

# De 0 a héroe con MLOps

#**pivo** and **code** 25

Juan Isern Ghosn | Innovasur

# juan@isern:/ \$ whoami

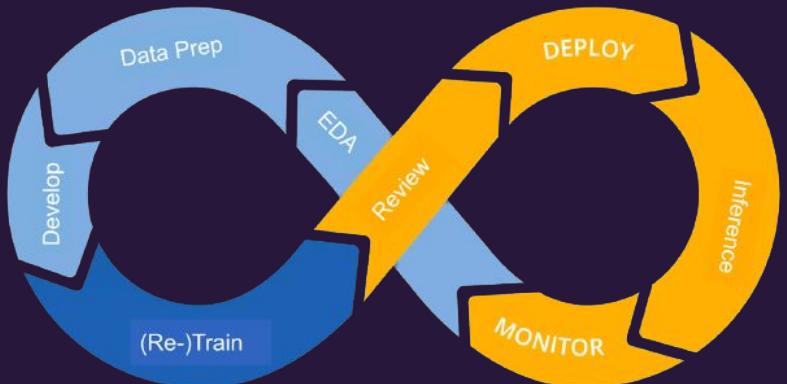
- Graduado en Ing. Informática + ADE
- Máster en Ciencia de Datos e Ingeniería de Computadores
- Doctorando en sistemas bioinspirados (UGR)
- Experto en visión artificial
- Optimización de modelos ML para sistemas de bajas prestaciones
- IA Manager en Innovasur



# MLOps

La combinación de prácticas de DevOps y aprendizaje automático (Machine Learning ~ ML). Facilita la colaboración y comunicación entre desarrolladores y equipos de operaciones.

1. Datos
2. Entrenamiento / Validación
3. Despliegue





Amazon SageMaker

gradient<sup>o</sup>  
by Paperspace

FLOYD

DOMINO  
DATA LAB

“All-in-one”



floure  
eight  
scale  
Aquarium  
Labeling



DAGSTER  
Processing

snowflake<sup>®</sup> databricks  
Data Lake / Warehouse

S3 Parquet  
Sources

TensorFlow fast.ai  
PYTORCH RAY K  
Frameworks & Distributed Training

Determined AI  
slurm workload manager  
Resource Management

Lambda CoreWeave  
Compute

SIGOPT Determined AI  
Weights & Biases tune Hyperparameter Tuning

Weights & Biases comet  
TensorBoard Neptune Machine Learning Lab  
mlflow Experiment Management

Jupyter Streamlit  
git  
Software Engineering

TACTON  
Feature Store

fiddler  
Monitoring

NVIDIA TensorRT  
TensorFlowLite  
ONNX  
Edge

SELDON  
MLflow  
ONNX  
Web

Buildkite great\_expectations  
CI / Testing

or or

Data

Training/Evaluation

Deployment

# ¿Qué vamos a ver?

Conjunto de herramientas básicas (Toolkit) a incluir en nuestros desarrollos de ML: versionado, gestión de los datos, etiquetado, etc.



