# A perspective on MLOps from experimentation: Weights and Bias implementations

**Globant**

**Daniel Jiménez**

**2024-05-17**

# Content

- Motivation
  - ✅ Machine Learning and Hardware
  - ✅ AI Computing
  - ✅ Challenges of experimentation

- Weights and Bias Platform

- Weights and Bias Example

# Motivation

**Thinking about machine learning in terms of three components:**

- Algos (LSTM, XGBoost, Random Forest)

- Data (From accessibility, and the cookies that you and I accept)

- Hardware (Gpus, CPus, architecture)

# Motivation

**Thinking about machine learning in terms of three components:**

Algos
Data            >    Computing    >    Resource demand
Hardware

# Motivation

## Thinking about machine learning in terms of three components:
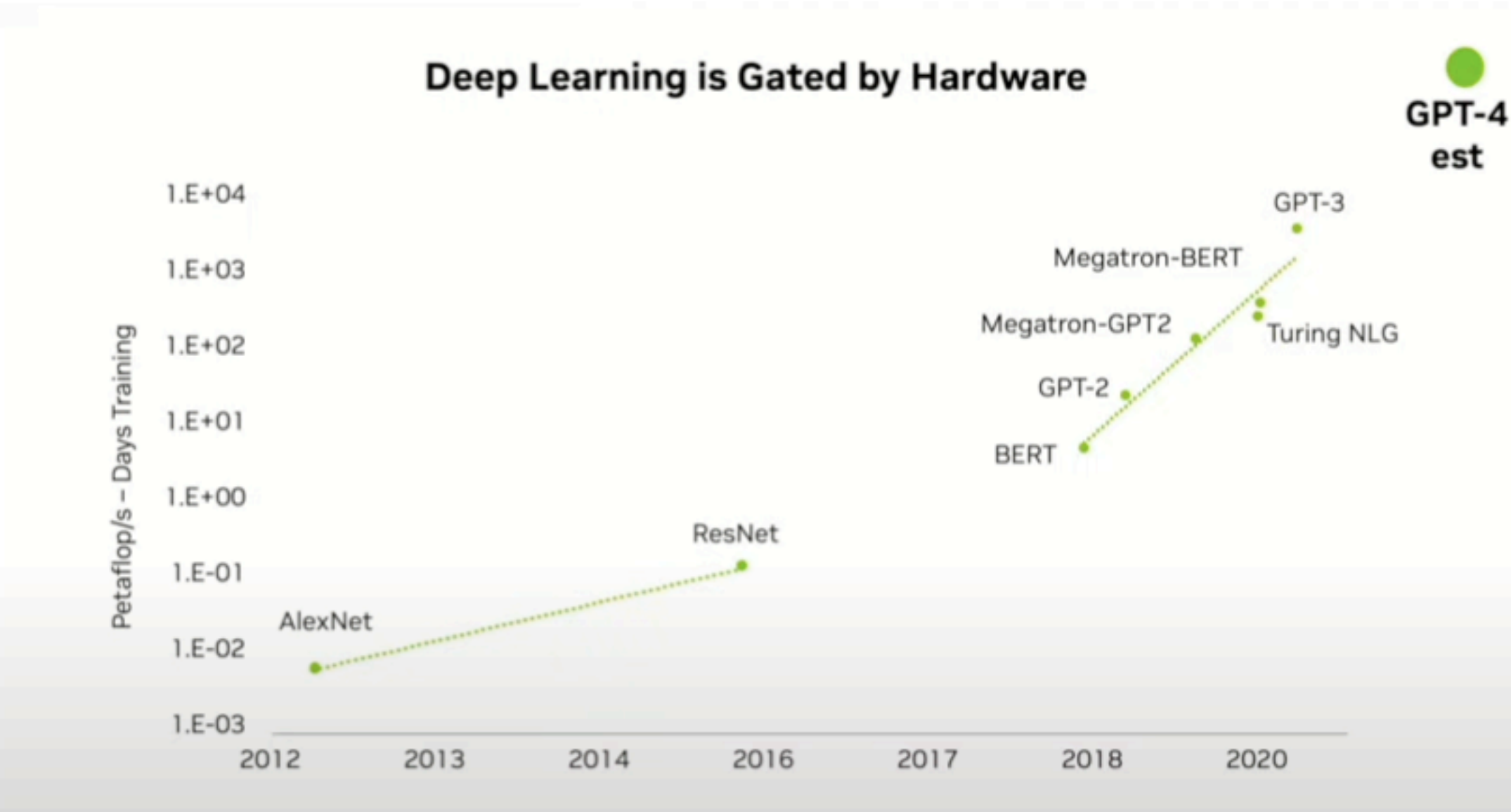
Algos
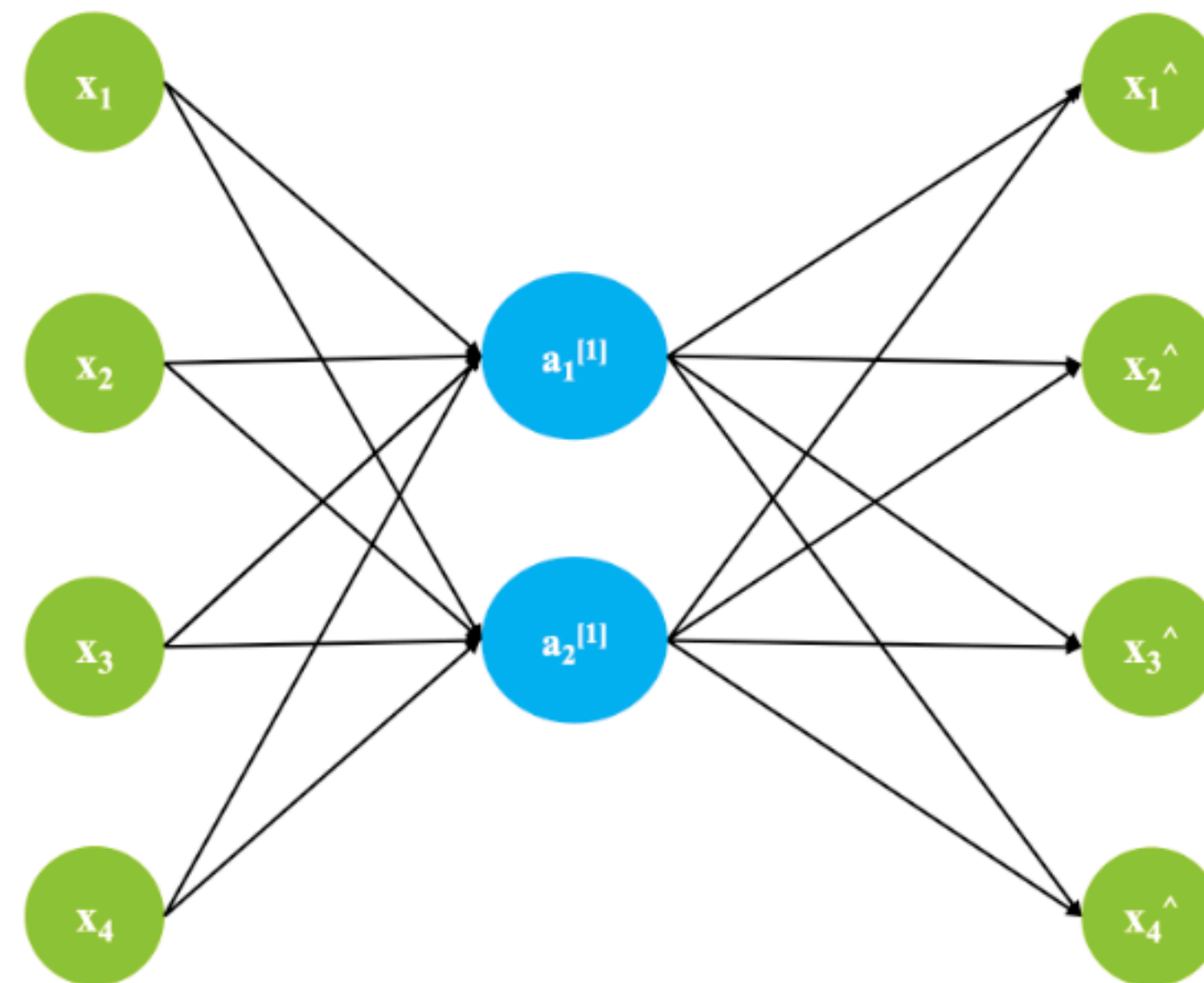Data       >     Computing     >     Resource demand
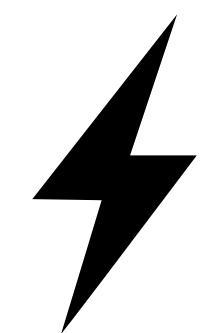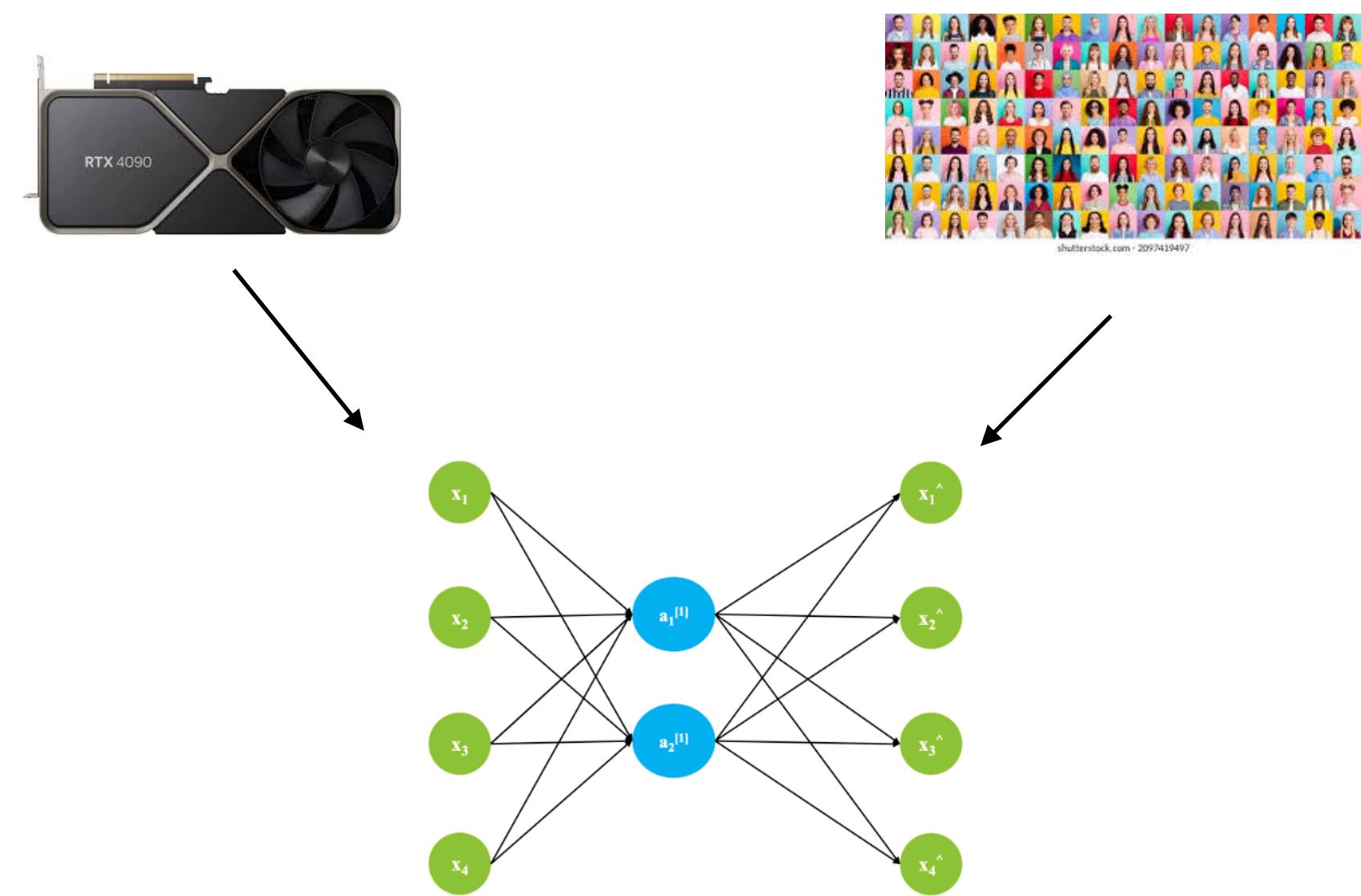Hardware

