicación• Complejidad de escalado

¿Dónde despliego mi modelo?

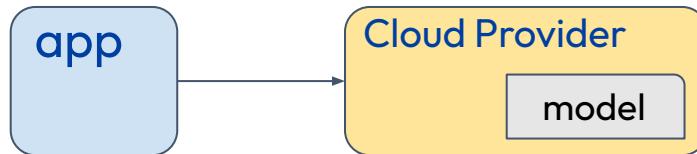
Como un web app independiente



Pros	Cons
<ul style="list-style-type: none">• Flexibilidad• Testeo• Versión de modelo independiente de la aplicación	<ul style="list-style-type: none">• Complejidad de escalado• Optimización de recursos (costes) compleja• Mantenimiento

¿Dónde despliego mi modelo?

Cloud Provider



Amazon SageMaker

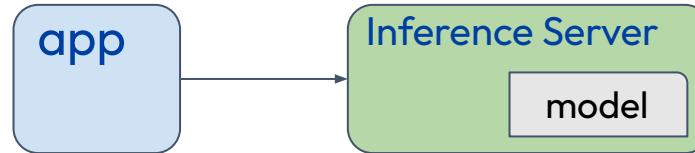


Microsoft
Ignite

Pros	Cons
<ul style="list-style-type: none">• Facilidad de uso• Auto escalado out of the box• Versión de modelo independiente de la aplicación	<ul style="list-style-type: none">• Vendor lock-in• Agilidad y flexibilidad del auto escalado• Dificultad para testear• Difícil de personalizar• Optimización de recursos (costes) compleja

¿Dónde despliego mi modelo?

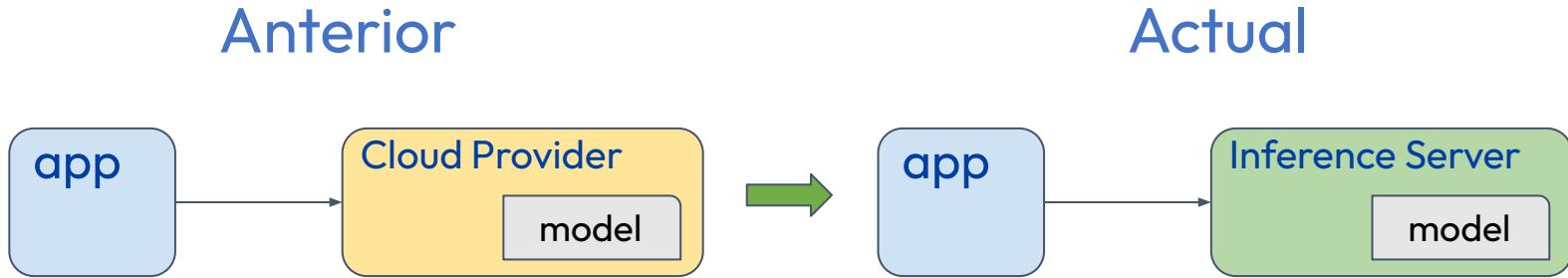
Con un Inference Server



Pros	Cons
<ul style="list-style-type: none">• Flexibilidad• Personalización• Testeo• Optimización de recursos (costes)	<ul style="list-style-type: none">• Testeo• Mantenimiento

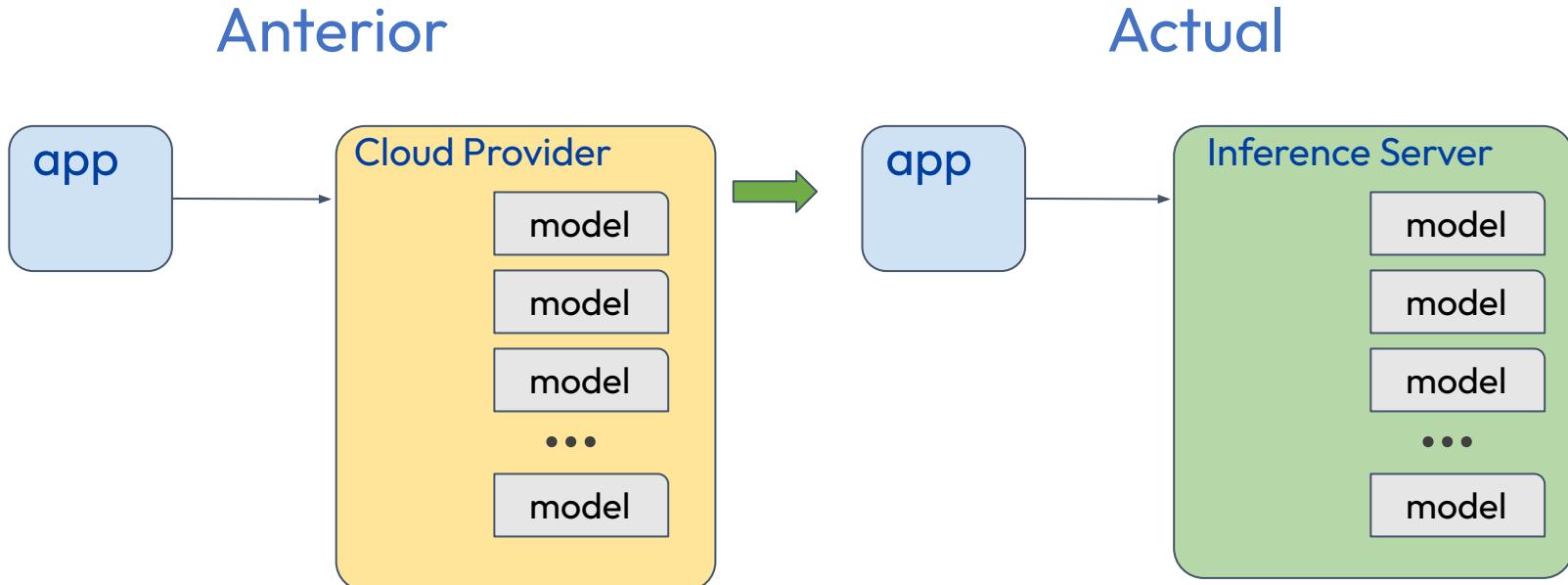
¿Dónde despliego mi modelo?