01

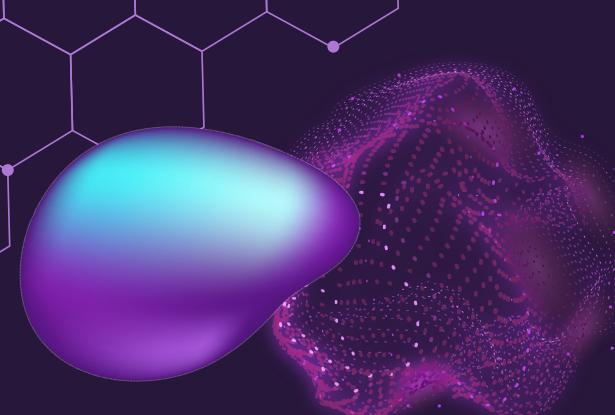
# LOS DATOS



x



+



x

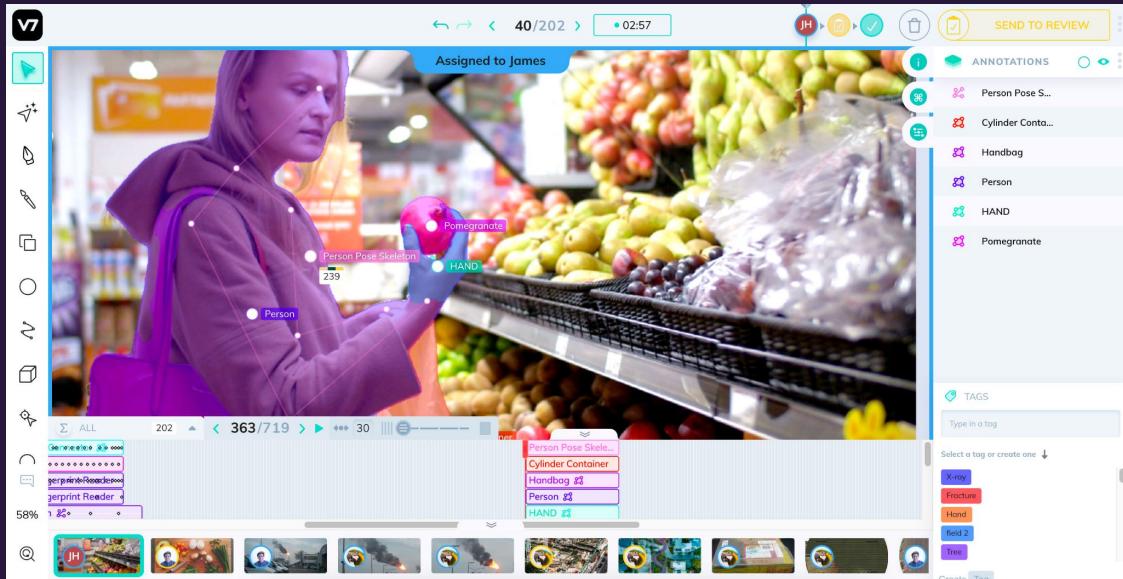


# MLOps en los Datos

Garantía del rendimiento de los modelos ML respecto al buen hacer en la ingesta, transformación y almacenamiento de los datos.

# 01. Datos: Etiquetado

El etiquetado de datos implica asignar etiquetas o etiquetas a datos sin procesar



# 01. Datos: Etiquetado

**Amazon**

SageMaker  
GroundTruth

Coste por etiqueta

500 objetos/mes gratis

Amazon MTurk

# 01. Datos: Etiquetado

**Amazon**  
SageMaker  
GroundTruth

Coste por etiqueta

500 objetos/mes gratis

Amazon MTurk

LabelStudio



Comunidad más amplia

Datos múlti-dominio:

- Texto
- Audio
- Video



Oh oh! Open source detected

# 01. Datos: Etiquetado

**Amazon**  
SageMaker  
GroundTruth

Coste por etiqueta  
500 objetos/mes gratis  
Amazon MTurk

LabelStudio



Comunidad más amplia  
Datos múlti-dominio:

- Texto
- Audio
- Video

CVAT

Más completo para CV  
Integración con ML para  
etiquetado asistido  
V. Community gratuita

# 01. Datos: Control de versiones

Mantener un registro de las distintas versiones de los datos utilizados para nuestros experimentos

- Introducción de nuevos datos
- Reducción de sesgos evidenciados
- Registro de datos



# 01. Datos: Control de versiones

**DVC**

interactive.ai



Git para datos / modelos

Conecta el  
almacenamiento con el  
repositorio de código

Registro

# 01. Datos: Control de versiones

**DVC**

interactive.ai



Git para datos / modelos

Conecta el  
almacenamiento con el  
repositorio de código

Registro

**Git LFS**



Extensión de Git

Datos de gran tamaño

**[Cuidado]** No incluir  
versionado de código  
conjuntamente

# 01. Datos: Control de versiones

**DVC**

interactive.ai



Git para datos / modelos

Conecta el  
almacenamiento con el  
repositorio de código

Registro

Git LFS



Extensión de Git

Datos de gran tamaño

**[Cuidado]** No incluir  
versionado de código  
conjuntamente

Pachyderm

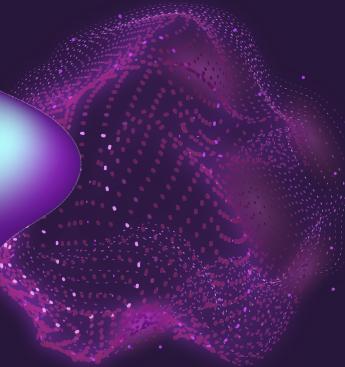
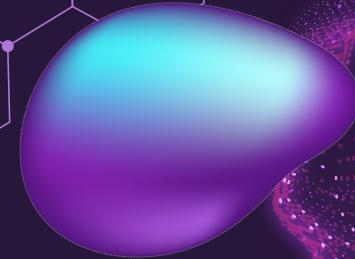
Version community  
“apañá”

