# Effortless AI Scaling: Red Hat OpenShift and LLM Integration Production-grade AI

Piros.

Indy Van Mol
Piros

# The explosion of AI and its integration into all business applications

# The LLM take over

**LLMs**

| Model | Version | Retirement date | Suggested replacements |
|---|---|---|---|
| `dall-e-2` | 2 | January 27, 2025 | `dalle-3` |
| `dall-e-3` | 3 | No earlier than April 30, 2025 | |
| `gpt-35-turbo` | 0301 | January 27, 2025<br><br>Deployments set to Auto-update to default will be automatically upgraded to version: `0125`, starting on November 13, 2024. | `gpt-35-turbo` (0125)<br><br>`gpt-4o-mini` |
| `gpt-35-turbo`<br>`gpt-35-turbo-16k` | 0613 | January 27, 2025<br><br>Deployments set to Auto-update to default will be automatically upgraded to version: `0125`, starting on November 13, 2024. | `gpt-35-turbo` (0125)<br><br>`gpt-4o-mini` |
| `gpt-35-turbo` | 1106 | No earlier than January 27, 2025<br><br>Deployments set to Auto-update to default will be automatically upgraded to version: `0125`, starting on November 13, 2024. | `gpt-35-turbo` (0125)<br><br>`gpt-4o-mini` |
| `gpt-35-turbo` | 0125 | No earlier than Feb 22, 2025 | `gpt-4o-mini` |
| `gpt-4`<br>`gpt-4-32k` | 0314 | June 6, 2025 | `gpt-4o` |
| `gpt-4`<br>`gpt-4-32k` | 0613 | June 6, 2025 | `gpt-4o` |
| `gpt-4` | 1106-preview | To be upgraded to `gpt-4` version: `turbo-2024-04-09`, starting no sooner than January 27, 2025 [1] | `gpt-4o` |
| `gpt-4` | 0125-preview | To be upgraded to `gpt-4` version: `turbo-2024-04-09`, starting no sooner than January 27, 2025 [1] | `gpt-4o` |
| `gpt-4` | vision-preview | To be upgraded to `gpt-4` version: `turbo-2024-04-09`, starting no sooner than January 27, 2025 [1] | `gpt-4o` |
| `gpt-4o` | 2024-05-13 | No earlier than March 20, 2025<br><br>Deployments set to Auto-update to default will be automatically upgraded to version: `2024-08-06`, starting on December 5, 2024. | |

**Closed-source vs. open-weight models**

Llama 3.1 405B closes the gap with closed-source models for the first time in history.