Nuestro Approach



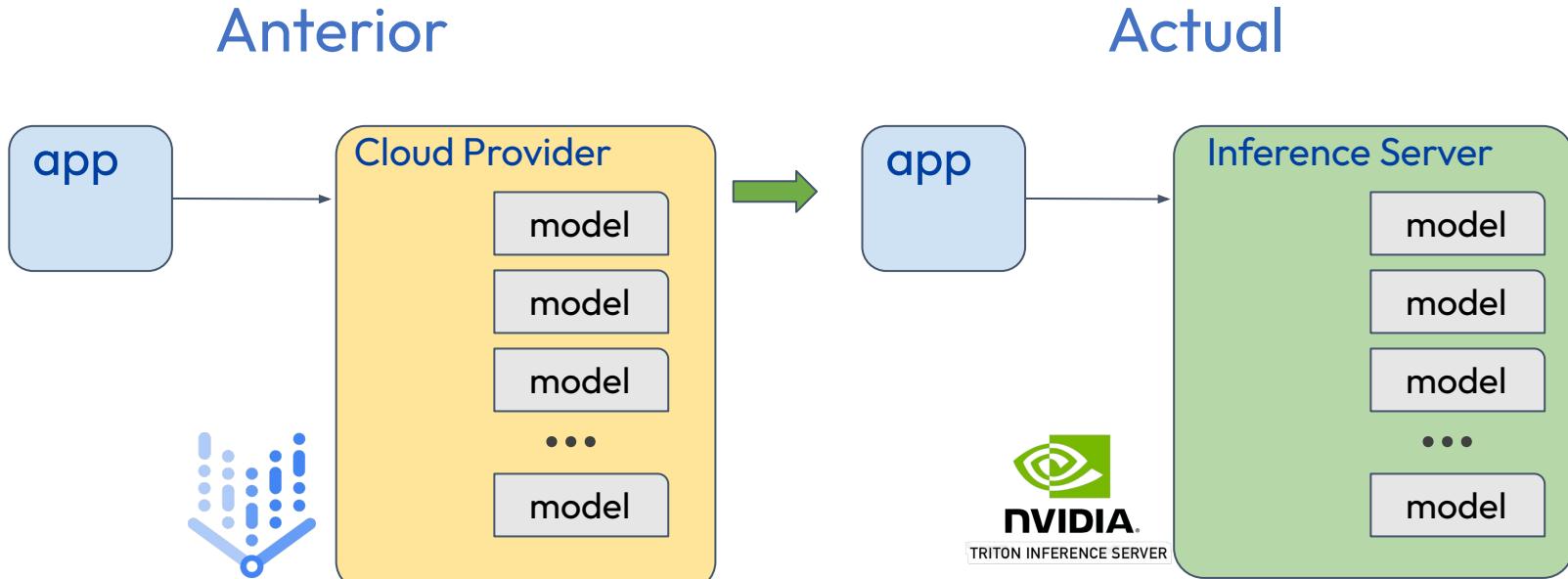
¿Dónde despliego mi modelo?

Nuestro Approach



¿Dónde despliego mi modelo?

Nuestro Approach



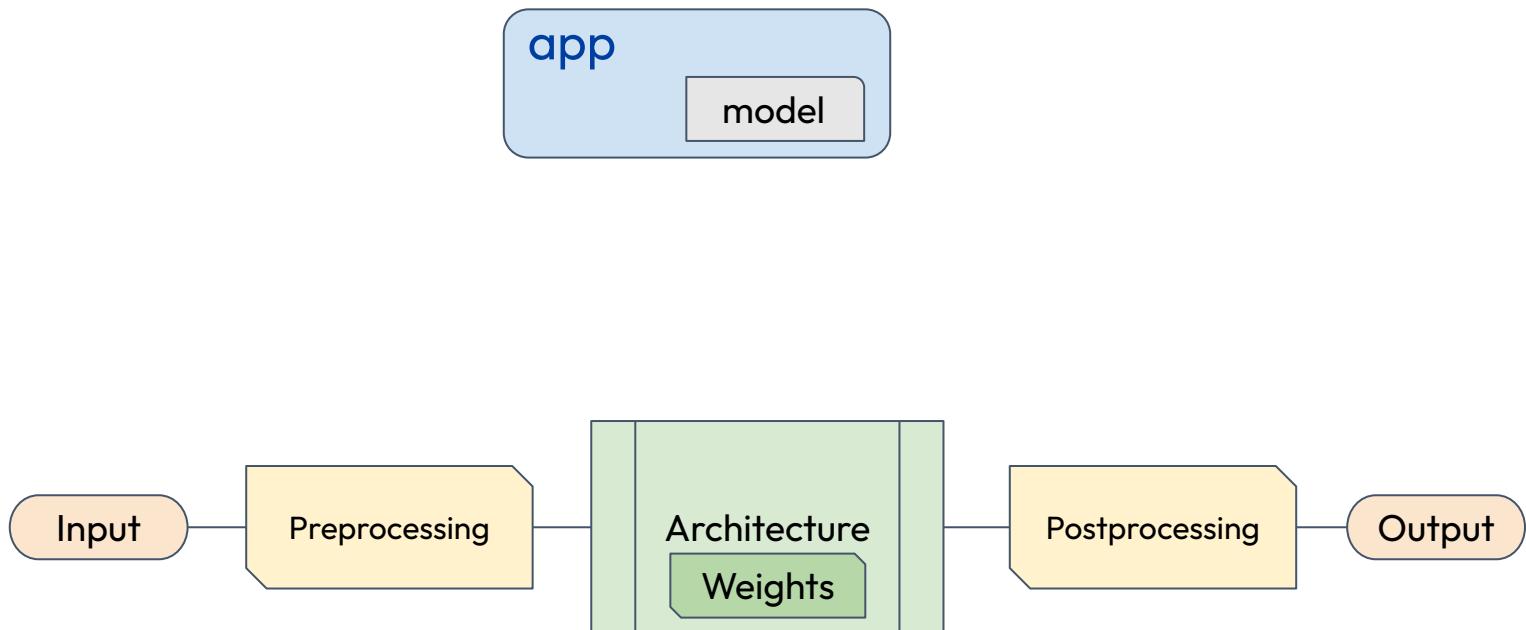
¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)