Pipelines de validación

Integración con  
frameworks principales

0

x



x

2

# ENTRENAMIENTO / VALIDACIÓN

+



# Entrenamiento

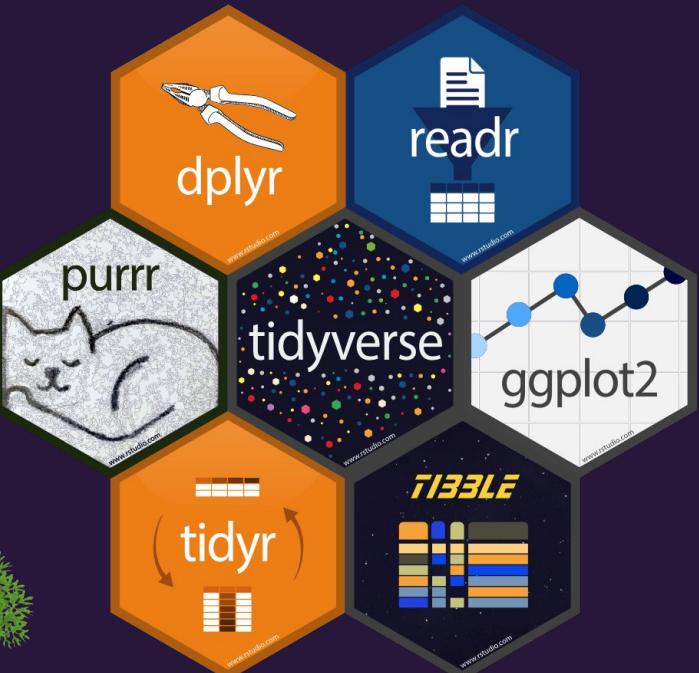
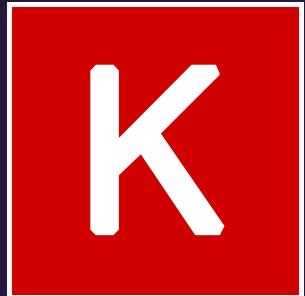
Aprender patrones o relaciones a partir de datos.

Ajuste de los parámetros internos del modelo.

Dependiente de una validación

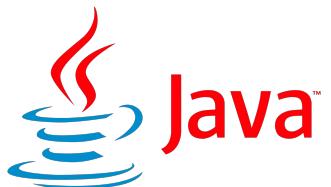
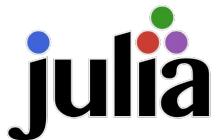
# 02. Entrenamiento: Frameworks

El de tu elección: Depende de tu caso de uso de aplicación ML



## 02. Entrenamiento: IDE

El de tu elección: Depende de framework / lenguaje



...



## 02. Entrenamiento: Control

### WandB

Cloud & self-hosted  
Vistas personalizadas  
Gestión de equipos &  
roles

**[Importante]** Caro

# 02. Entrenamiento: Control

## WandB

Cloud & self-hosted  
Vistas personalizadas  
Gestión de equipos &  
roles  
**[Importante]** Caro

## MLFlow



Self-hosted  
Lo justo para hacer de  
todo  
Control de versiones de  
modelos & registro

# 02. Entrenamiento: Control

## WandB

Cloud & self-hosted  
Vistas personalizadas  
Gestión de equipos & roles  
**[Importante]** Caro

## MLFlow

Self-hosted  
Lo justo para hacer de todo  
Control de versiones de modelos & registro

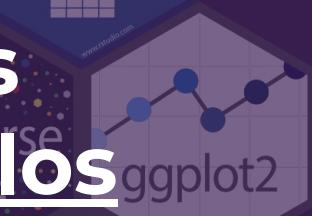
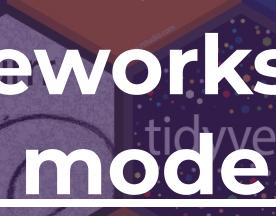


## Comet

Despliegues automatizados  
Accesibilidad a terceras partes  
Entorno colaborativo

