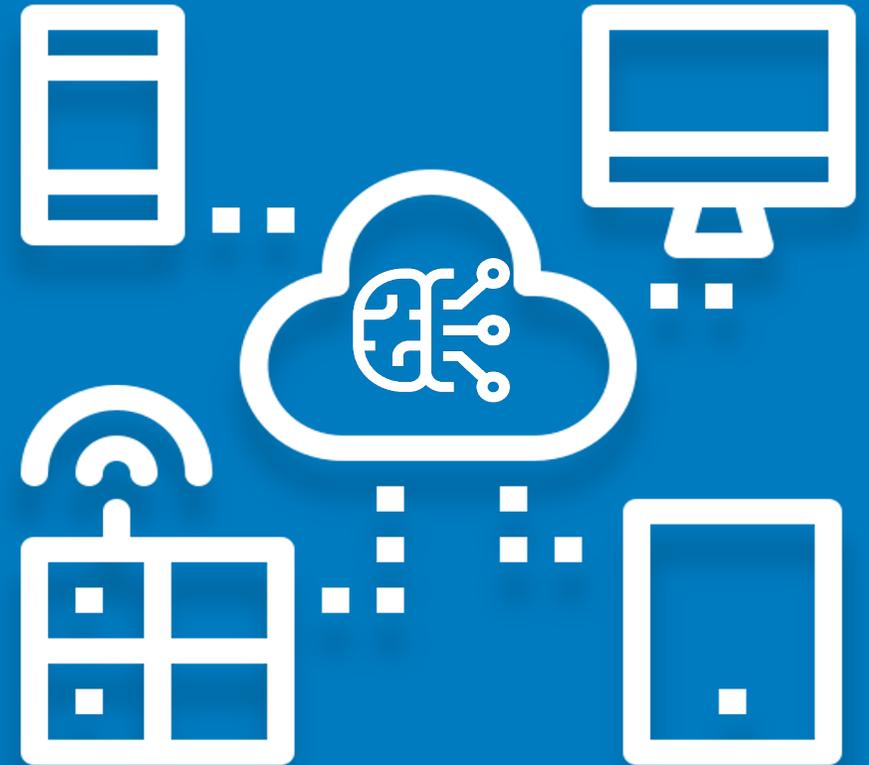


# Inteligencia artificial

*MLOps, la operación de los procesos de Machine Learning*



Cristina Jerez  
UPCnet - Universitat Politècnica de Catalunya

David Garcia  
UPCnet - Universitat Politècnica de Catalunya

**Transformers**

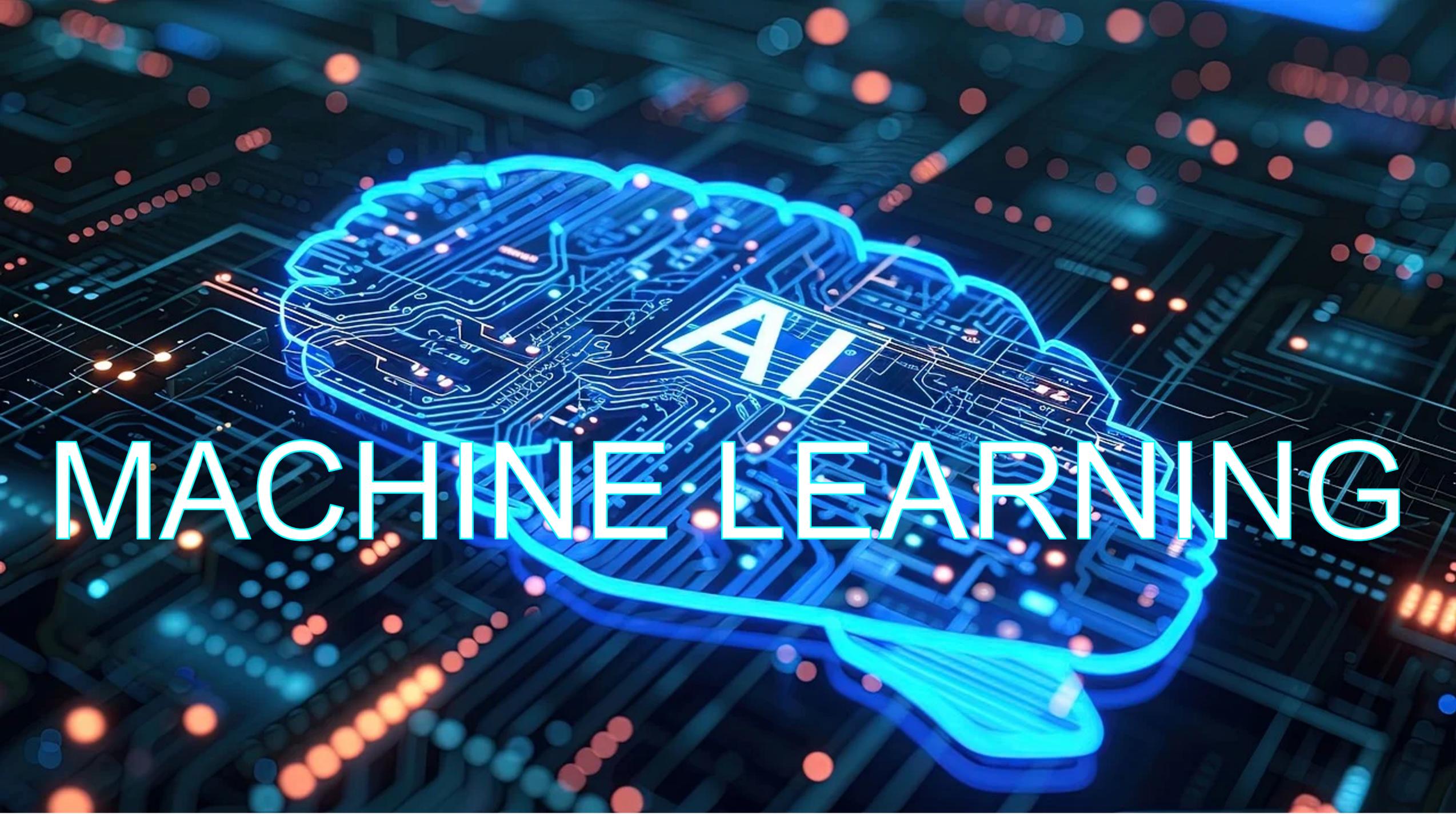
Entrenamiento

Alucinación

Embedding Inferencia

Fine-Tuning





# MACHINE LEARNING

# Data Analysts

Los equipos de data analysts se encargan de recopilar, procesar, analizar e interpretar gran cantidad de datos para poder responder a preguntas o resolver problemas en diversos campos.