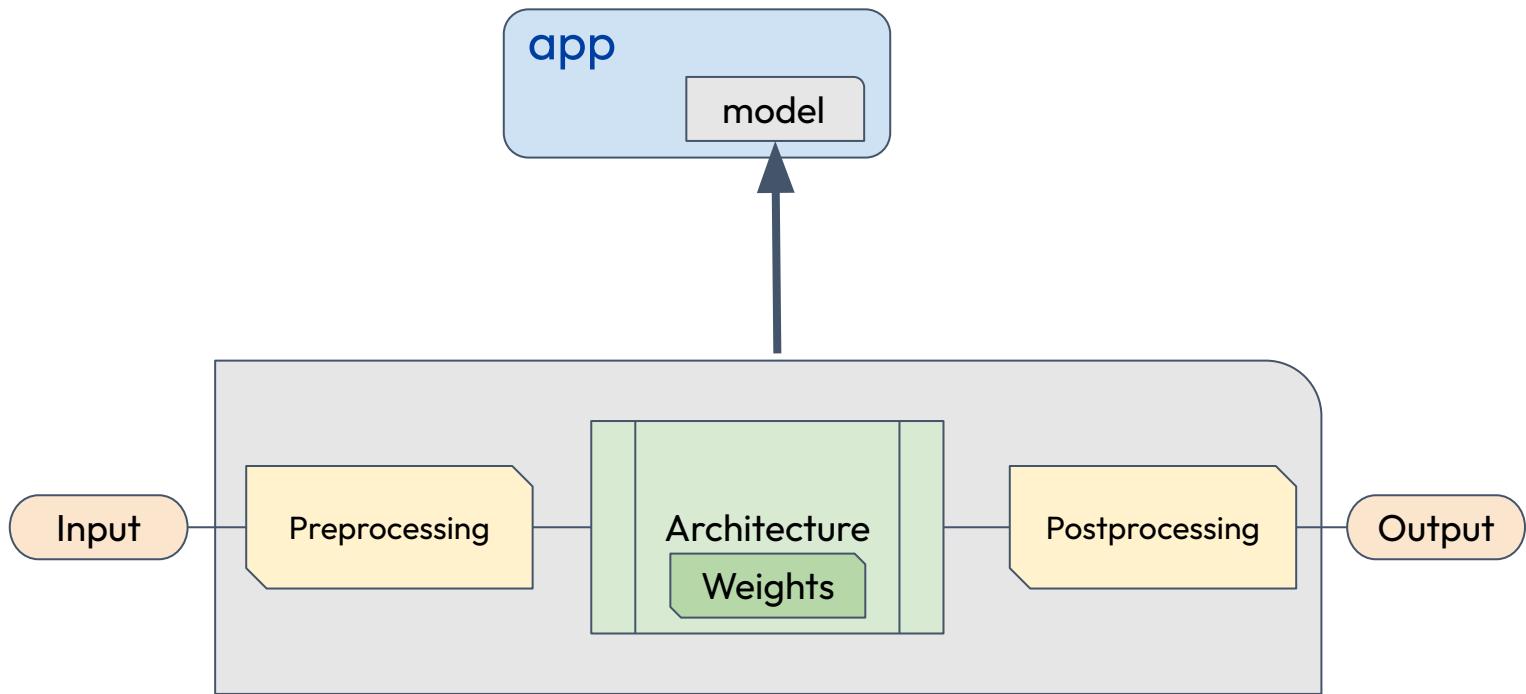
Nuestra definición de modelo



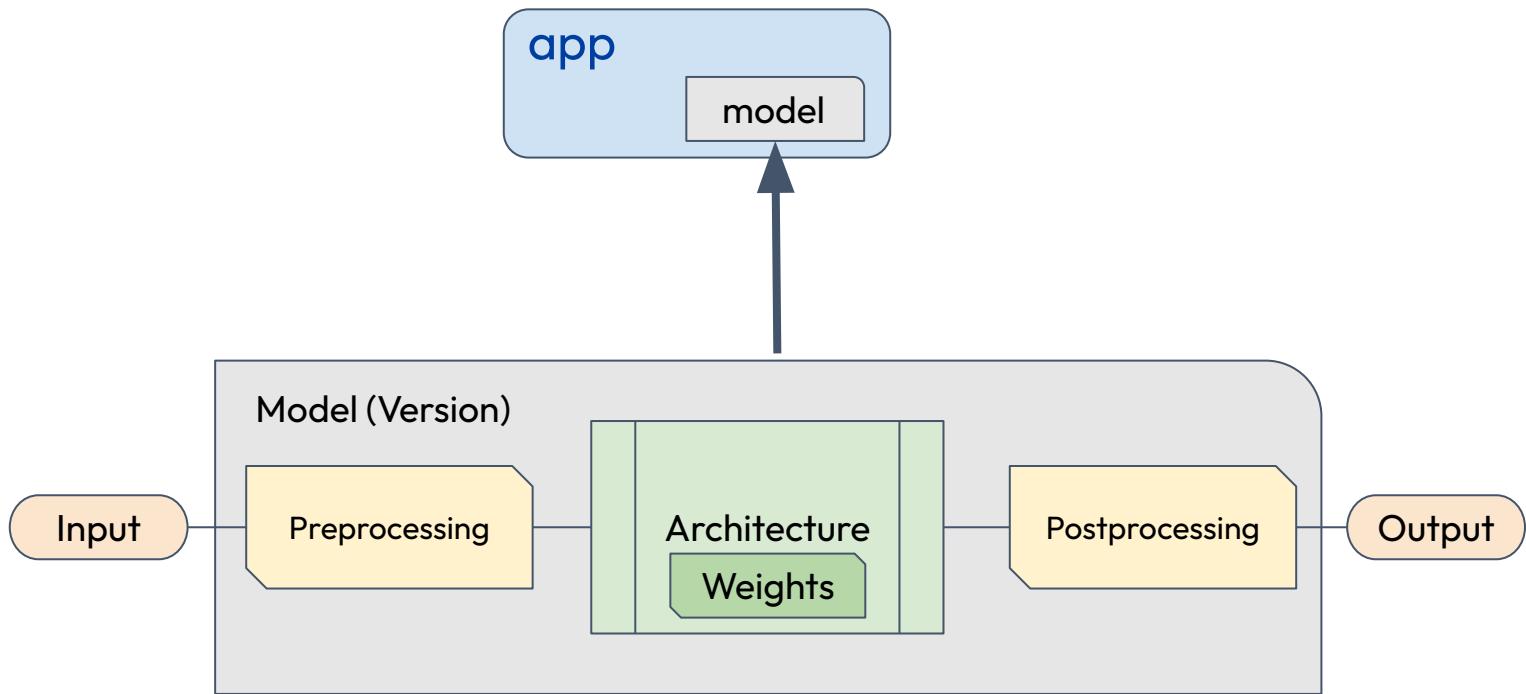
Nuestra definición de modelo



Nuestra definición de modelo



Nuestra definición de modelo



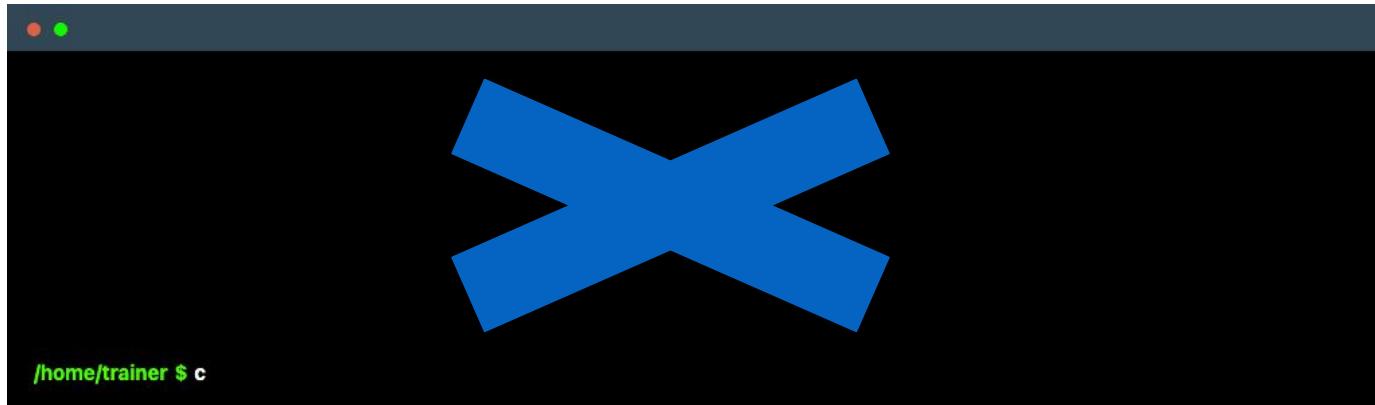
¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)
- Automatizado (sin errores)



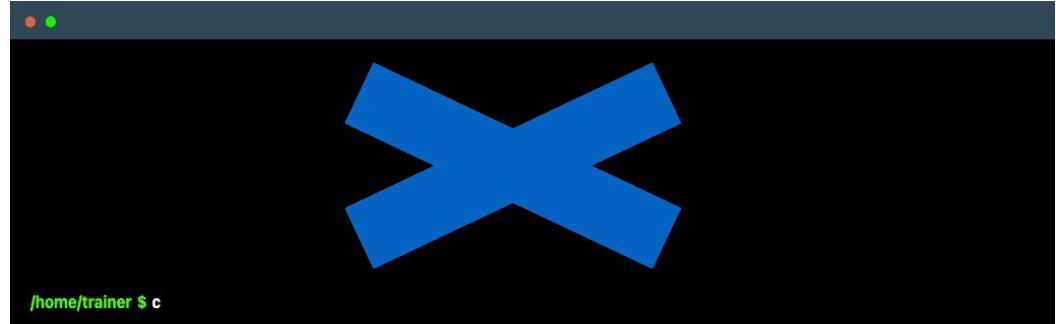
¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)
- Automatizado (sin errores)



¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)
- Automatizado (sin errores)
 - Scripting
 - Github Actions
 - MLOps Tools