Inference

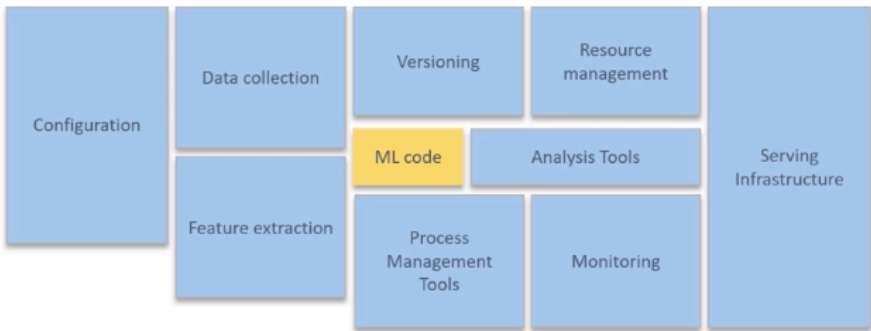# Motivation



Nvidia Reference



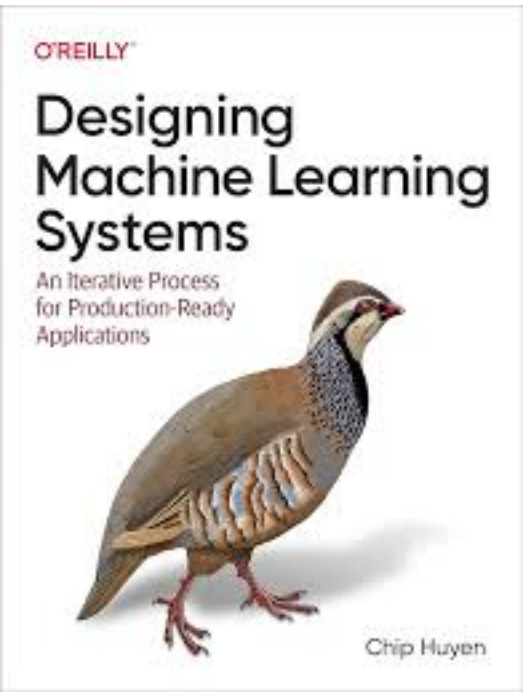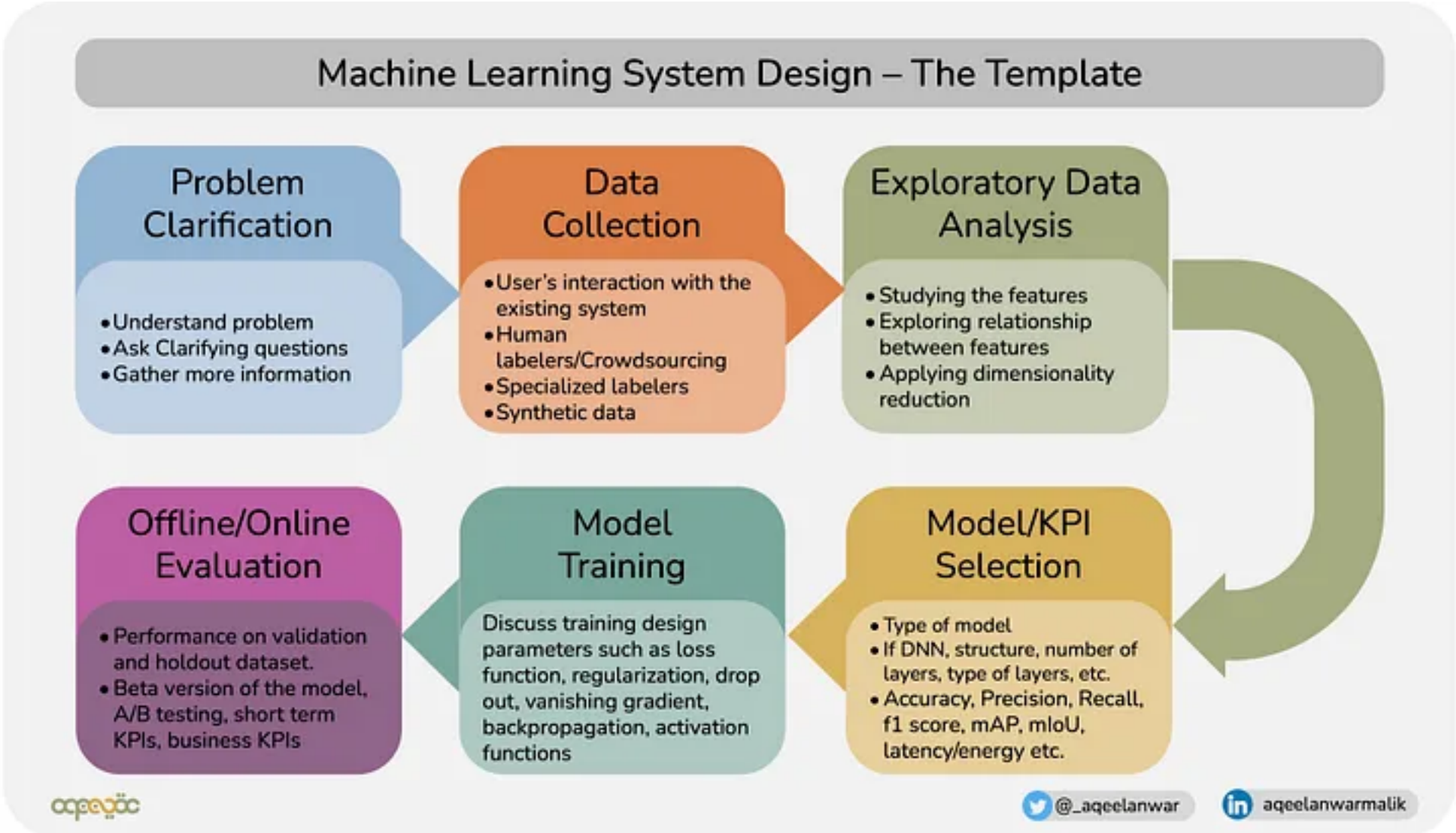