## Necesidades

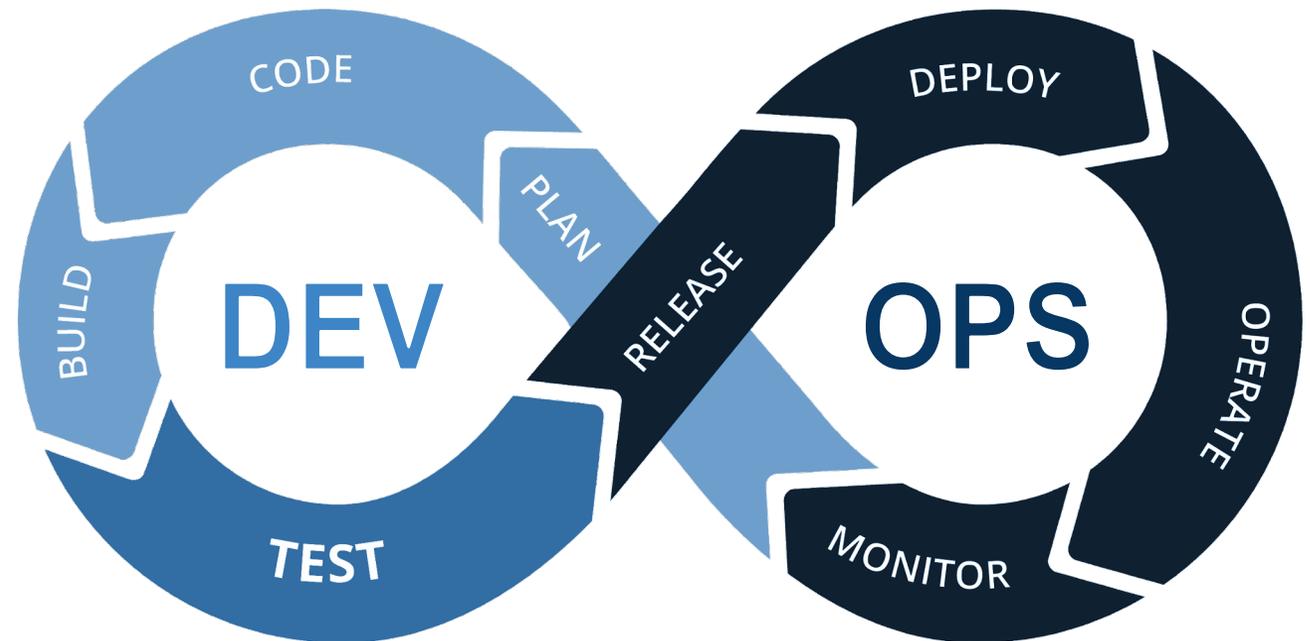
- **Recopilación y limpieza** de los datos con los que trabajar
- **Entornos de Programación** adaptados para trabajar con este ecosistema de datos
- **Facilitar el testeado del modelo** para la evaluación del correcto funcionamiento con **infraestructuras adaptables** a los requisitos de los modelos

# DevOps...

Los equipos de Operaciones están ya en la rueda del funcionamiento de DevOps, trabajando con equipos de desarrollo generando un flujo de confianza

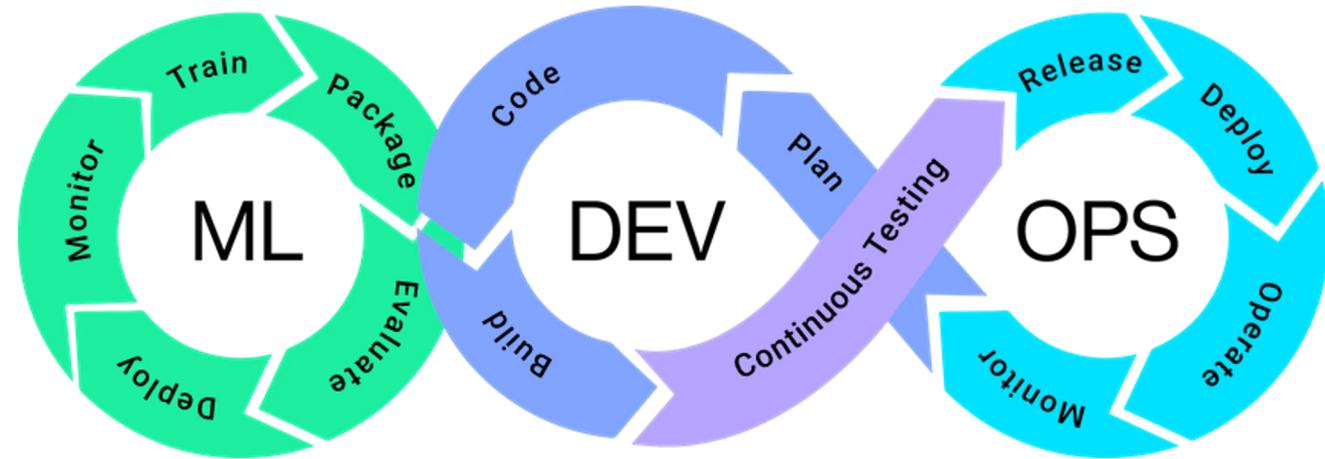
## Nuevos retos

- **Data Analysts...** nuevas necesidades
- Entornos de **desarrollo** ... con nuevas **tecnologías**
- Entornos **productivos** ... con nuevas **capacidades**: GPU, accesibilidad...
- **Operación** de estos entornos ... con nuevas preocupaciones vinculadas al **uso y seguridad**
- **Monitorización** de estos entornos ... con **nuevos procesos**