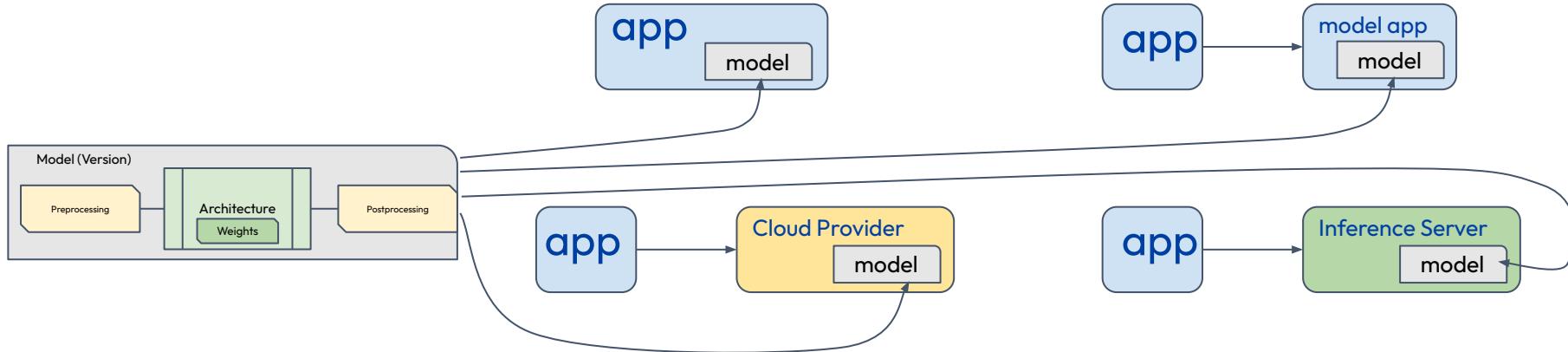


¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)
- Automatizado (sin errores)
- Transparente para la aplicación (agilidad)

¿Cómo despliego mi modelo?

- Versionado del modelo (reproducibilidad)
- Automatizado (sin errores)
- Transparente para la aplicación (agilidad)



Erase una vez ...

Escenario real en el que debido a un problema en la aplicación, ha sido necesario hacer un rollback a una versión anterior

Un modelo que estaba en producción desde hace meses deja de funcionar correctamente

No puede ser... Ayer funcionaba perfectamente y no hemos tocado nada del modelo

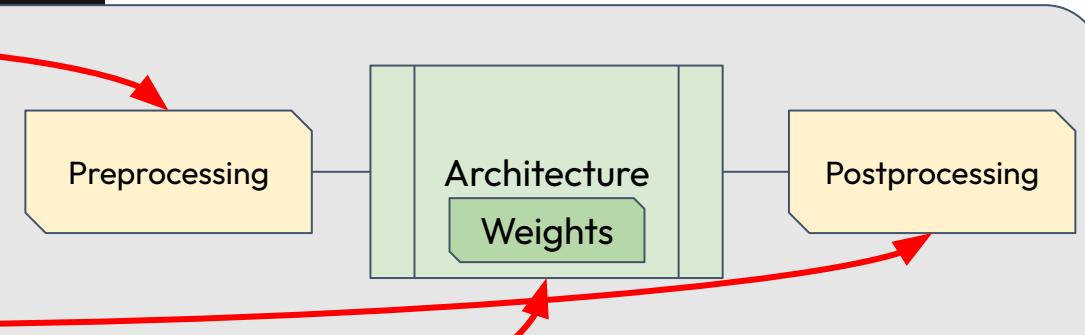
Erase una vez ...

```
class ClassifyController(UseCase):
    @staticmethod
    def build(model_name: str, preprocessor_name: str, postprocessor_name: str ) -> Classifier:
        if preprocessor_name == "FakeVSReal":
            preprocessor = StandardizeInput()
        elif preprocessor_name == "FakeVSRealV1":
            preprocessor = NormalizeInput(mean=0.5, std=1.0)
        elif preprocessor_name == "FakeVSRealV2":
            preprocessor = NormalizeInput(mean=127.5, std=127.5)
        elif preprocessor_name == "FakeVSRealV2a":
            preprocessor = StandardizeInput2()
        else:
            raise ValueError(f"Unknown preprocessor {preprocessor_name}")

        if postprocessor_name == "FakeVSReal":
            postprocessor = FakeVSRealPostprocessor(th=0.234)
        elif postprocessor_name == "FakeVSRealV1":
            postprocessor = FakeVSRealPostprocessor(th=0.766)
        elif postprocessor_name == "FakeVSRealV2":
            postprocessor = FakeVSRealPostprocessor()
        elif postprocessor_name == "FakeVSRealV2a":
            postprocessor = FakeVSRealPostprocessor(invert=True)
        else:
            raise ValueError(f"Unknown postprocessor {postprocessor_name}")

        if model_name == "FakeVSReal":
            model = load_model("resources/models/fake_vs_real/20220417_fake_vs_real_arch_1")
        elif model_name == "FakeVSRealV1":
            model = load_model("resources/models/fake_vs_real/20220422_fake_vs_real_arch_1")
        elif model_name == "FakeVSRealV2":
            model = load_model("resources/models/fake_vs_real/20220422_fake_vs_real_arch_2")
        elif model_name == "FakeVSRealV2a":
            model = load_model("resources/models/fake_vs_real/20220425_fake_vs_real_arch_2")
        else:
            raise ValueError(f"Unknown postprocessor {model_name}")

        return FakeVSRealClassifier(model=model,
                                    preprocessor=preprocessor,
                                    postprocessor=postprocessor)
```



Erase una vez ...

```
class ClassifyController(UseCase):
    @staticmethod
    def build(model_name: str, preprocessor_name: str, postprocessor_name: str ) -> Classifier:
        if preprocessor_name == "FakeVSReal":
            preprocessor = StandardizeInput()
        elif preprocessor_name == "FakeVSRealV1":
            preprocessor = NormalizeInput(mean=0.5, std=1.0)
        elif preprocessor_name == "FakeVSRealV2":
            preprocessor = NormalizeInput(mean=127.5, std=127.5)
        elif preprocessor_name == "FakeVSRealV2a":
            preprocessor = StandardizeInput2()
        else:
            raise ValueError(f"Unknown preprocessor {preprocessor_name}")

        if postprocessor_name == "FakeVSReal":
            postprocessor = FakeVSRealPostprocessor(th=0.234)
        elif postprocessor_name == "FakeVSRealV1":
            postprocessor = FakeVSRealPostprocessor(th=0.766)
        elif postprocessor_name == "FakeVSRealV2":
            postprocessor = FakeVSRealPostprocessor()
        elif postprocessor_name == "FakeVSRealV2a":
            postprocessor = FakeVSRealPostprocessor(invert=True)
        else:
            raise ValueError(f"Unknown postprocessor {postprocessor_name}")

        if model_name == "FakeVSReal":
            model = load_model("resources/models/fake_vs_real/20220417_fake_vs_real_arch_1")
        elif model_name == "FakeVSRealV1":
            model = load_model("resources/models/fake_vs_real/20220422_fake_vs_real_arch_1")
        elif model_name == "FakeVSRealV2":
            model = load_model("resources/models/fake_vs_real/20220422_fake_vs_real_arch_2")
        elif model_name == "FakeVSRealV2a":
            model = load_model("resources/models/fake_vs_real/20220425_fake_vs_real_arch_2")
        else:
            raise ValueError(f"Unknown postprocessor {model_name}")

        return FakeVSRealClassifier(model=model,
                                    preprocessor=preprocessor,
                                    postprocessor=postprocessor)
```



Objetivos



Escalabilidad



Resiliencia



Coste

La evaluación incorrecta de nuestros modelos puede generar un alto coste económico 💰.