@maximelabonne

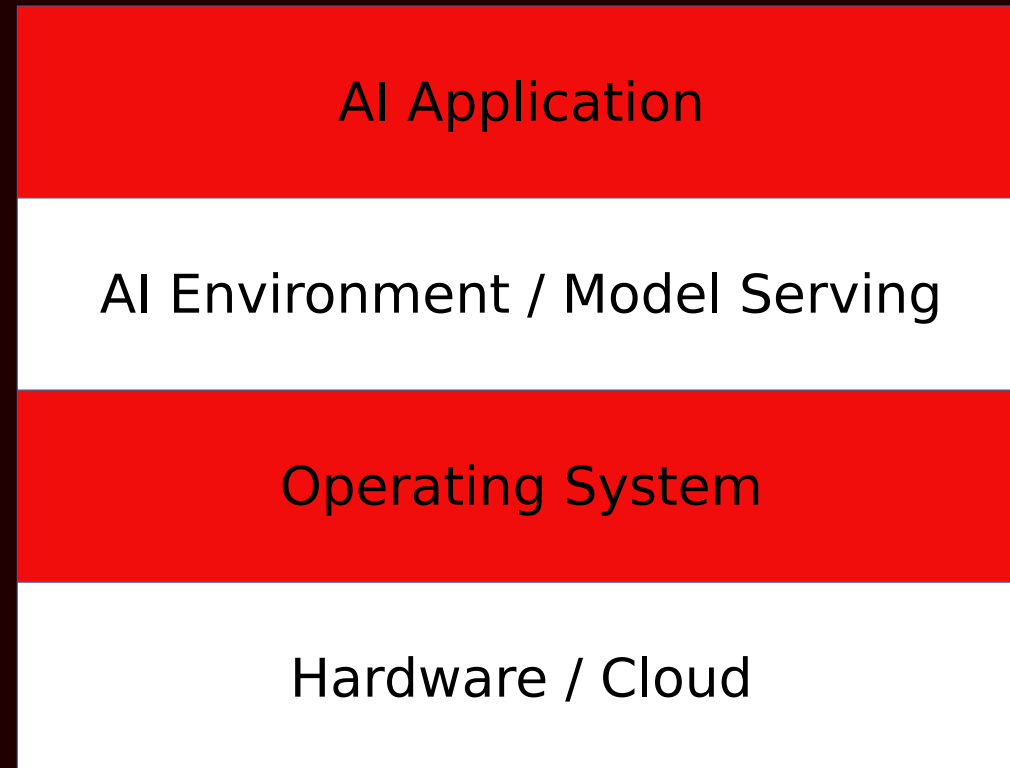# How to build you AI envoirment

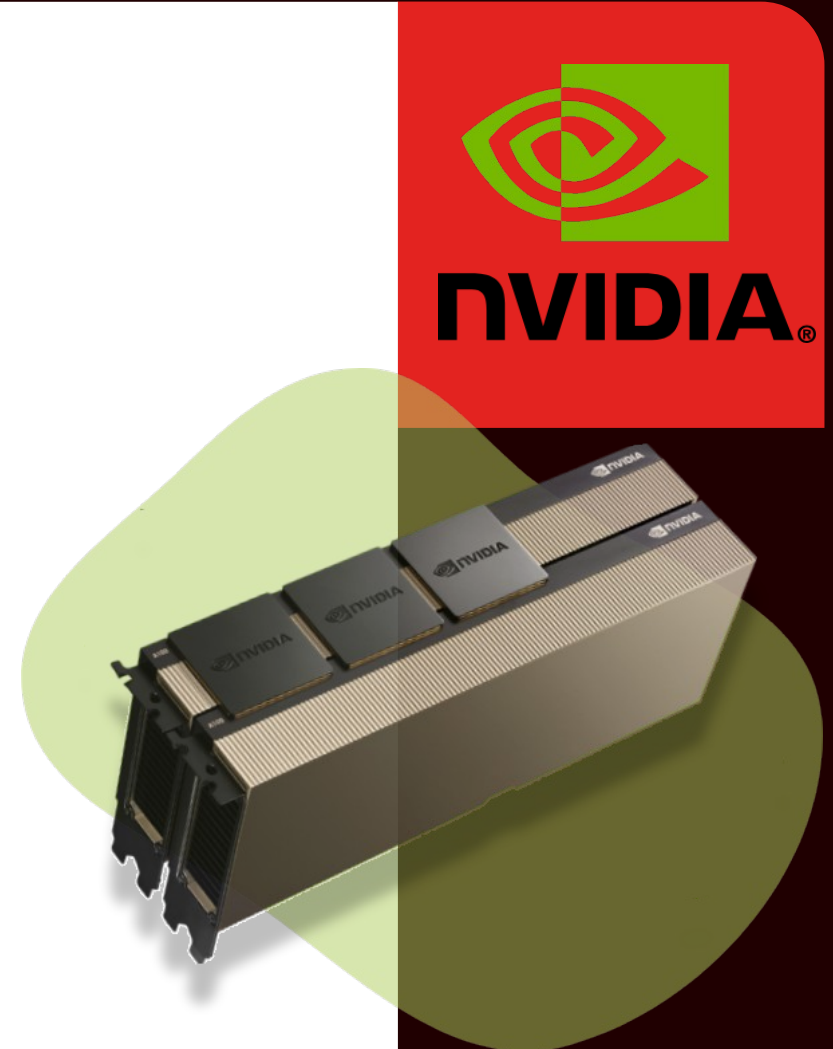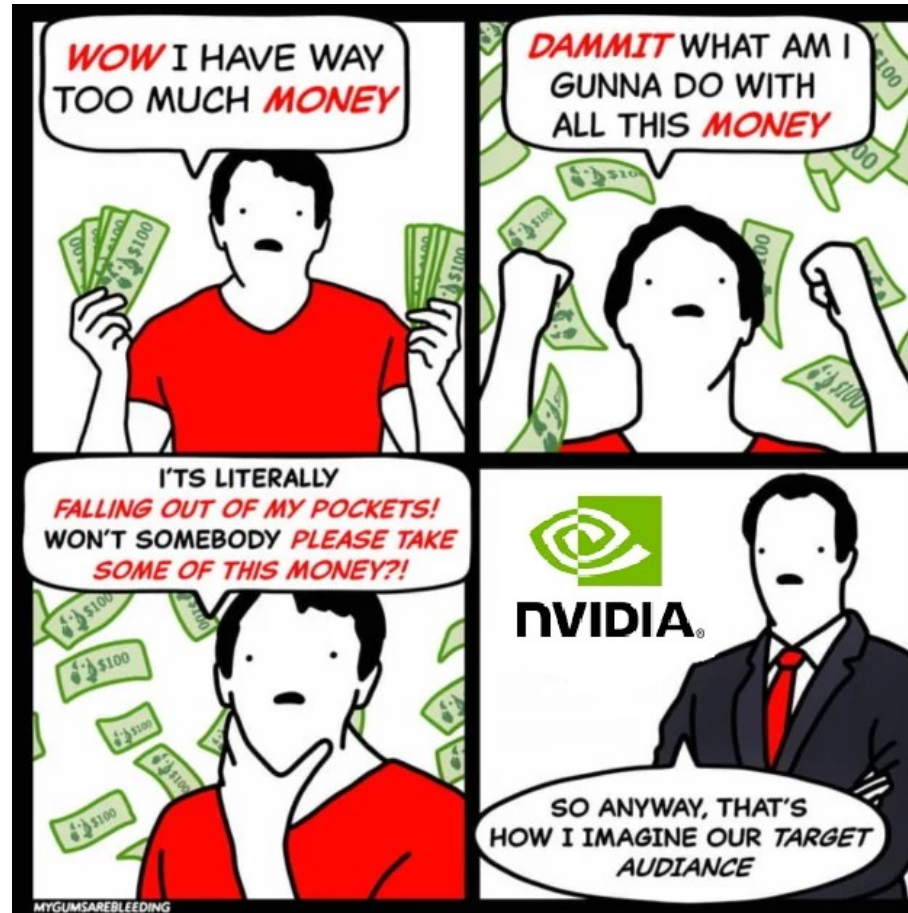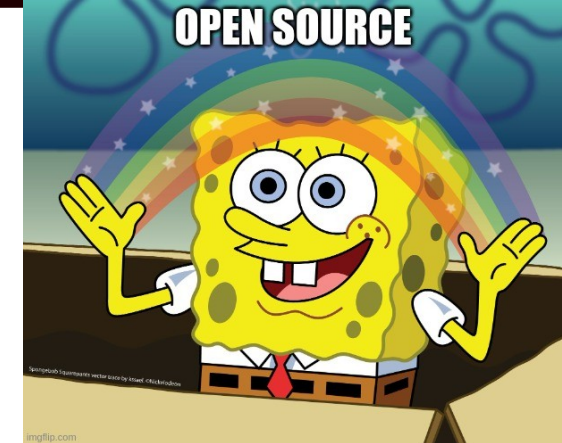| AI Application |
|:---:|
| AI Environment / Model Serving |
| Operating System |
| Hardware / Cloud |