Carlos Santana Tweet
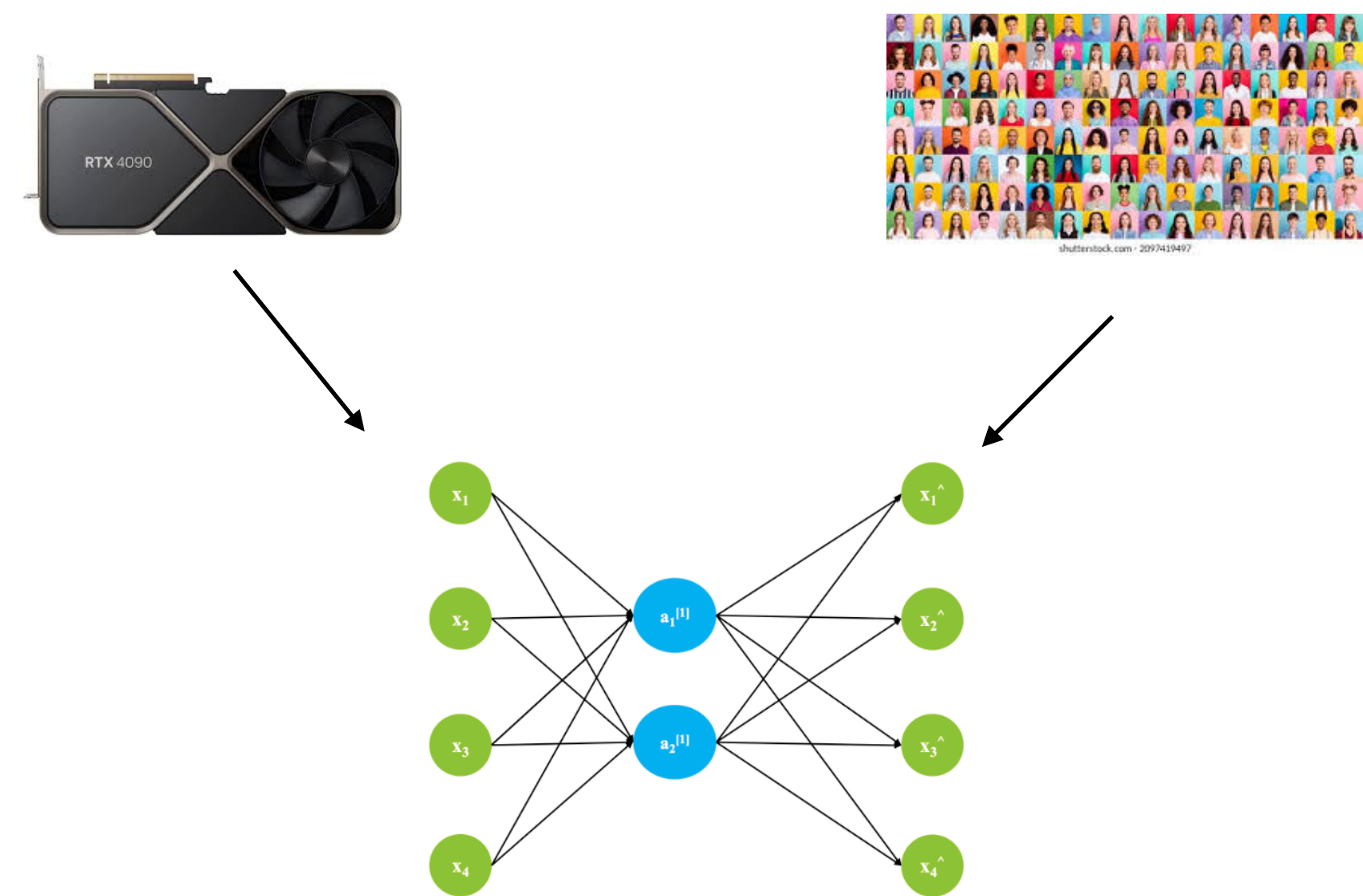
# Motivation
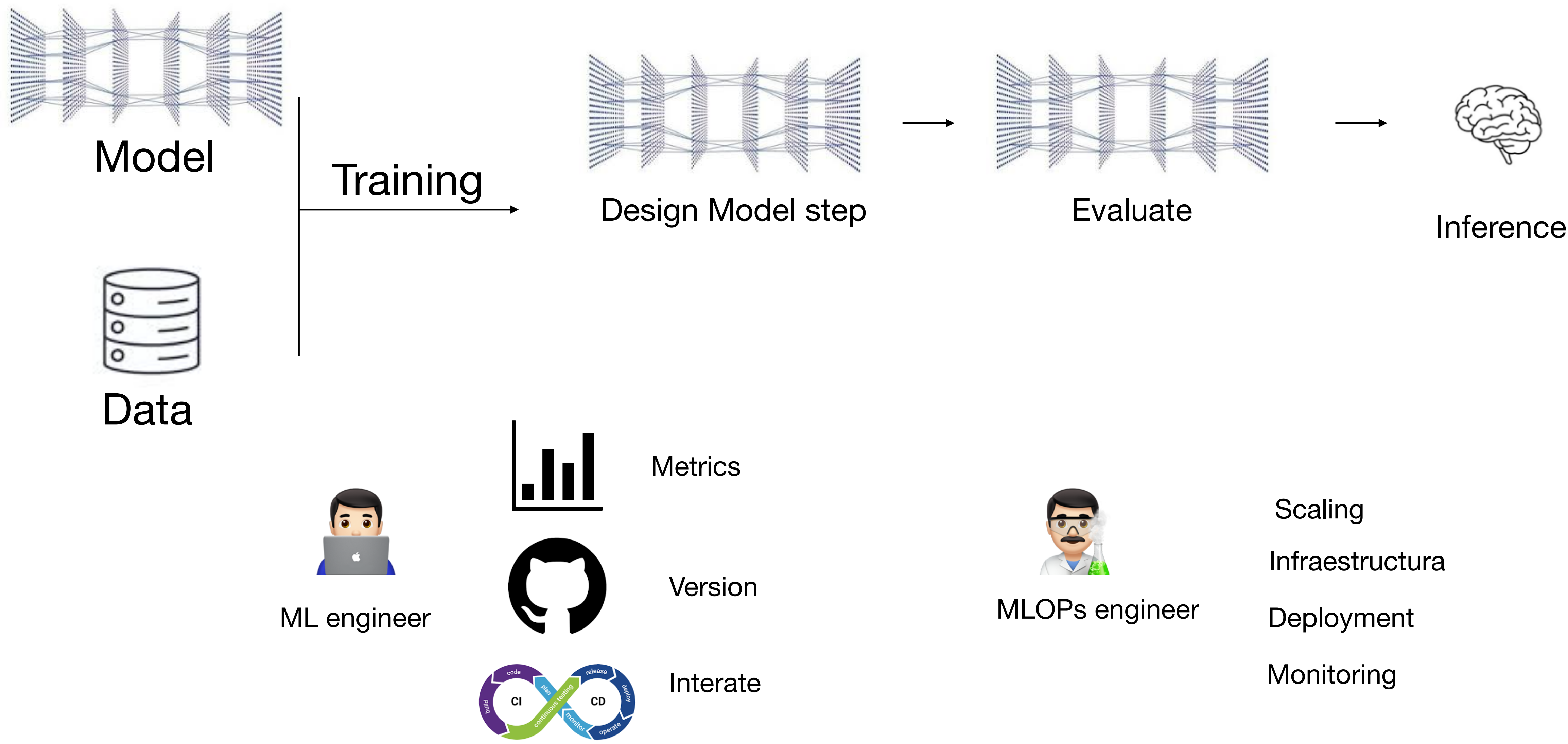
# Motivation