Automatizar y versionar el despliegue de modelos **acelera el flujo de desarrollo**, evitando contratiempos y permitiendo mantener el **foco en la complejidad intrínseca del problema a resolver**. 

Un buen diseño y arquitectura hace que tus software sea mantenible, flexible y confiable ( también en experimentación y MLOps ).

Objetivos



Escalabilidad



Resiliencia



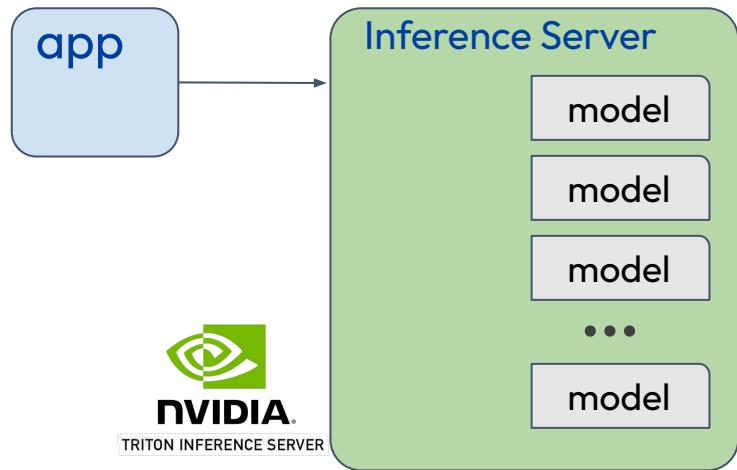
Coste

La evaluación incorrecta de nuestros modelos puede generar un alto coste económico 💰.

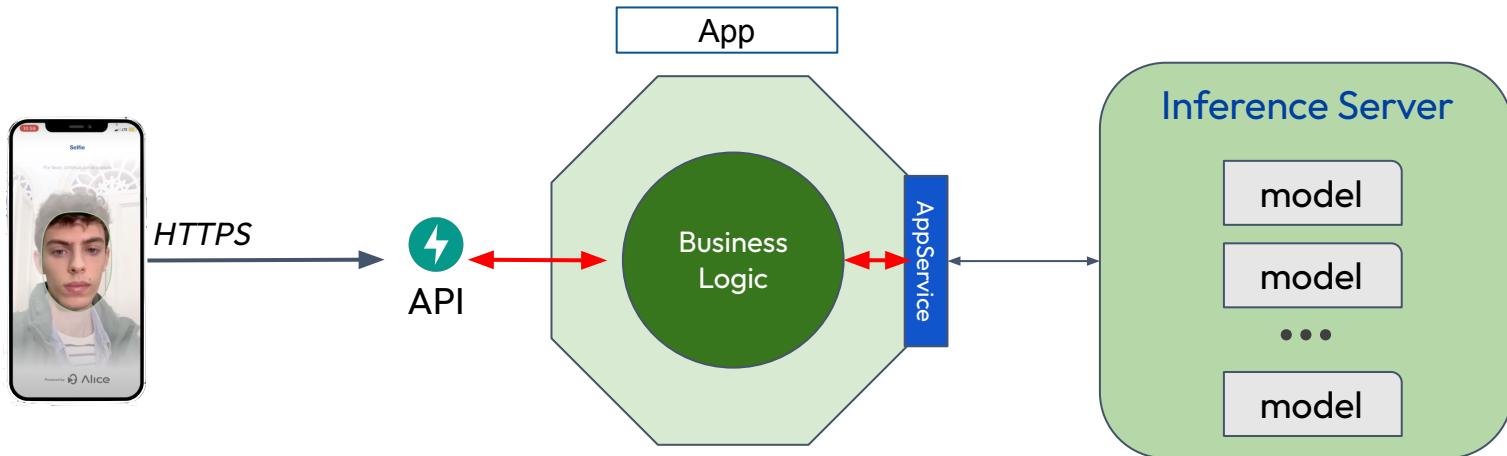
Automatizar y versionar el despliegue de modelos acelera el flujo de desarrollo, evitando contratiempos y permitiendo mantener el foco en la complejidad intrínseca del problema a resolver. 

Un buen diseño y arquitectura hace que tu software sea mantenible, flexible y confiable (🙌 también en **experimentación** y **MLOps** 🙌).

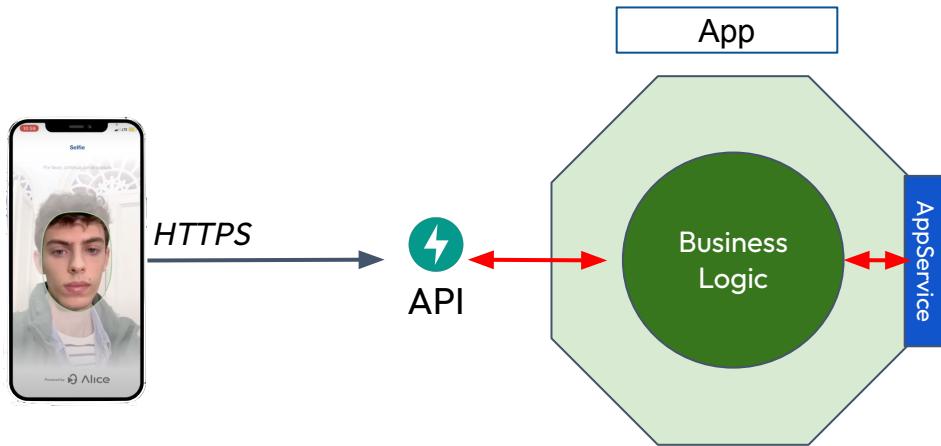
Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



¿Este selfie es fake o es real?



Lógica de negocio

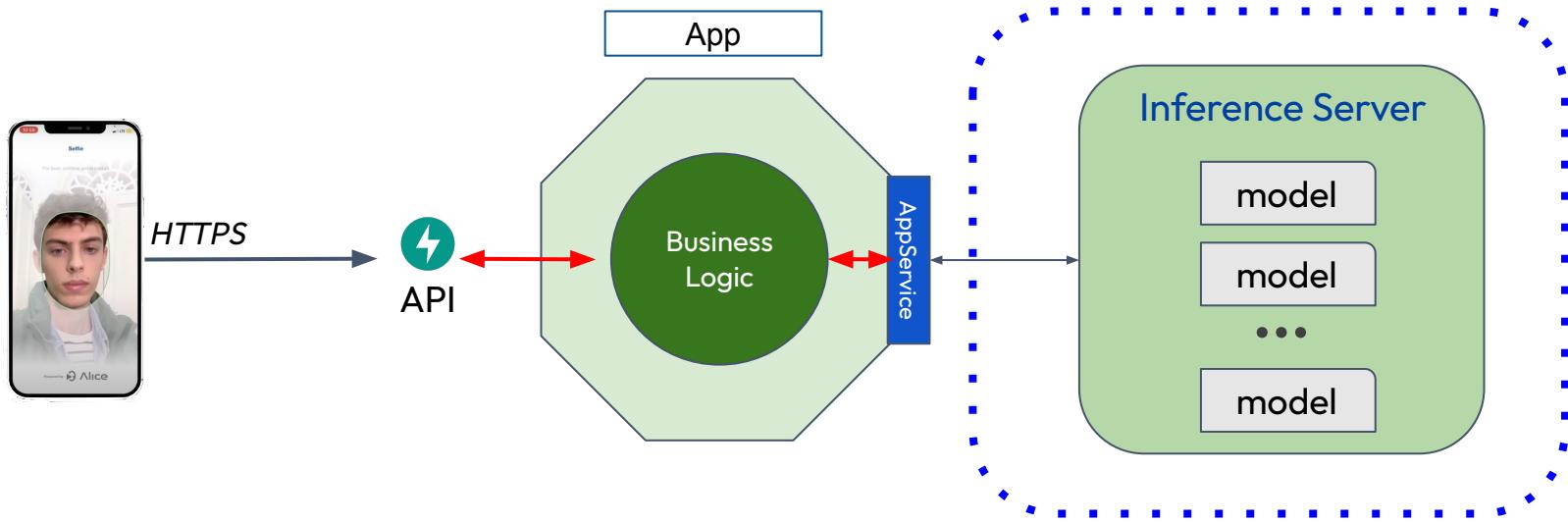
Ejecución de un modelo de clasificación que diferencia entre selfies fake y reales