# MLOps

Debemos analizar los sistemas de AI desde una perspectiva de Operaciones de aprendizaje automático (MLOps) para aprender y definir un marco de común en el flujo de entrenamiento, desarrollo y operación para crear crear sistemas de aprendizaje automático robustos, escalables y seguros.



# Necesidades vs. Problema vs. Solucion

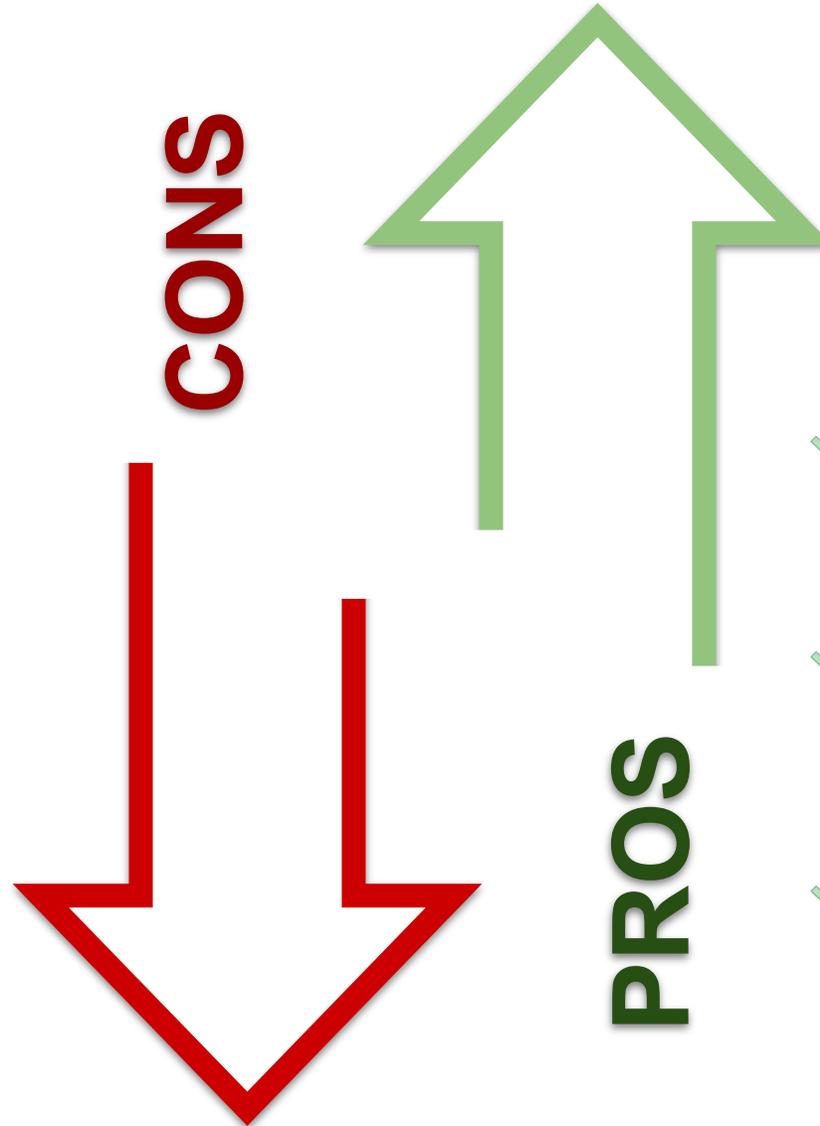




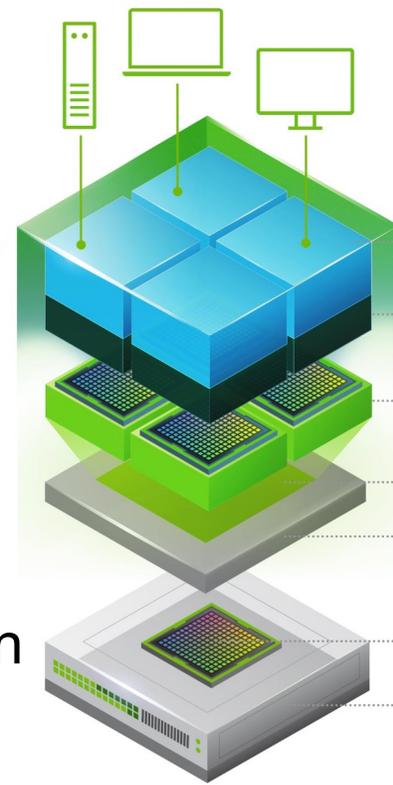