## 02. Entrenamiento: Frameworks

Diversidad de frameworks  
Distintos formatos de modelos



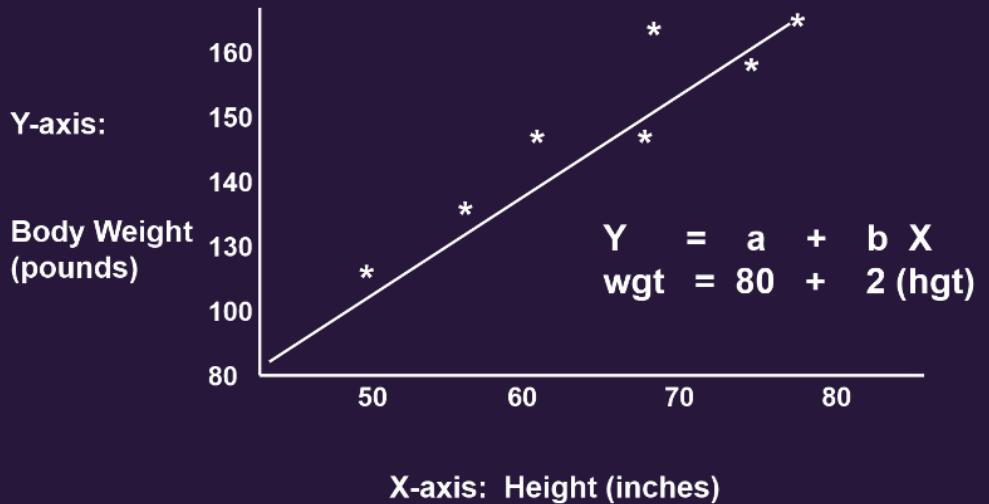
## 02. Entrenamiento: Interoperabilidad

Necesidad de un estándar:

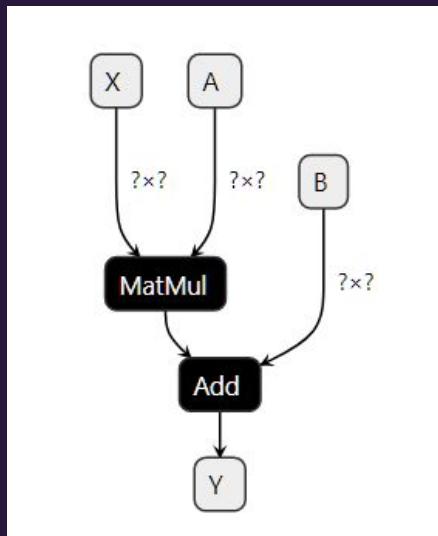
- Interoperabilidad entre sistemas/frameworks
- Adaptado a la innovación/ Soporte de la comunidad



## 02. Entrenamiento: ONNX



```
def onnx_linear_regressor(X):
    "ONNX code for a linear regression"
    return onnx.Add(onnx.MatMul(X, coefficients), bias)
```



# 02. Entrenamiento: Interoperabilidad

## Frameworks



Native support

Converters



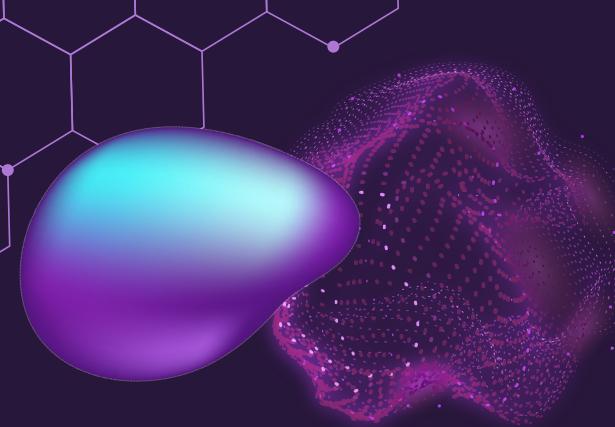
ONNX Model

0

# 3 Despliegue



x



x

+



# Despliegue

Poner el modelo ML a disposición del usuario final / grupos de interés.