# AI Accelerators

- A100
- H100
- H200
- Blackwell
- L40s

# Nvidia Nvidia Nvidia

# What are the challenges that are impeding your progress?
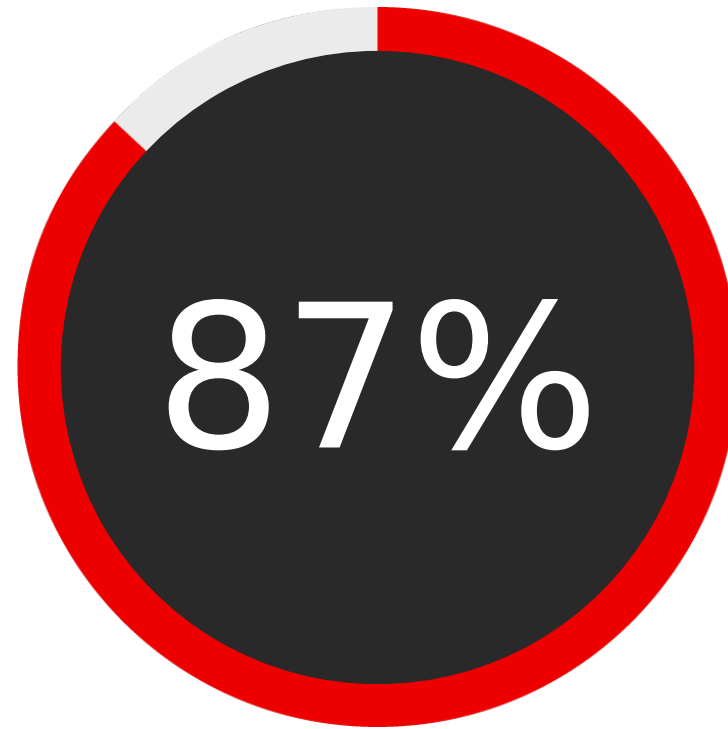
**Bridging environments**
How do you bridge the data scientist and app developer environments?