# Arquitectura ML

# Arquitecturas. Primera aproximación

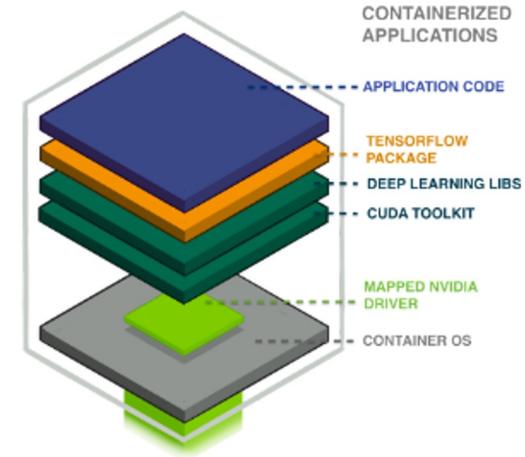
- ✗ Rigidez elección del hardware
- ✗ Gestión cambio poco práctica
- ✗ No escalable
- ✗ Difícil tener más de 1 modelo



- ✓ Implementación más sencilla
- ✓ Adquisición de conocimiento
- ✓ Facilidad para entrenar y hacer inferencia en un mismo entorno



# Arquitecturas. Contenedores!



✗ Complejidad para encontrar documentación y requisitos software de los modelos

✗ Imágenes muy pesadas (~7-8Gb)