Monitorizar el rendimiento del modelo

Automatizar el despliegue



## 03. Despliegue: Parcial



Evaluación por parte del equipo de desarrollo

Difusión interna/ marketing del modelo

Servicio web



**Ejemplos**

# 03. Despliegue: Control

Monitorizar modelo en producción

WandB

Cloud & self-hosted  
Vistas personalizadas  
Gestión de equipos &  
roles  
**[Importante]** Caro

MLFlow

Self-hosted  
Lo justo para hacer de  
todo  
Control de versiones de  
modelos & registro



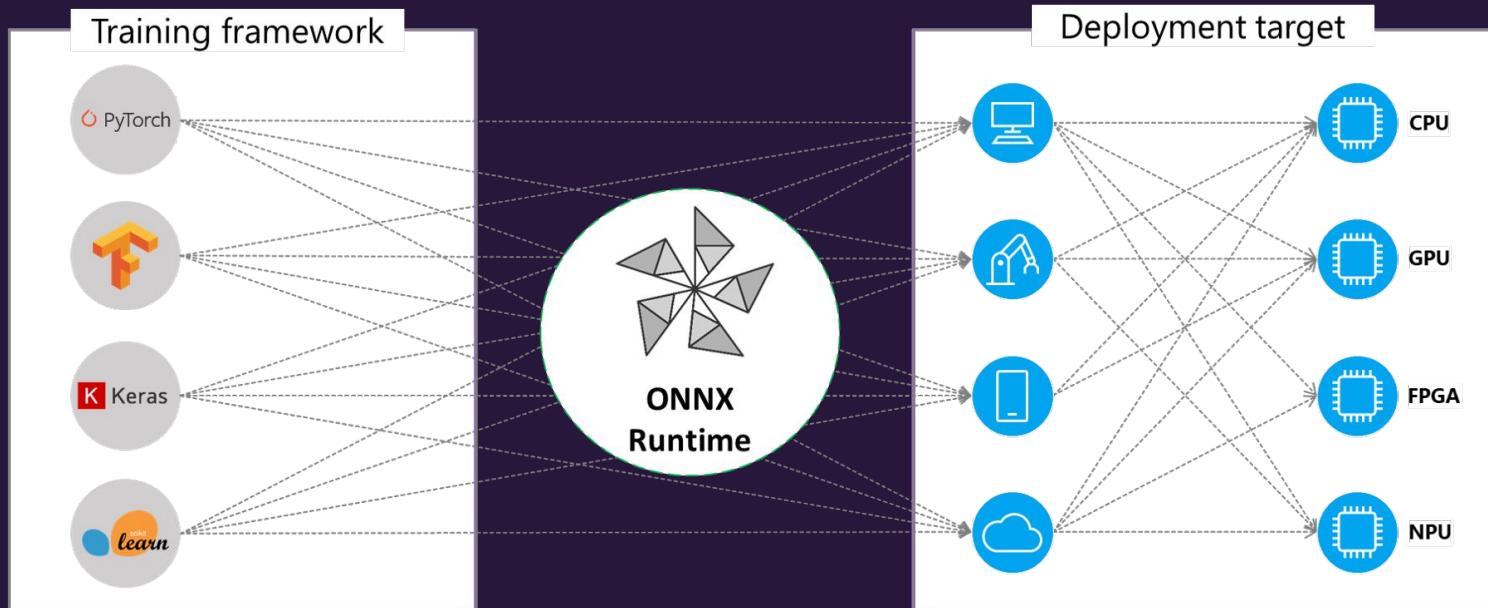
Comet

Despliegues  
automatizados  
Accesibilidad a terceras  
partes  
Entorno colaborativo

# 03. Despliegue: Arquitectura

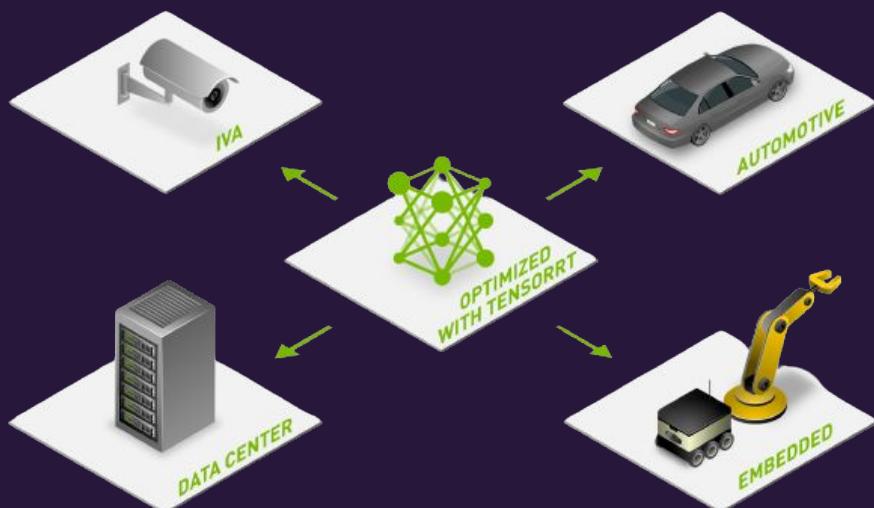


Elección del dispositivos de procesamiento. ONNX como nexo.