Consume

Architecture

<u>Machine Learning Systems</u>

# Motivation





Machine Learning System Design — The template

# Motivation



Model

Data

Training

Design Model step

Evaluate

Inference

Metrics

Version

Interate

ML engineer

MLOPs engineer

Scaling

Infraestructura

Deployment

Monitoring

# Motivation



Model

Training → Design Model step → Evaluate → Inference → Deploy

Data

Metrics

Version

Interate

ML engineer

MLOPs engineer

Scaling

Infraestructura

Deployment

Monitoring

Manager

# Motivation



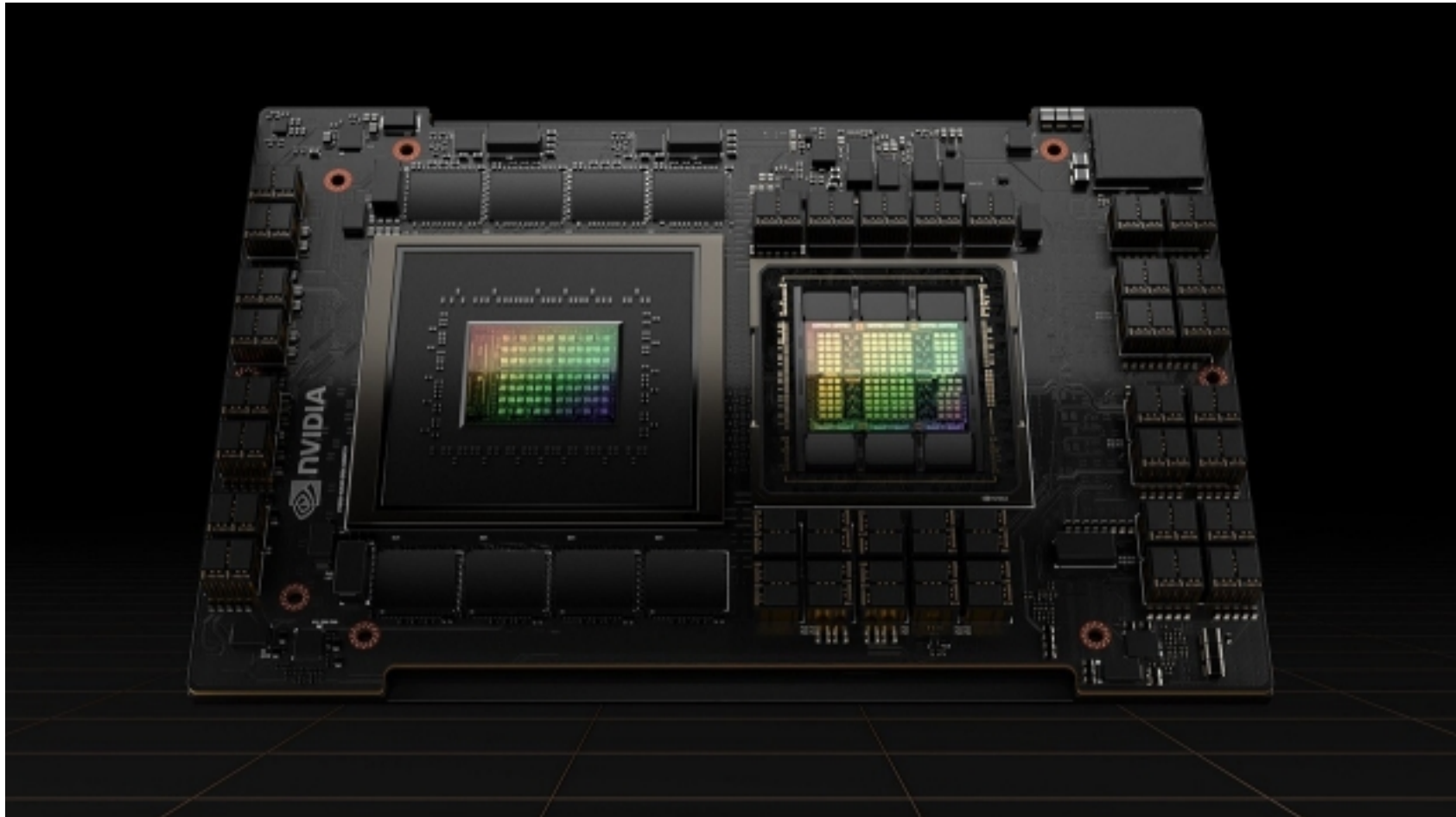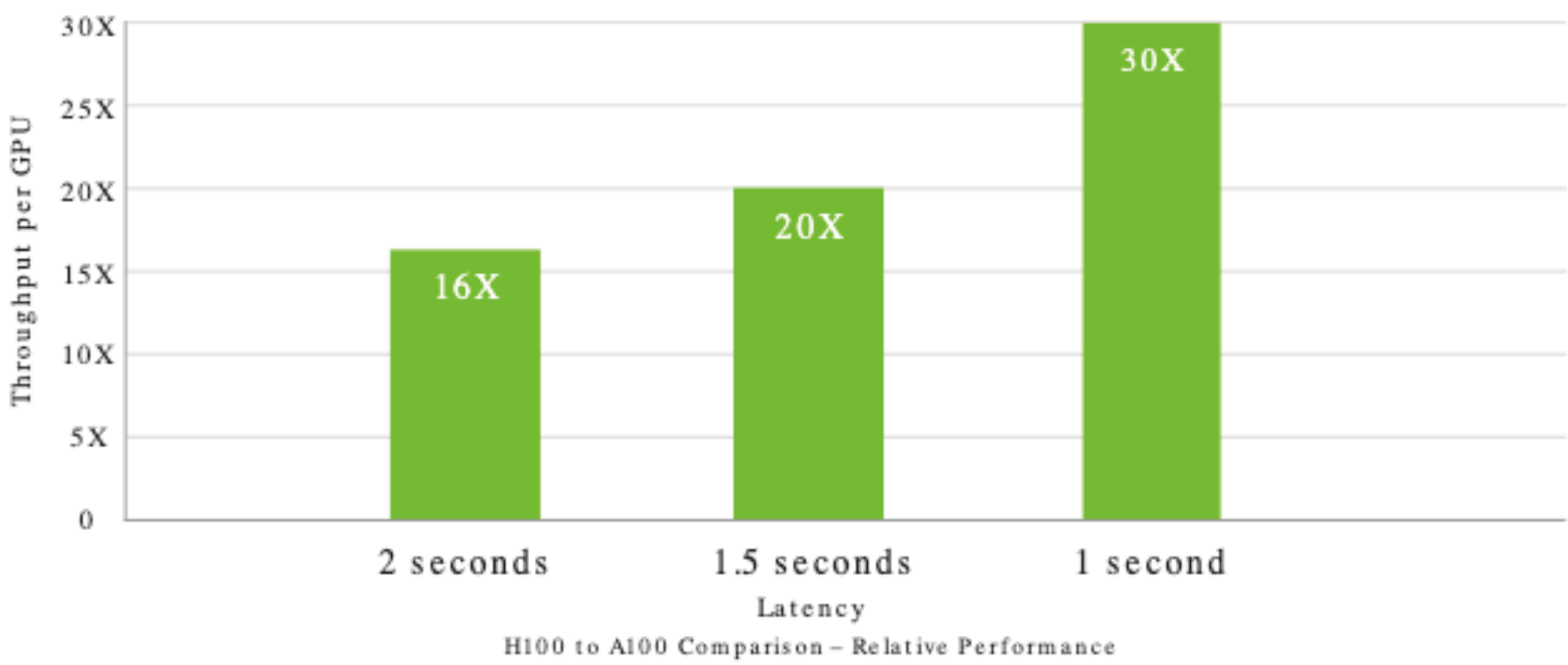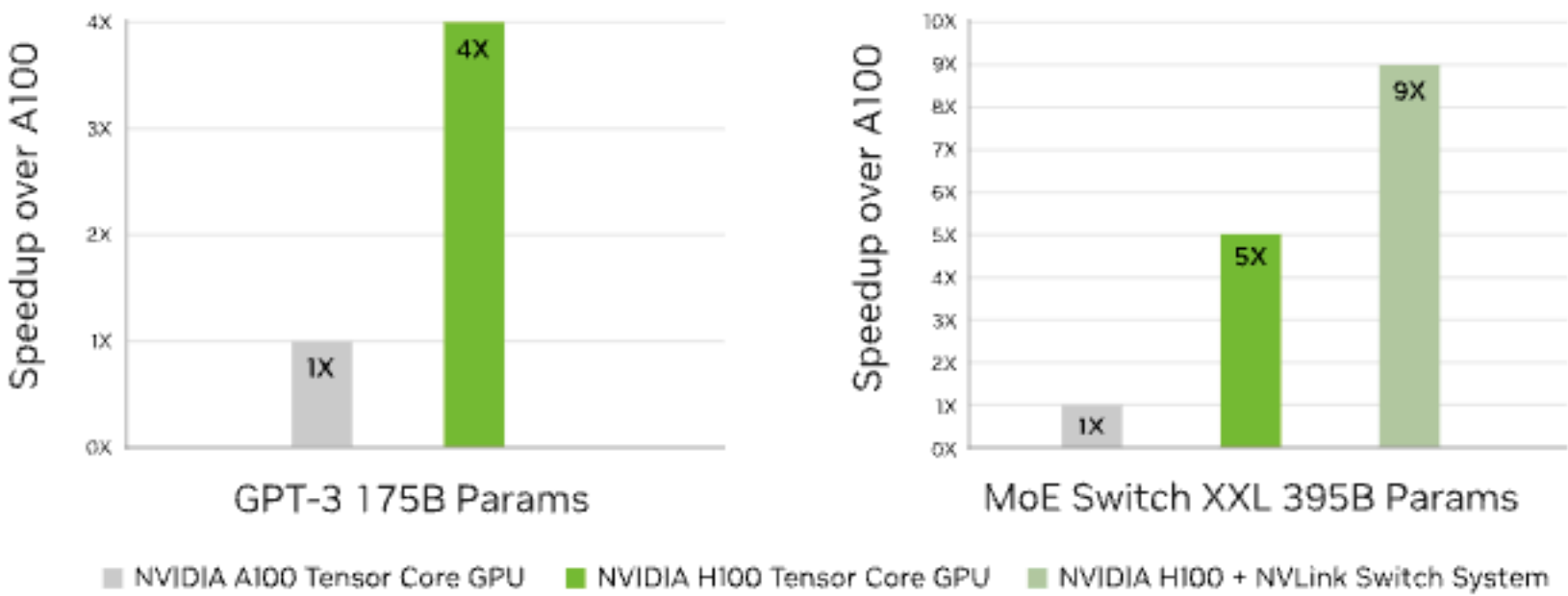Hopper 100