✗ Complejidad para hacer pruebas en local por requisitos hardware

**CONS**

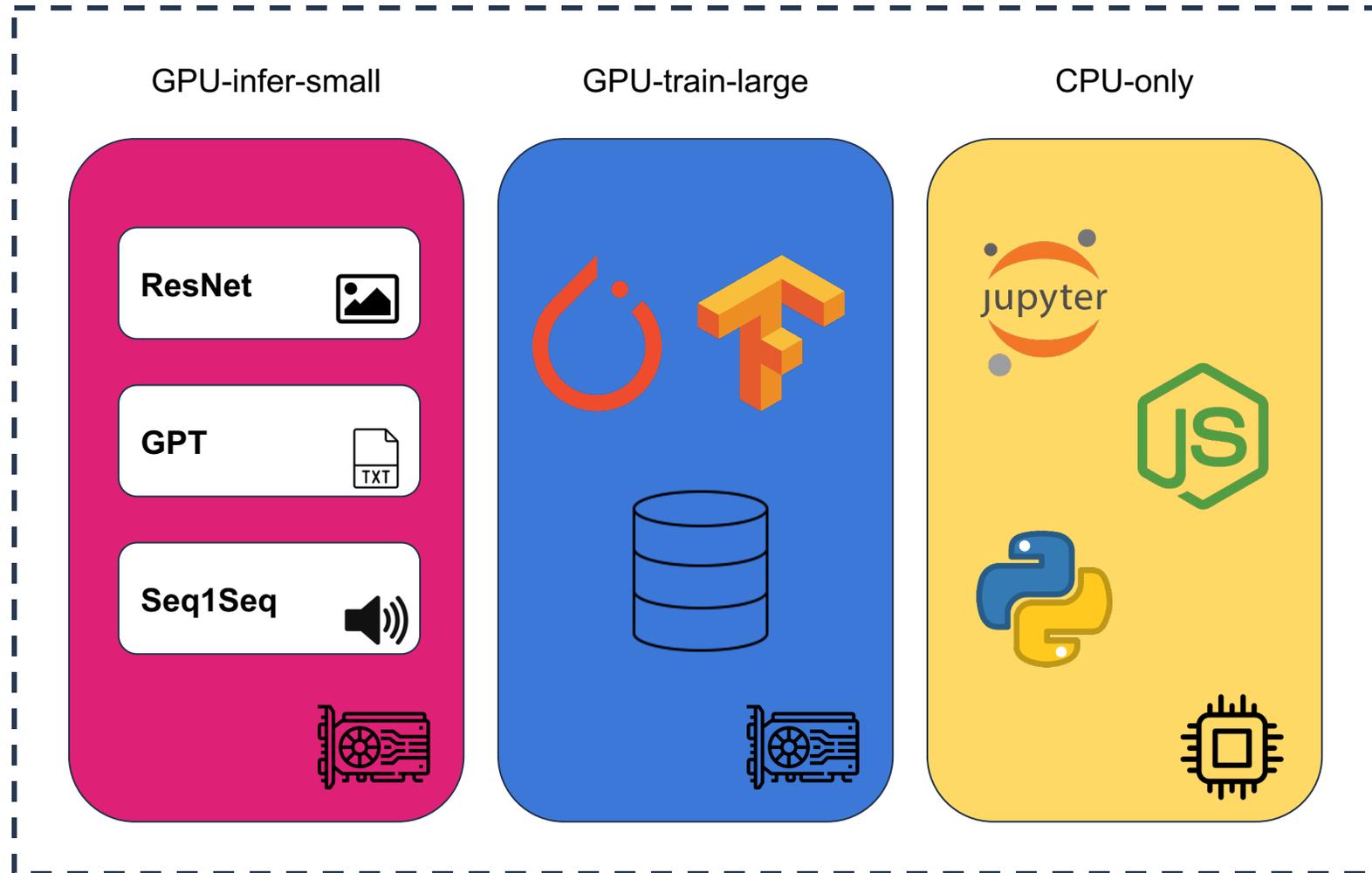
- ✓ Escalabilidad
- ✓ Disponibilidad de imágenes de contenedores diferentes
- ✓ Flexibilidad de requisitos hardware y entornos de ejecución

**PROS**



# Arquitecturas. Entornos productivos

Worker Nodes



# AWS IA Services

AWS dispone de diversos servicios que nos pueden dotar de capacidades de IA

## SAGE MAKER STUDIO



Nos permite desarrollar modelos de forma ágil y sencilla.

Nos ofrece una suite de herramientas que permite ejecutar todos los pasos de creación y entrenamiento de un modelo.

## SAGEMAKER



Nos ofrece una plataforma en la que ejecutar nuestros modelos para hacer predicciones.

Es la base para disponer de un entorno gestionado para operar modelos de inteligencia artificial.

## JUMPSTART



Este servicio nos permite ejecutar modelos fundacionales para realizar pruebas o hacer fine-tuning.

AWS ofrece un catálogo de modelos ya preparados para desplegarse de forma sencilla en SageMaker.