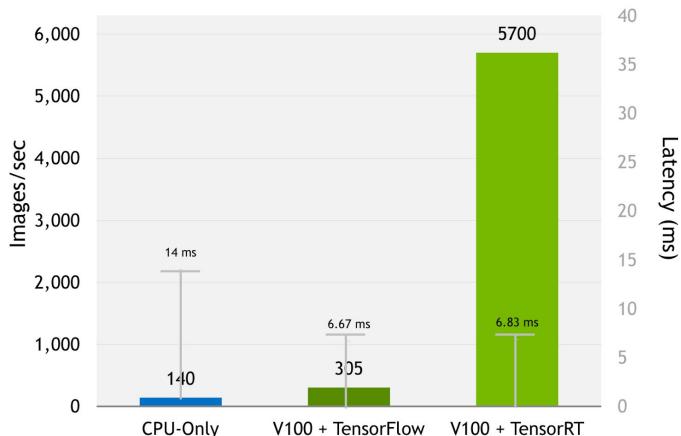


# 03. Despliegue: NVIDIA TensorRT

Framework de NVIDIA para NVIDIA



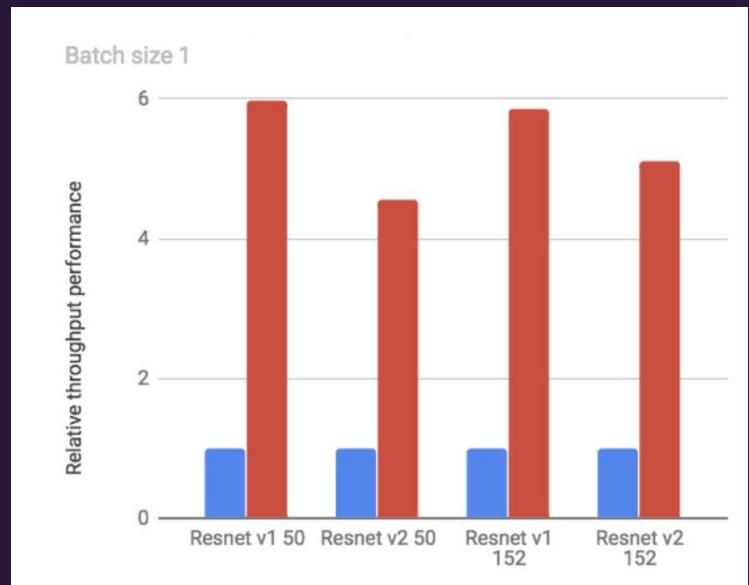
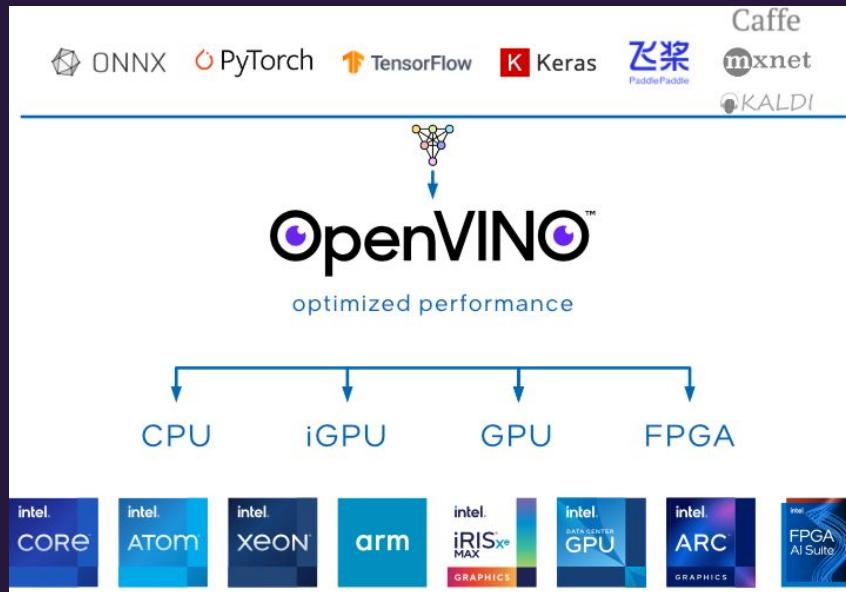
Up to 40x Faster CNNs on V100 vs. CPU-Only Under 7ms Latency (ResNet50)