**Rendimiento de Inferencia de IA hasta 30 Veces Mayor en los Modelos Más Grandes**

Inferencia del chatbot Megatron (530 mil millones de parámetros)



Rendimiento proyectado sujeto a cambios. Inferencia en el chatbot basado en el modelo de parámetros Megatron 530B para longitud de secuencia de entrada = 128, longitud de secuencia de salida = 20 | Clúster A100: red HDR IB | Clúster H100: sistema de conmutador NVLink, NDR IB

**Entrenamiento de IA Hasta 4 Veces Superior en GPT-3**



Rendimiento proyectado sujeto a cambios. Entrenamiento GPT-3 175B Clúster A100: red HDR IB, clúster H100: red NDR IB | Variante del Switch-XXL del transformador de capacitación de mezcla de expertos (MoE) con parámetros 395B en un conjunto de datos de token de 1T, clúster A100: red HDR IB, clúster H100: red NDR IB con sistema de conmutador NVLink donde se indique.

# AI Computing

- Mathematically intensive process to calculate machine learning algorithms using accelerated hardware.

- Three Elements to AI Computing

  - ELT

  - Model Design

  - Inference

# Challenges of experimentation

- Making the pipeline of models for inference efficient

- Tracking model experimentation

- GPU-CPU monitoring

- Understanding time consumption per algorithm phase

- Resource configuration for the problem being addressed (Machine Learning system)

# Weights and Bias Platform

# Weights and Bias

- https://wandb.ai/site

- Build models faster, fine-tune LLMs, develop GenAI applications with confidence, all in one system of record developers are excited to use. (https://wandb.ai/site)

# Weights and Bias VS MLflow

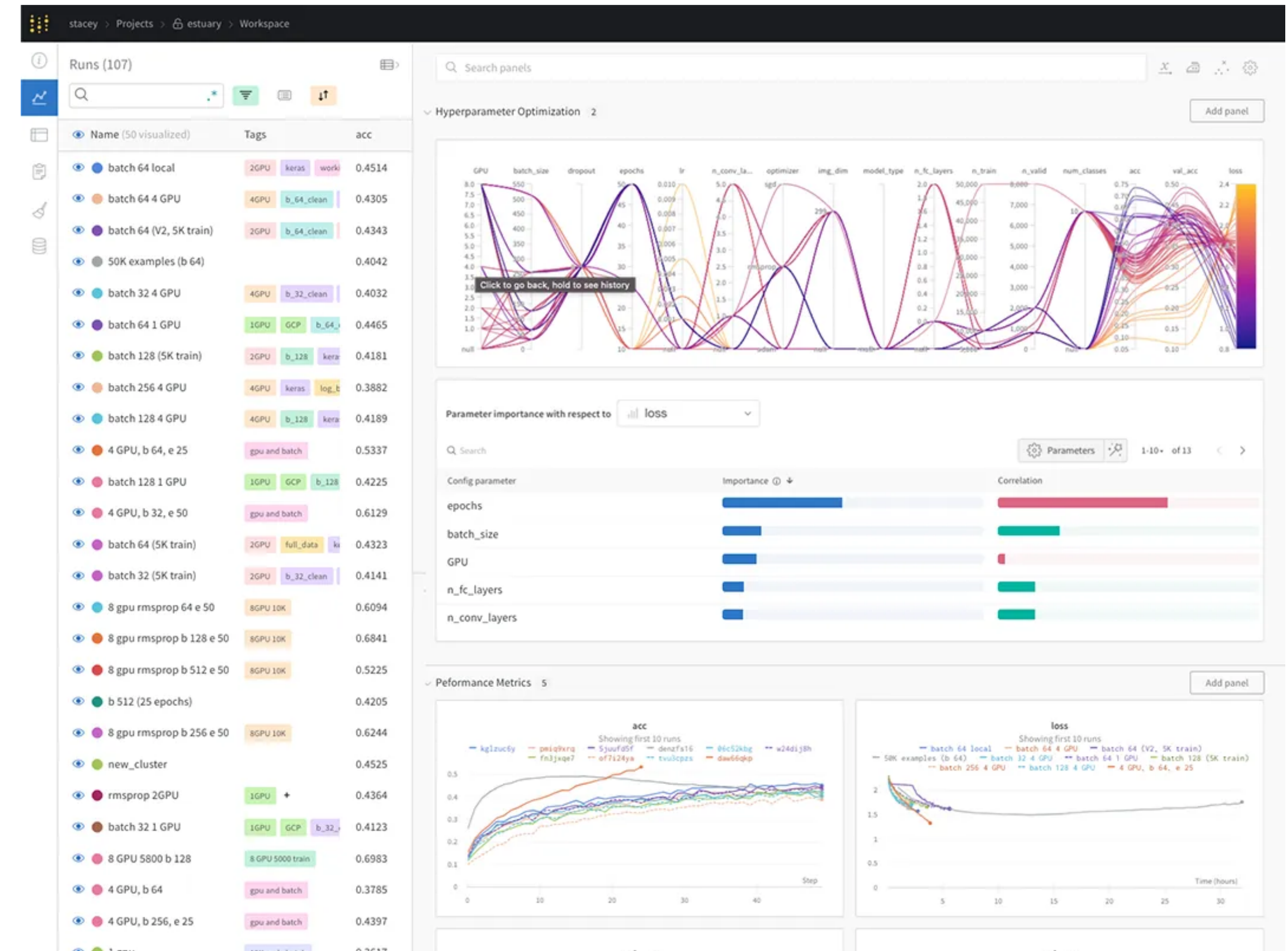Weights and Bias : It is specifically designed for tracking machine learning experiments.

MLflow: It follows the complete lifecycle of machine learning, including experiment tracking, model management, and deployment

# Weights and Bias Platform

| Function | Tool |
|---|---|
| **Code Tracking** | Github, Gitlab, Bitbucket |
| **Experiment tracking** | Weights and bias , MLflow |
| **Artifact tracking** | Weights and bias |
| **Model repository** | Weights and bias |
| **ML pipelines** | MLflow |
| **Environment isolation** | Conda - Docker |
| **Orchestration** | MLflow |
| **Modelling** | Tensorflow, Pytorch, scikit-learn |

# Weights and Bias Platform

- Work is done under artifacts (best practices).