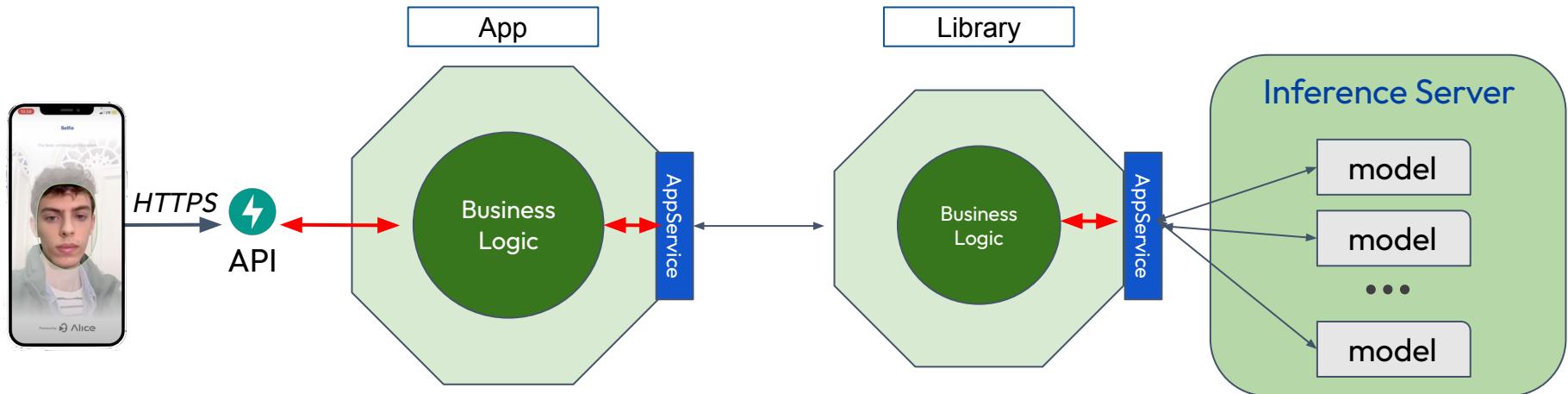
Infraestructura

Inferencia del modelo

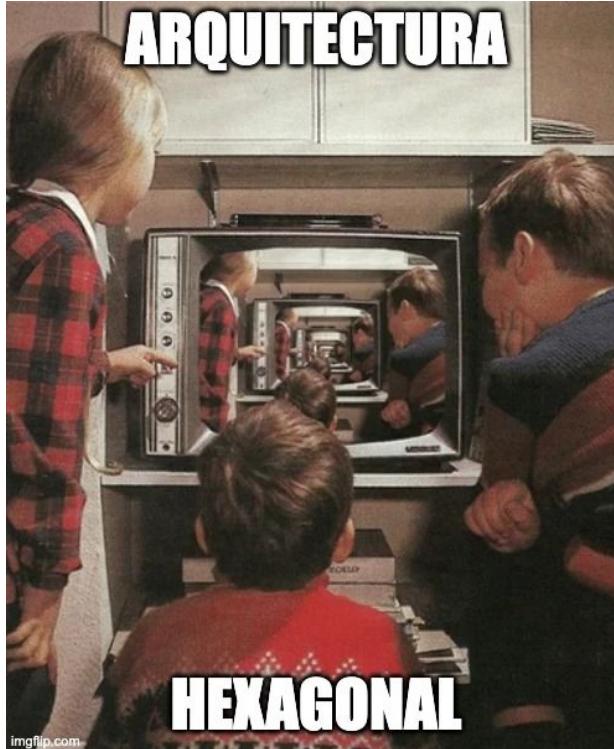
Buen diseño en todas las capas [ml 🎉]



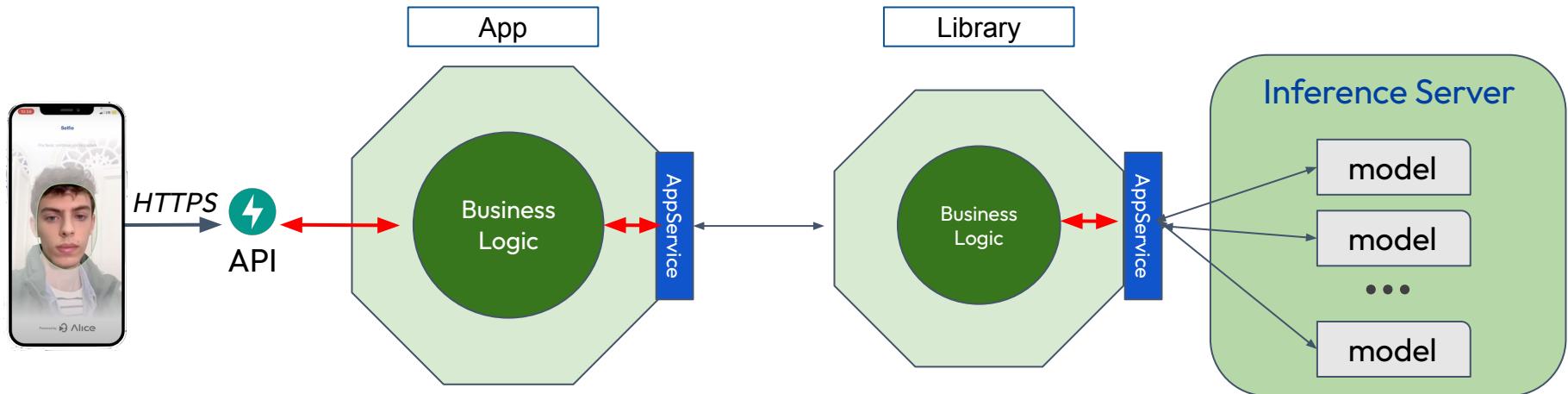
Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]

🎯 Dominio

Inferencia de un modelo

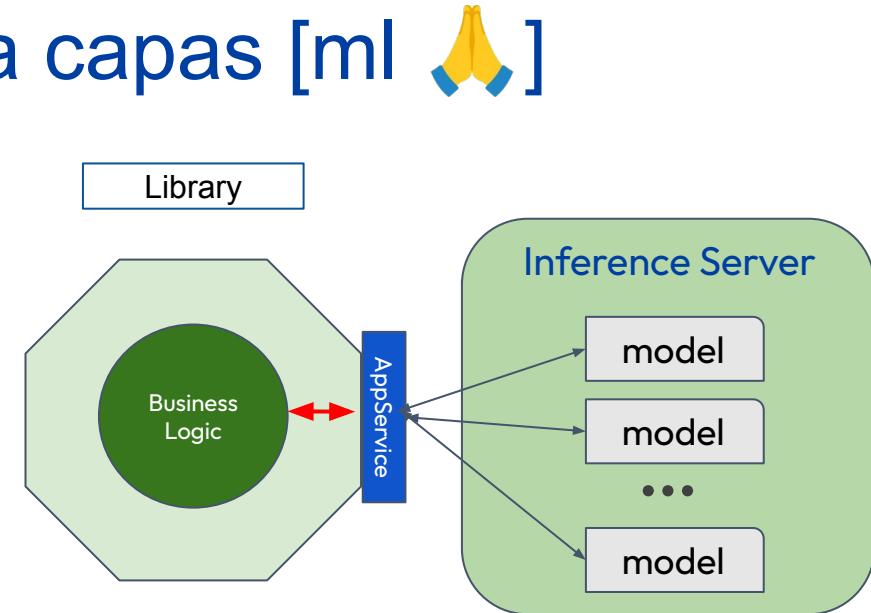
🤝 Lógica de negocio

Tarea de clasificación

Tamaño de entrada (h, w, c)