Inference throughput (images/sec) on ResNet50. **V100 + TensorRT:** NVIDIA TensorRT (FP16), batch size 39, Tesla V100-SXM2-16GB, E5-2690 v4@2.60GHz 3.5GHz Turbo (Broadwell) HT On. **V100 + TensorFlow:** Preview of volta optimized TensorFlow (FP16), batch size 2, Tesla V100-PCIE-16GB, E5-2690 v4@2.60GHz 3.5GHz Turbo (Broadwell) HT On. **CPU-Only:** Intel Xeon-D 1587 Broadwell-E CPU and Intel DL SDK. Score doubled to comprehend Intel's stated claim of 2x performance improvement on Skylake with AVX512.

# 03. Despliegue: Intel OpenVINO

Framework de Intel para Intel





# 03. Despliegue: CI/CD

Nueva tecnología, mismos conocidos:

Github  
actions

Gitlab CI/CD

Jenkins 

Pipelines de acciones ante eventos de desencadenado:

- Nuevos datos de entrenamiento
- Nuevo modelo más efectivo
- Nuevo modelo más eficiente

04

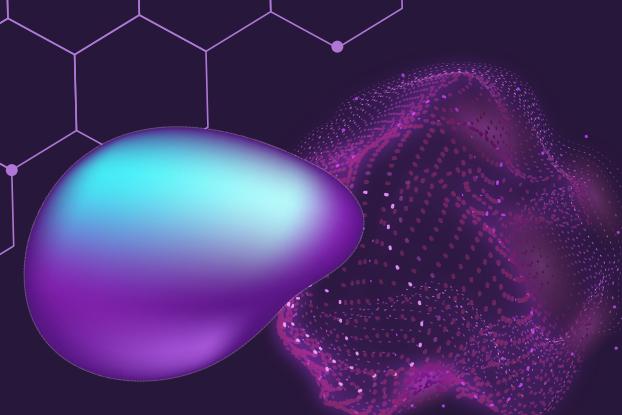
# Flujo completo



x



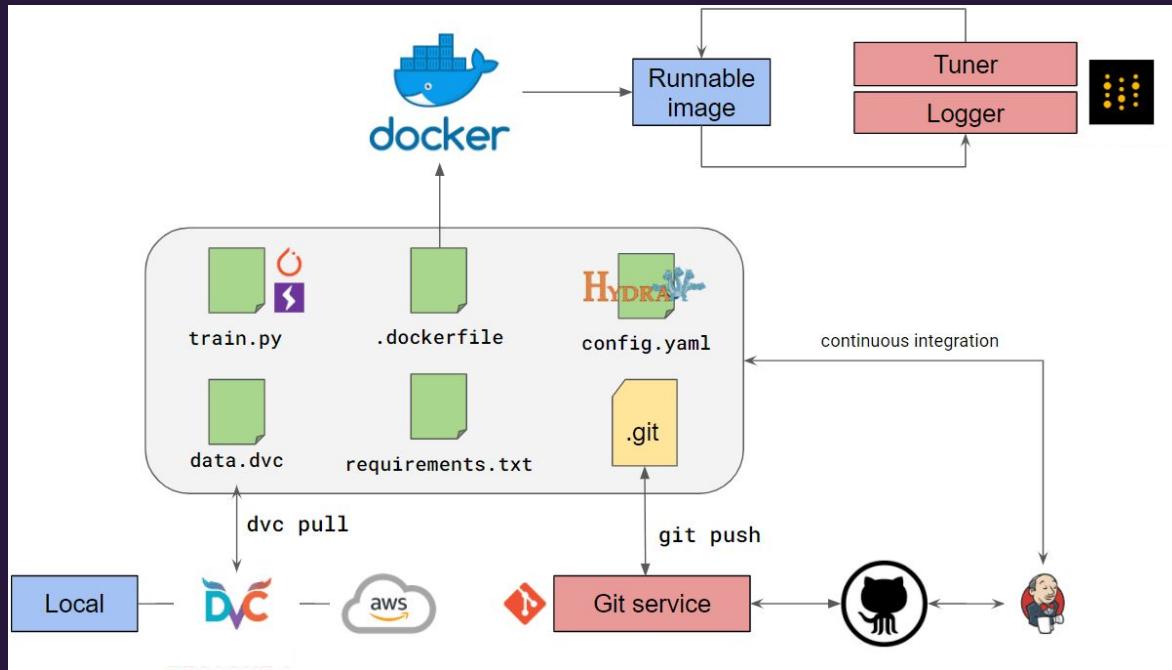
+

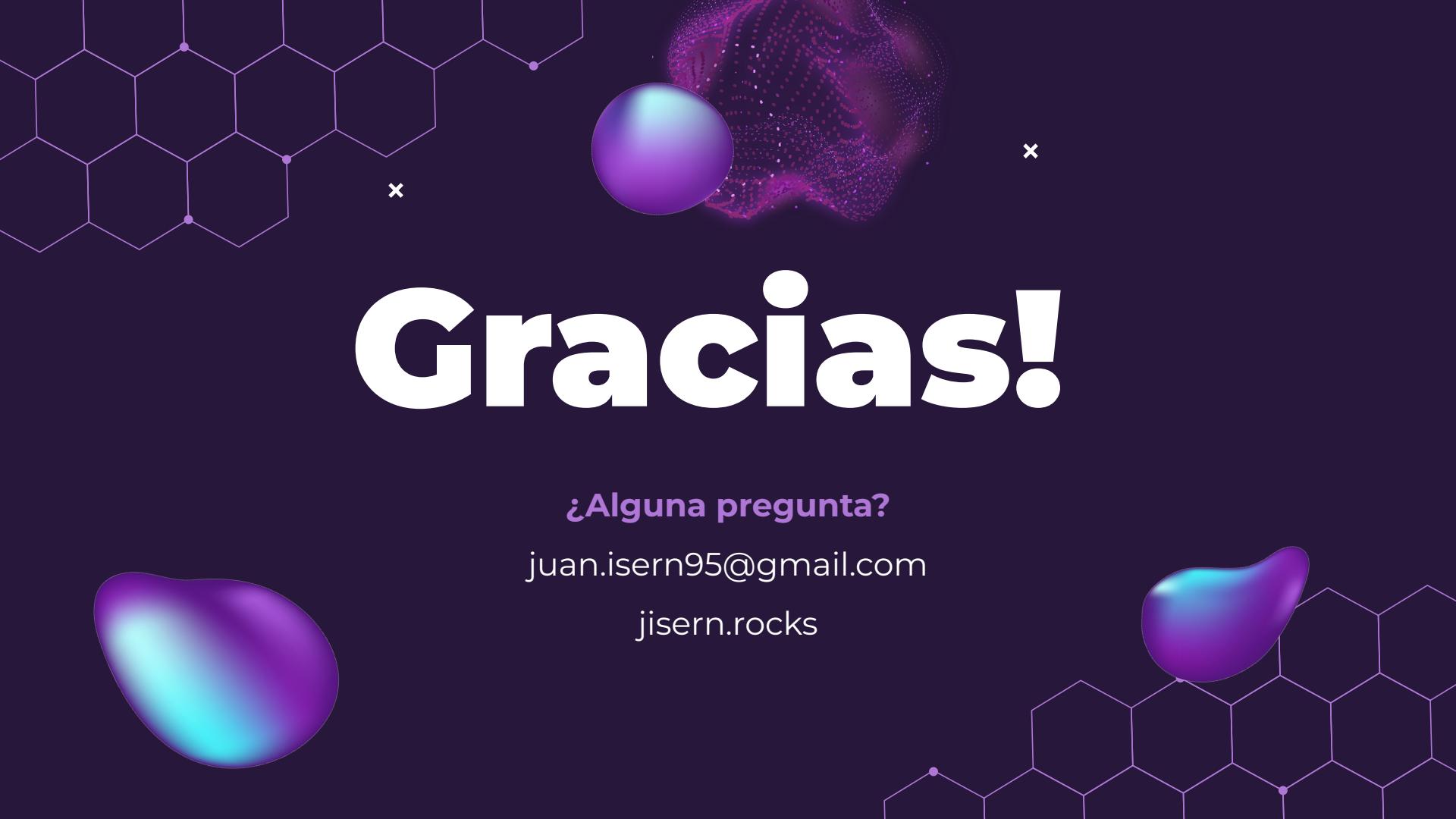


x



# 04. Ejemplo operativo





# Gracias!

¿Alguna pregunta?

[juan.isern95@gmail.com](mailto:juan.isern95@gmail.com)

[jisern.rocks](http://jisern.rocks)

# 04. Otros: IaaS-ML