**Infrastructure complexity**
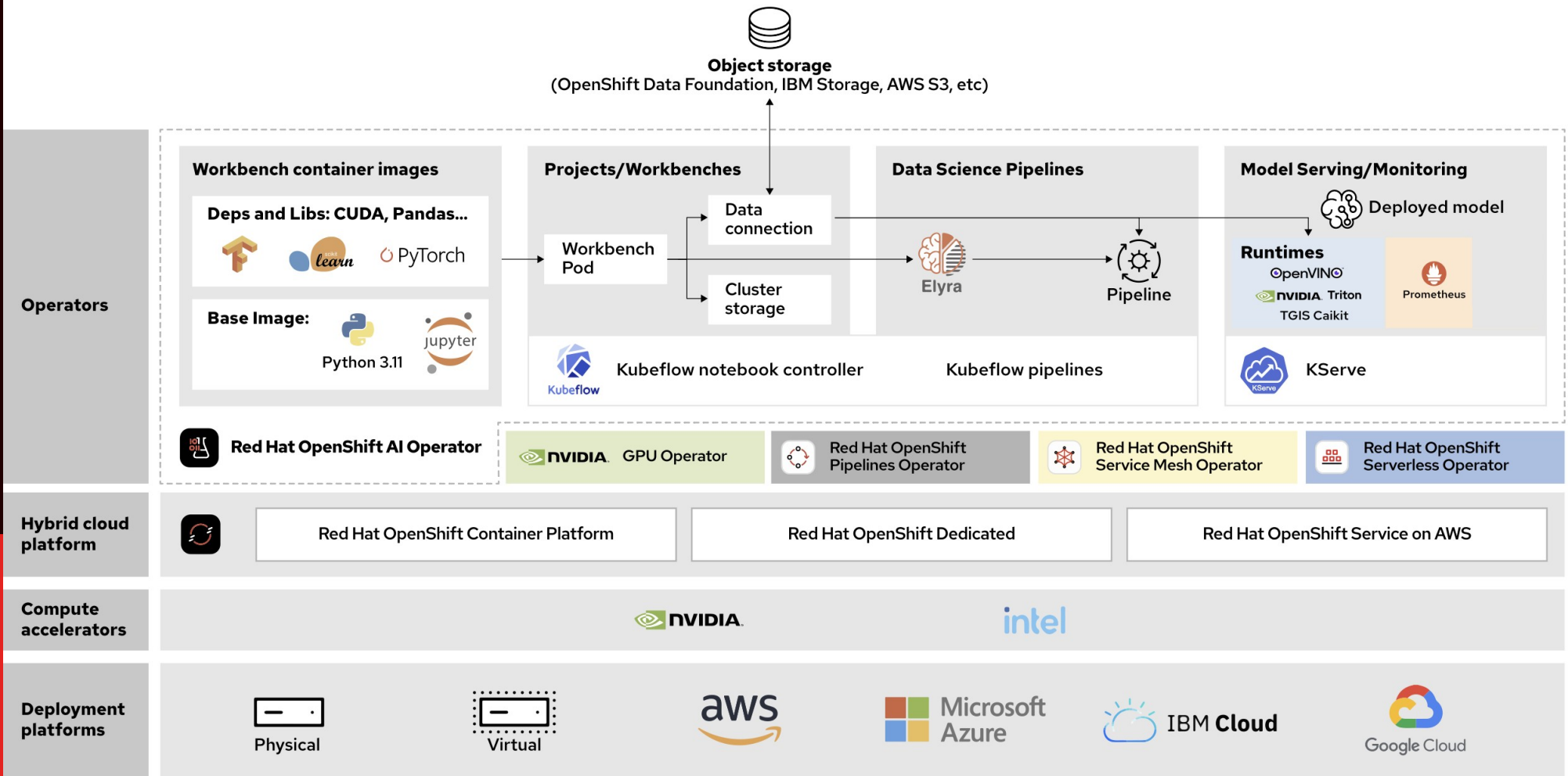How do you make the complexity of infrastructure invisible?