## BEDROCK

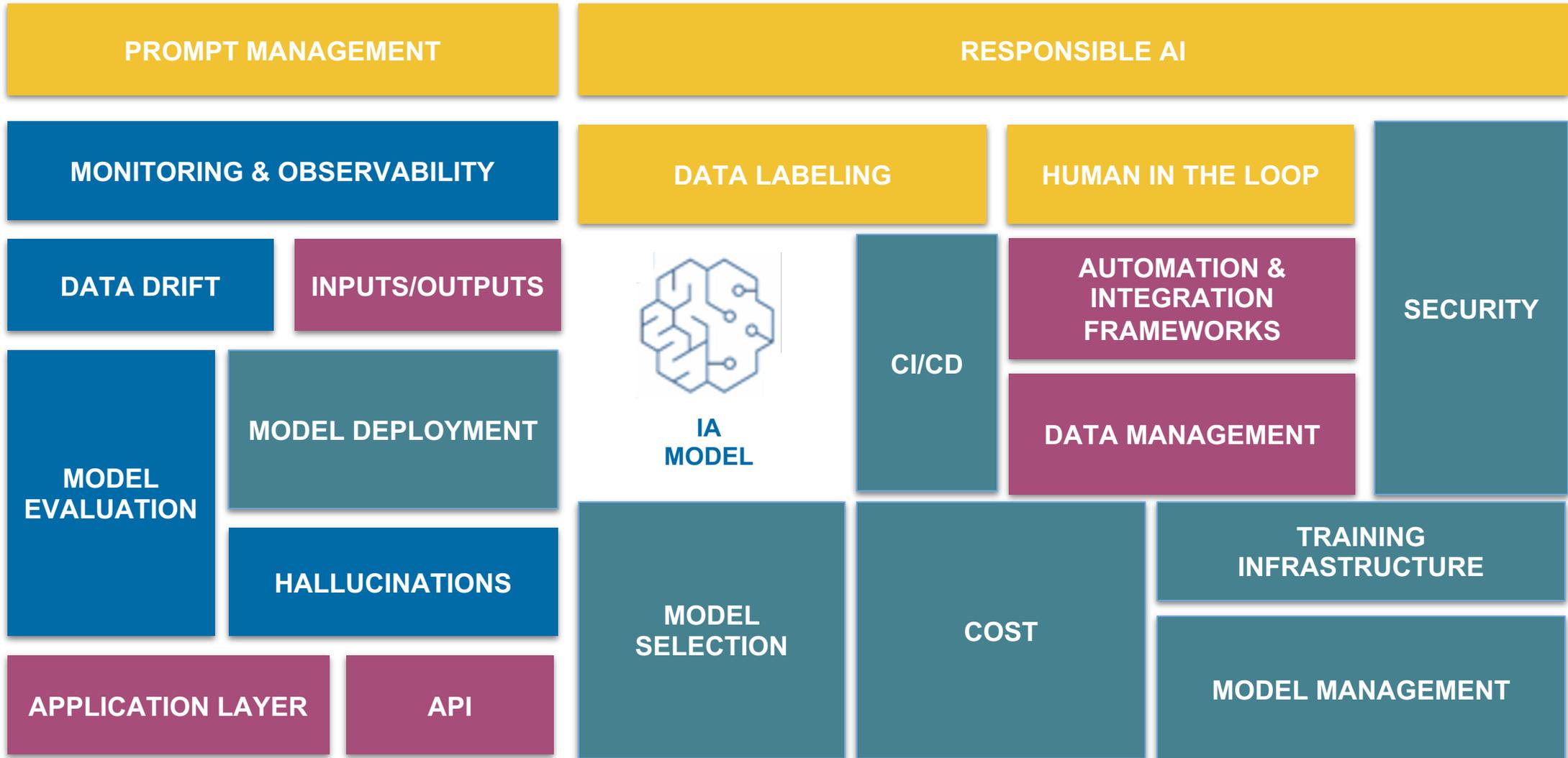


El servicio serverless de SageMaker.

Nos permite disponer de varios modelos fundacionales siempre disponibles y sólo pagaremos por los tokens de entrada y salida.



# ML Services



# Conclusiones y Recomendaciones





Cristina Jerez Alvarez  
[cristina.jerez-alvarez@upcnet.es](mailto:cristina.jerez-alvarez@upcnet.es)



David Garcia Domenech  
[david.garcia.domenech@upcnet.es](mailto:david.garcia.domenech@upcnet.es)