-  Pipelines should yield different artifacts as results.
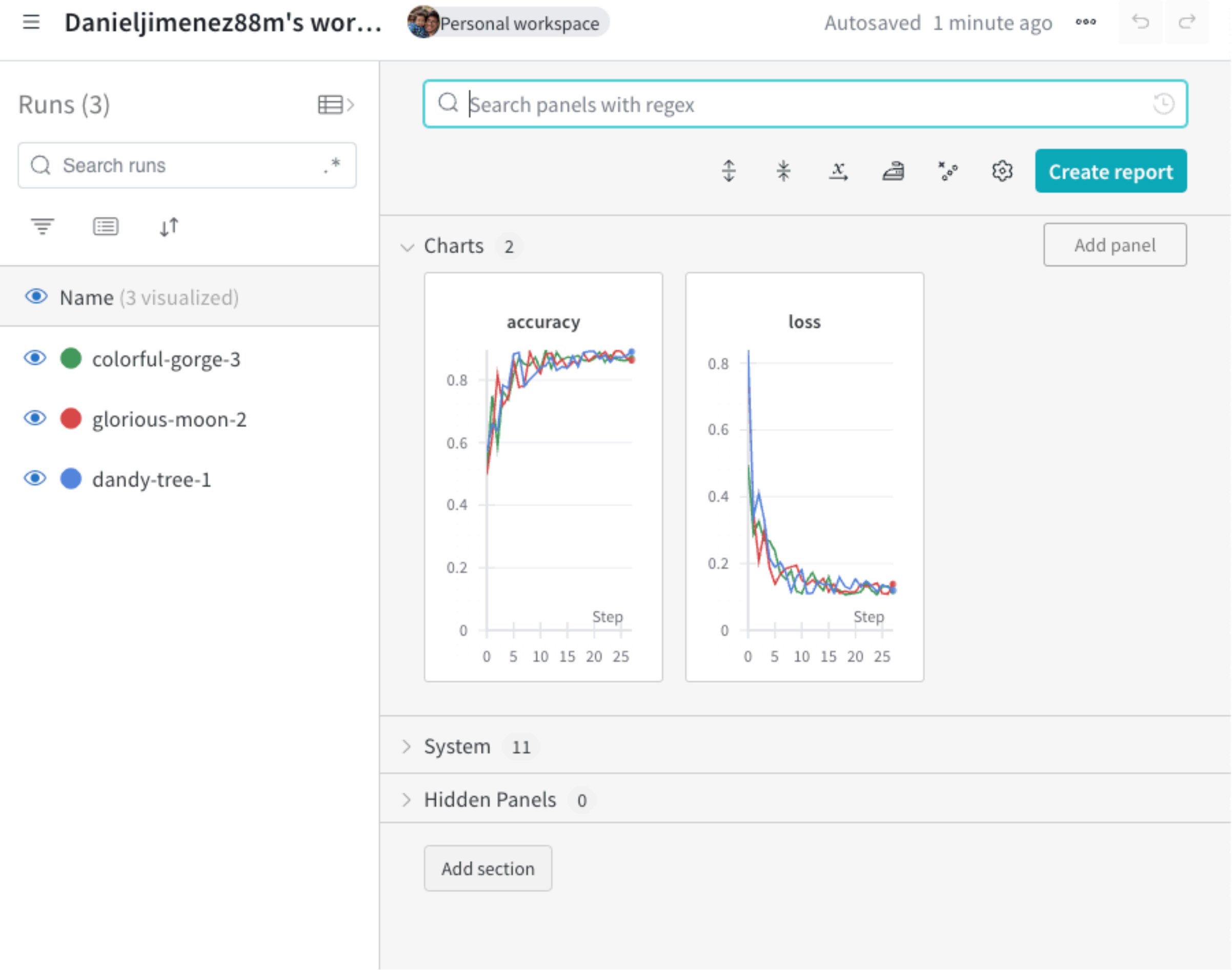
-  Artifacts must be versioned, tracked, and saved.
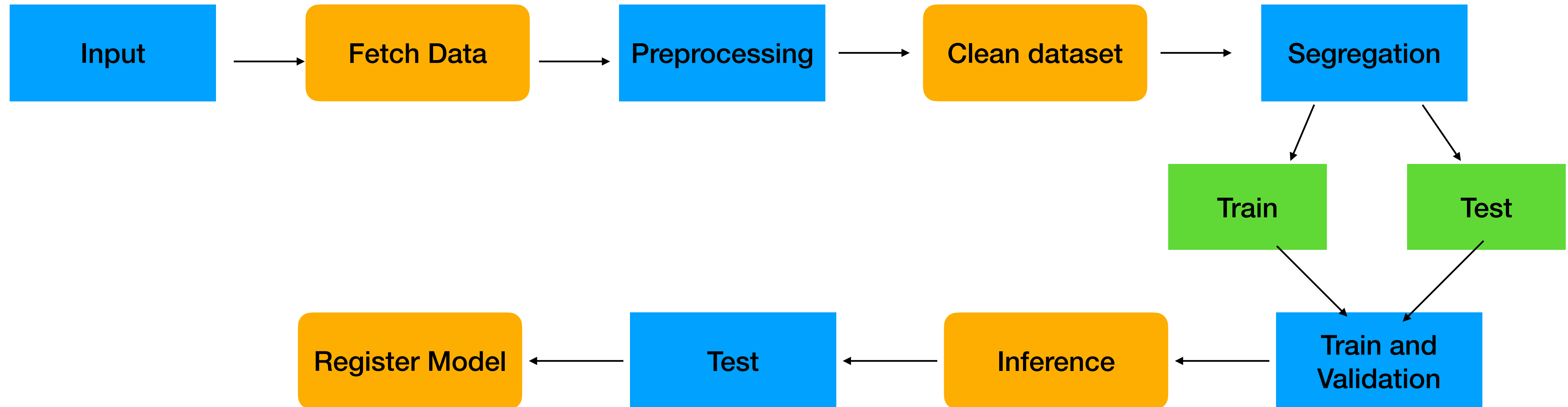
Experiments names

Metrics logs

System Performs
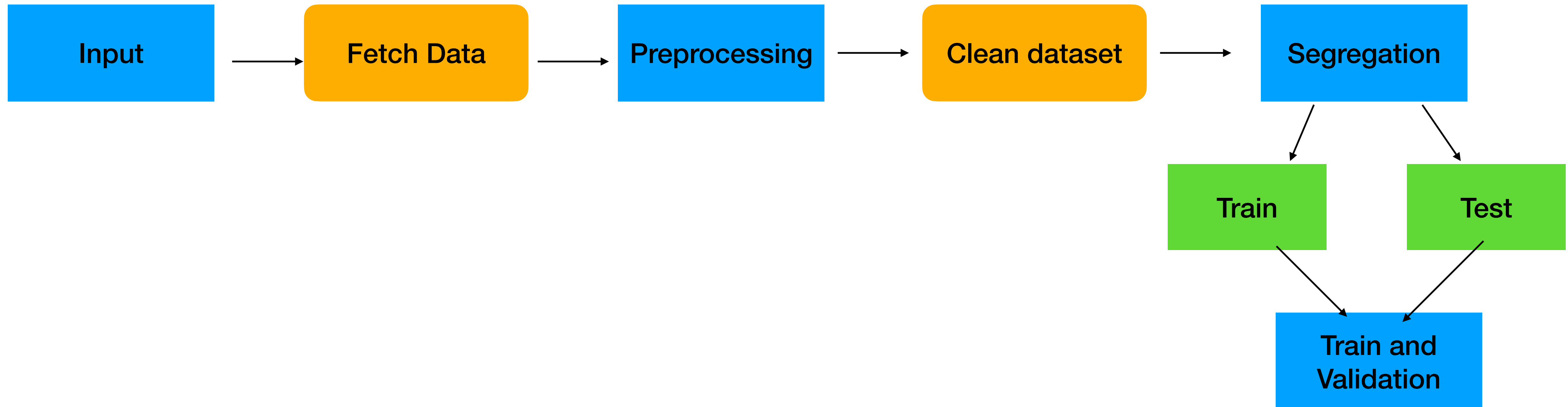
System Performs

General Options

# Weights and Bias Propose Roadmap

- Generate a Infrastructure code like a system with Artifacts

- We will follow a design path focused on data science engineering, although we will not build an API.