Normalizado con $\text{mean}=0.5$ y $\text{std}=0.5$

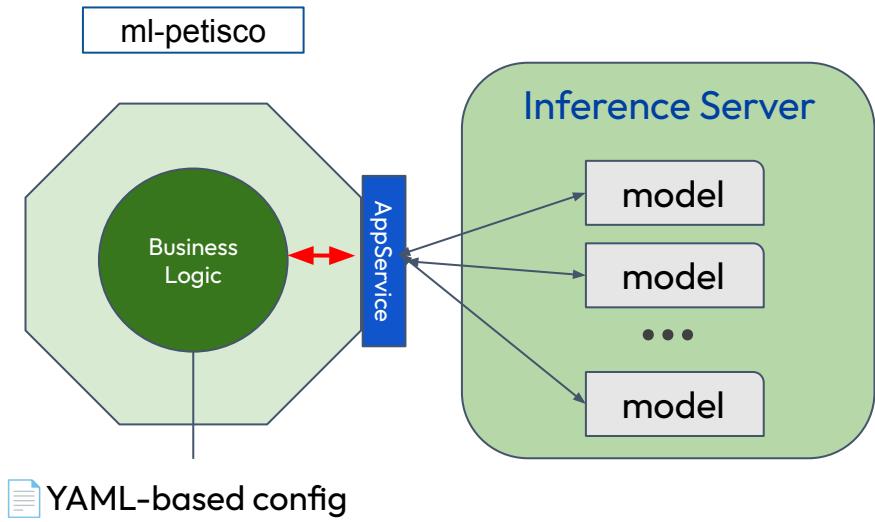
Fake si salida > 0.5 . Si no, es real



🏗 Infraestructura

Ejecuta el modelo X con los recursos Y

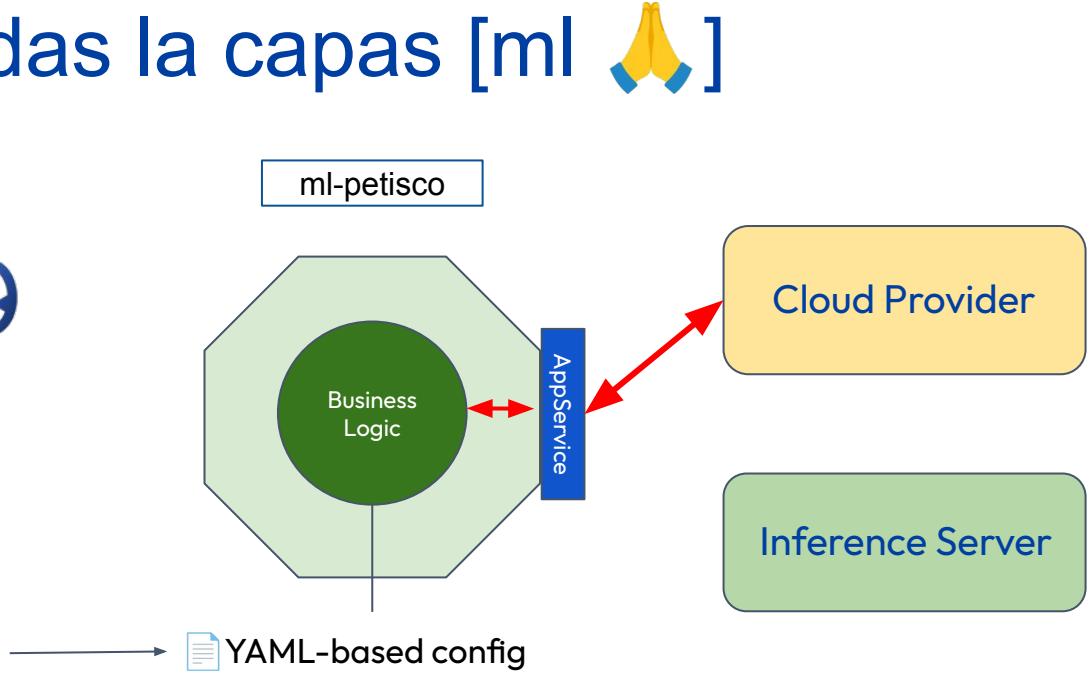
Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



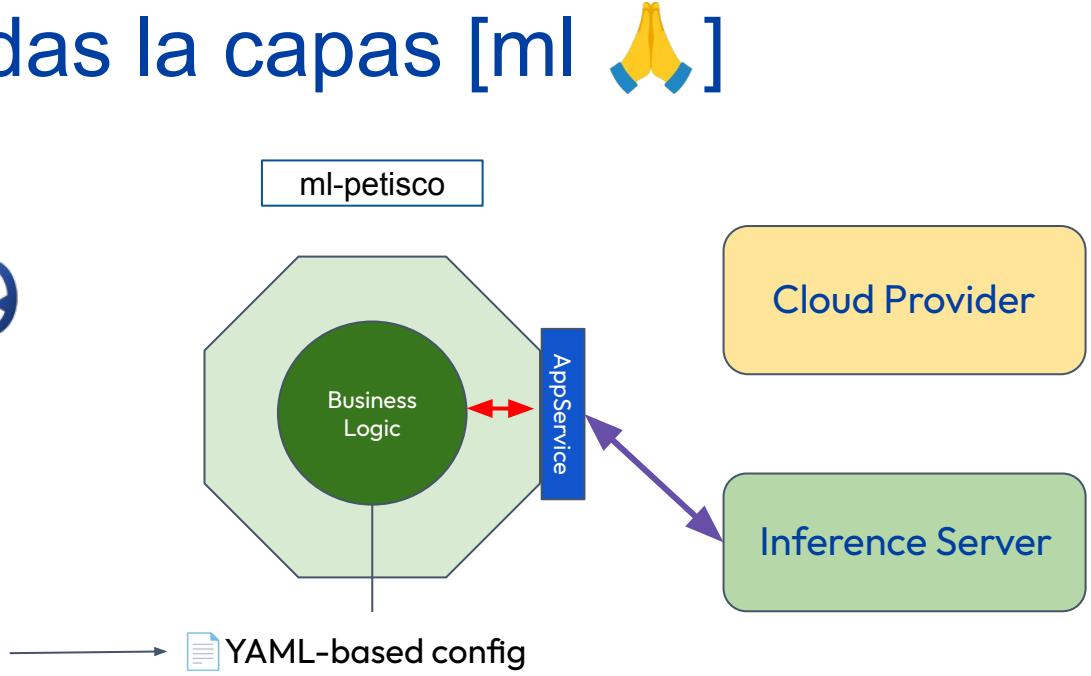
```
- id: selfie_fake_vs_real
  encoder:
    name: BufferMediaImageEncoder
    args:
      image_format: ".png"
      resize_mode: "fix_size"
      size: 224
      resize_type: 1
  model:
    runner:
      type: vertexai
      model: "220784813566623"
    parameters:
      return_with_threshold: True
    postprocess:
      with_softmax: True
  formatter:
    name: ClassificationFormatter
```