**Realizing value**
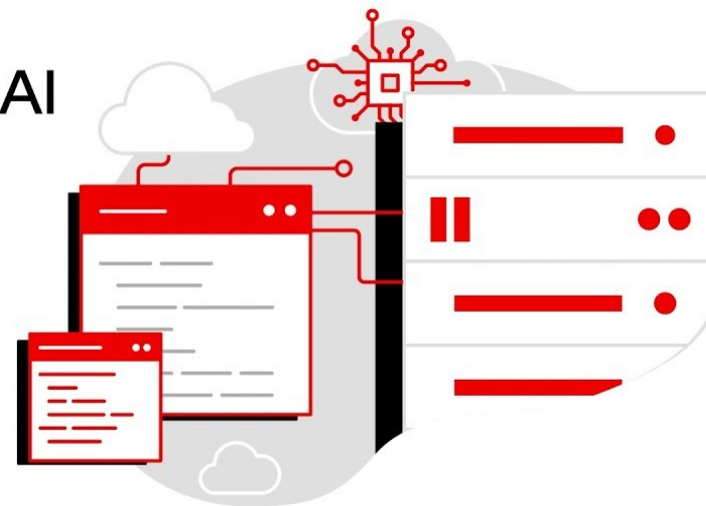How do you safely realize the business benefits of AI innovation?

## 87%

**87% of AI projects never make it into production**

**Object storage**
(OpenShift Data Foundation, IBM Storage, AWS S3, etc)

**Operators**

**Workbench container images**

**Deps and Libs: CUDA, Pandas...**

**Base Image:**
Python 3.11

**Projects/Workbenches**

Workbench Pod

Data connection

Cluster storage

Kubeflow notebook controller

**Data Science Pipelines**

Elyra

Pipeline

Kubeflow pipelines

**Model Serving/Monitoring**

Deployed model

**Runtimes**
OpenVINO
NVIDIA Triton
TGIS Caikit

Prometheus

KServe

**Red Hat OpenShift AI Operator**

NVIDIA GPU Operator

Red Hat OpenShift Pipelines Operator

Red Hat OpenShift Service Mesh Operator

Red Hat OpenShift Serverless Operator

**Hybrid cloud platform**

Red Hat OpenShift Container Platform

Red Hat OpenShift Dedicated

Red Hat OpenShift Service on AWS

**Compute accelerators**

NVIDIA

intel

**Deployment platforms**

Physical

Virtual

aws

Microsoft Azure

IBM Cloud

Google Cloud

**Red Hat OpenShift Container Platform**

**Red Hat OpenShift AI**

A **hybrid, unified platform** and **comprehensive set of tools**

**The Development Platform**

**Kubeflow**

VSCode

**Pipelines**

**Object storage**
(OpenShift Data Foundation, IBM Storage, AWS S3, etc)

**Operators**

**Workbench container images**

**Deps and Libs: CUDA, Pandas...**

**Base Image:**
Python 3.11

**Projects/Workbenches**

Workbench Pod

Data connection

Cluster storage

**Data Science Pipelines**

Elyra

Pipeline

**Model Serving/Monitoring**

Deployed model

**Runtimes**
OpenVINO
NVIDIA Triton
TGIS Caikit

Prometheus

Kubeflow notebook controller

Kubeflow pipelines

KServe

**Red Hat OpenShift AI Operator**

NVIDIA GPU Operator

Red Hat OpenShift Pipelines Operator

Red Hat OpenShift Service Mesh Operator

Red Hat OpenShift Serverless Operator

**Hybrid cloud platform**

Red Hat OpenShift Container Platform

Red Hat OpenShift Dedicated

Red Hat OpenShift Service on AWS

**Compute accelerators**

NVIDIA

intel

**Deployment platforms**

Physical

Virtual

aws

Microsoft Azure

IBM Cloud

Google Cloud

**Production Interference**

# Wrap Up

- Buy Some Hardware with Nvidia GPUs
- OpenShift for Scaleability
- Develop your LLM solution on top of the OpenShift AI platform
- Move quickly from development to Production with OpenShift AI endpoint
- Congratulations you have integrated Open Source AI that adds value to your Business

# We Are Piros

# Let's talk about your challenges

**Get in touch**

📞 CALL US ———
+32 (0)475 757 998

📍 FIND US ———
Gaston Geenslaan 11, 3001 Leuven
Piros.be