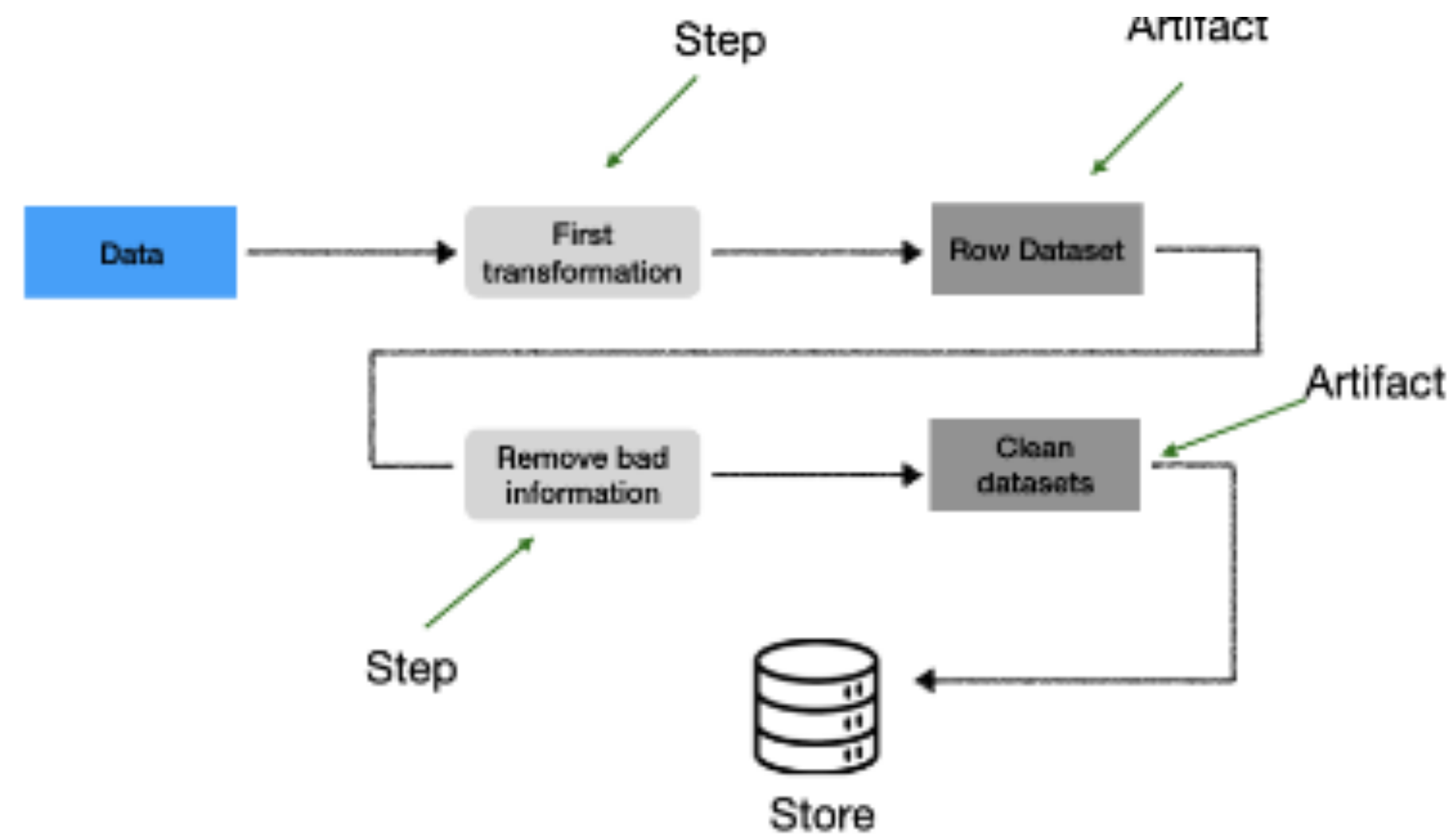- It will end with the registration of a model, to be used in inference.

# Weights and Bias Propose Roadmap

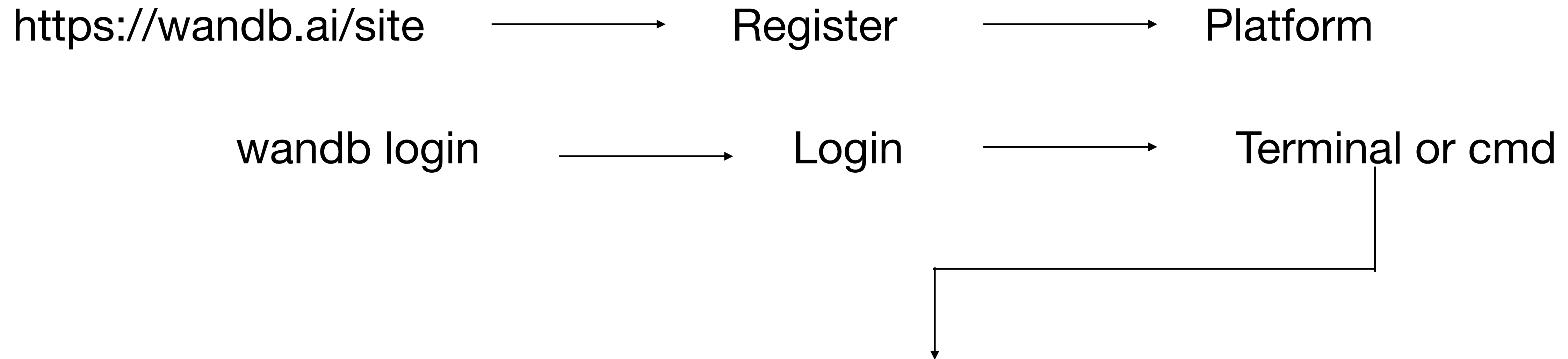# Weights and Bias Propose Roadmap segmentation

# Weights and Bias Platform



**Artifacts:** These are the products resulting from the components. Artifacts can be used as inputs in subsequent components, thus creating connections between the different steps of the pipeline. Artifacts should be carefully tracked and versioned to ensure traceability and efficient management of the different model development versions.

# Weights and Bias Example

# Roadmap

https://wandb.ai/site  ⟶  Register  ⟶  Platform

wandb login  ⟶  Login  ⟶  Terminal or cmd

Generally, they send a link where there's a code that you must paste into the terminal in order to access

# Roadmap

# Roadmap



Python Scripting with Argparse

# Python Scripting with Argparse

- Automating tasks

- Using the command line for customization of situations in front of model design

- Flexible code configuration

- By having a script with Argparse that contains all the necessary steps to train and evaluate your model, other users can easily reproduce your experiments and results

# Python Scripting with Argparse

Colocar el link del código

# Roadmap into WandB

wand.init()

wand.log()

wand.log_artifact()

# Roadmap into WandB

```python
wand.init(
project='name',
job_type='some amazing process',
group='experiment_final_final_ya_no_va_más'
)
```

```python
wandb.Artifact(
name=args.artifact_name,
type=args.artifact_type,
description=args.artifact_description,
metadata={'original_url': args.file_url})
```

# Now let's write code!!! That's what we came here for."

https://github.com/Carlos-Jimenez-mlops/wandb-mlops