Buen diseño en todas las capas [ml 🙏]

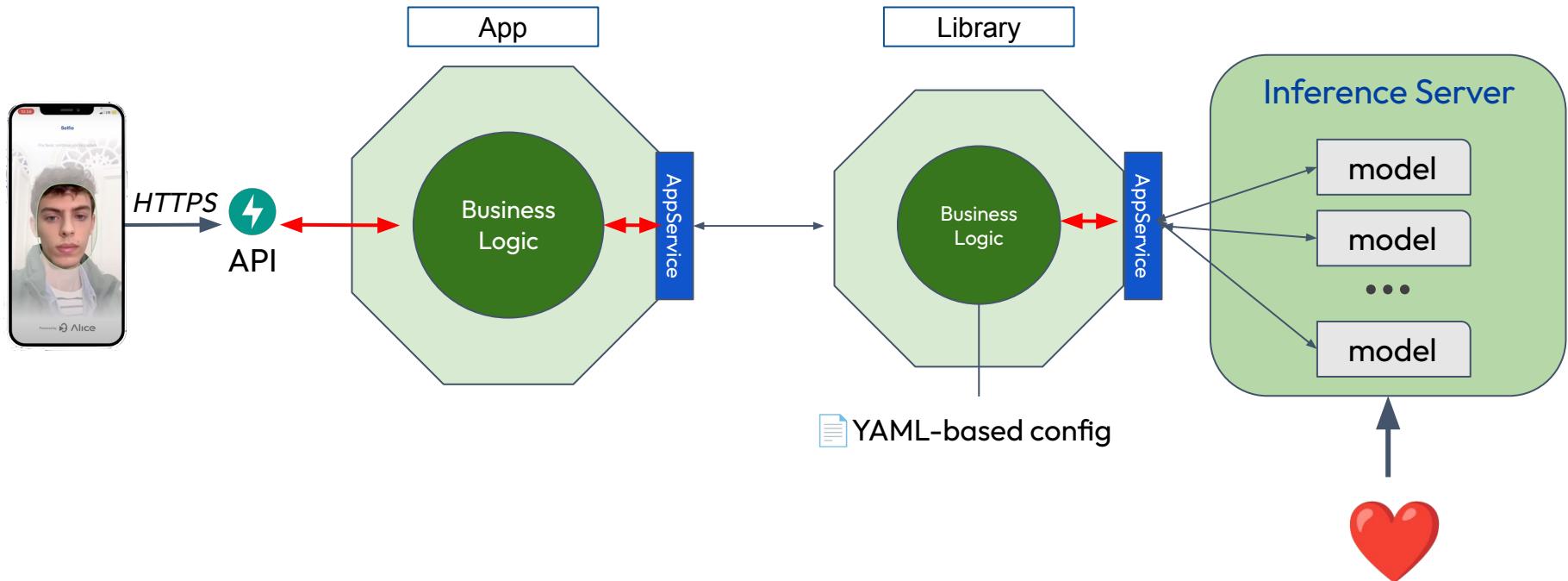


```
- id: selfie_fake_vs_real
  encoder:
    name: BufferMediaImageEncoder
    args:
      image_format: ".png"
      resize_mode: "fix_size"
      size: 224
      resize_type: 1
  model:
    runner:
      type: triton
      model: selfie_fake_vs_real_mbnet
    meta:
      bls_type: classification
      bls_version: "1"
      model_version: "3"
    parameters:
      return_with_threshold: True
      postprocess:
        with_softmax: True
    formatter:
      name: ClassificationFormatter
```

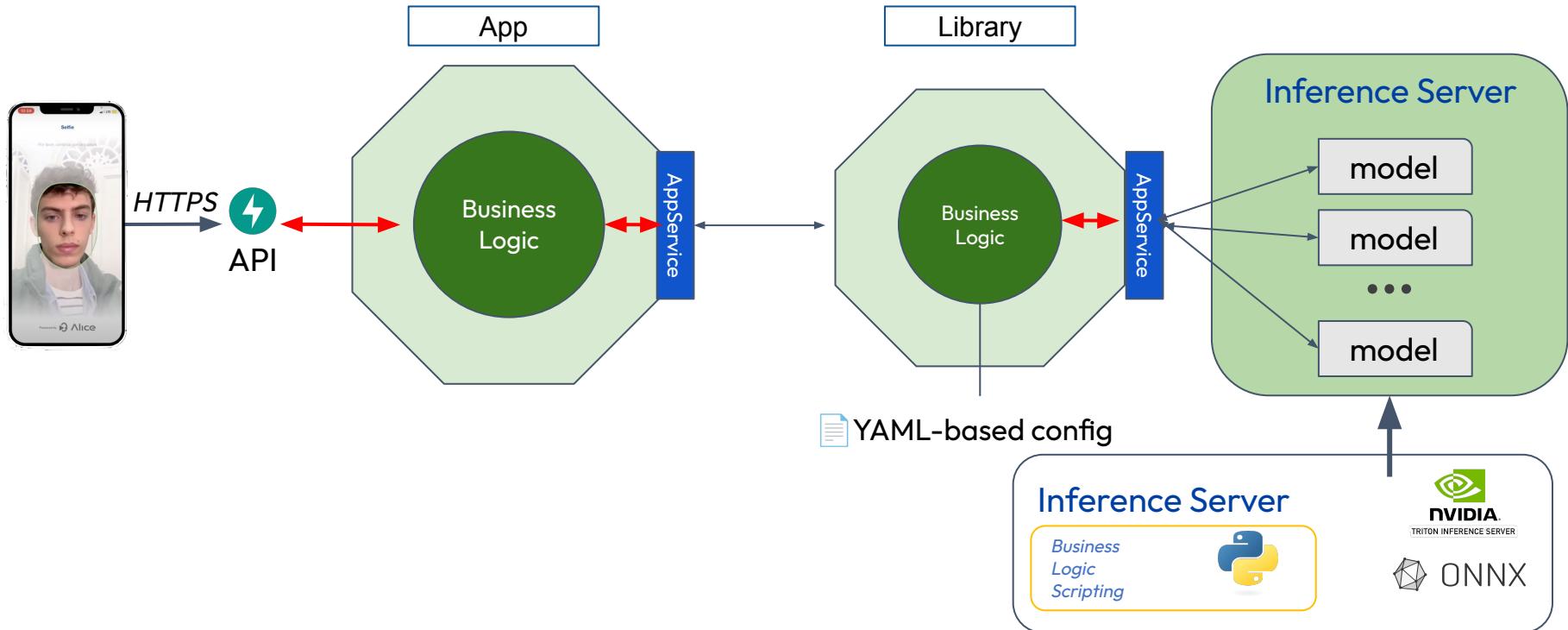


→ **YAML-based config**

Buen diseño en todas las capas [ml 🙏]



Buen diseño en todas las capas [ml 🙏]



Objetivos



Escalabilidad



Resiliencia



Coste

La **evaluación incorrecta** de nuestros modelos puede generar un **alto coste económico** 💰.

Automatizar y versionar el despliegue de modelos **acelera el flujo de desarrollo**, evitando contratiempos y permitiendo mantener el **foco en la complejidad intrínseca del problema a resolver**. 

Un buen diseño y arquitectura hace que tu software sea mantenible, flexible y confiable (👉 tambiéen en **experimentación** y **MLOps** 🙏).

