

從教室到產業， 平台思維的世代交替：

與Red Hat攜手共育次世代數位人才

Ato Lin | 林嘉彥

Sr. Solution Architect

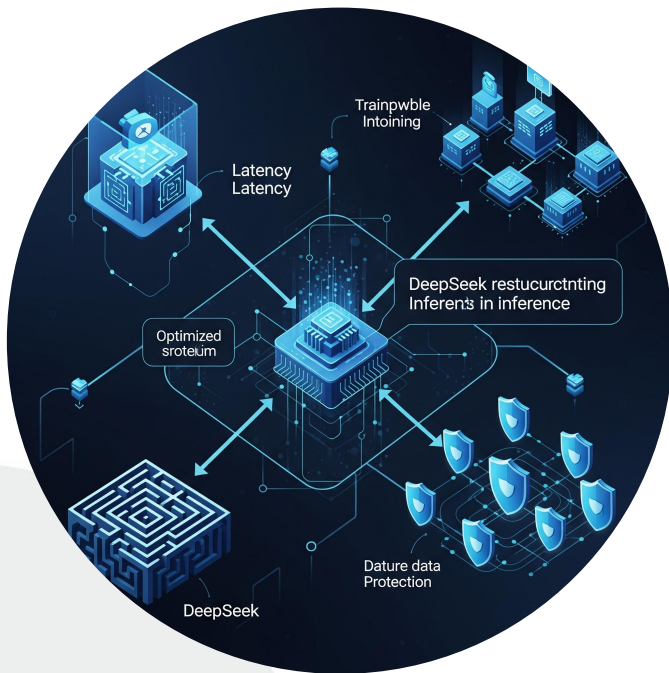
Red Hat

2025:「推理之年」 - AI 從訓練競賽轉向推理革命



面臨核心痛點：

- ▶ **推理效能低**：算力緊張、推理速度慢、延遲高、影響使用者體驗和業務效率
- ▶ **企業級安全與合規保障困難**：企業對模型私域部署有強烈需求，但在保障資料安全、隱私合規方面面臨挑戰
- ▶ **雲端化部署與彈性不足**：企業希望AI應用具備雲端化與彈性伸縮能力，但現方案在彈性與彈性方面不足



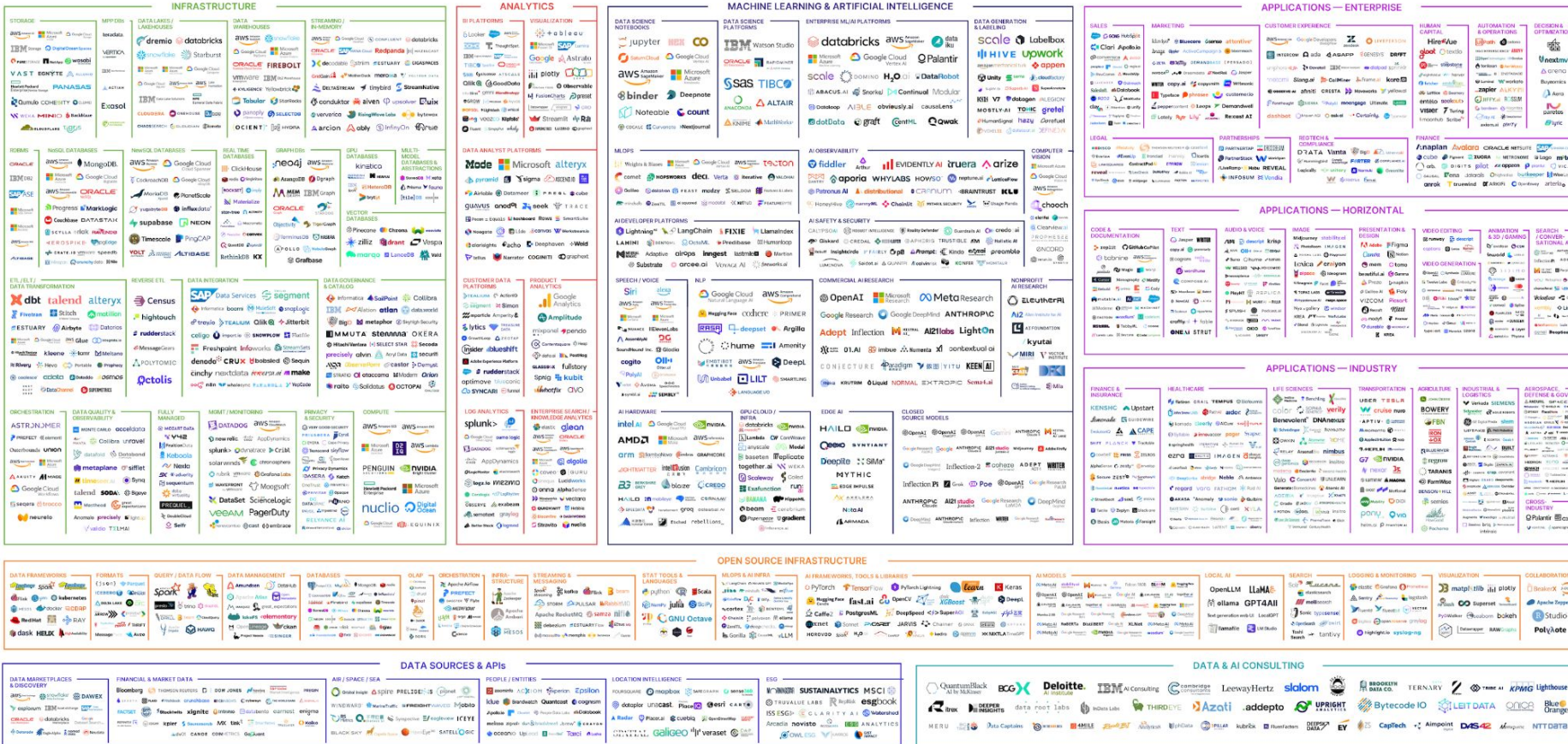
GPU 的使用狀況

- ▶ GPU 資源分配不均
 - 實驗室？各自採購？
- ▶ GPU 使用率低
 - 暑假, 不好申請, 管理？
- ▶ RTX 到高階卡 難以整合
 - 預算, 不同時期採購

人工智慧佈局

上游動態複雜性

THE 2024 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE



智慧應用

AI 工作負載

MLOps | LLMops | Inference (推論)

企業級 Kubernetes

AI 加速器

優化的作業系統

基礎架構



符合您需求的人工智慧基礎架構

協調、優化和規模化您的整個AI 生命週期

適合用途的模型

運用開源模型

使用小型語言模型進行最佳化

最大化 GPU 資源效率,
提升成本效益

降低複雜性

保持選擇彈性, 同時不損掌控力

部署統一的基礎架構, 以實現效能
規模化和優化成本

透過整合工具和自動化工作流程,
簡化 MLOps

提高靈活性

一次建置 AI, 隨處部署

將人工智慧融入您的數據

降低供應鏈風險

人工智慧平台

晶片多樣性 (GPU、TPU、NPU 等)

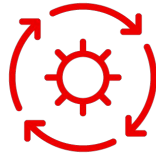
混合雲功能

使用紅帽建立人工智慧解決方案的價值

以自己的方式運用 AI



自帶模型



任何 GPU 上皆可運行
實現高性價比的 Gen AI



部署靈活性



優化使用 &
自助服務

Red Hat AI 提供經過驗證與量化，並附帶基準測試數據的第三方模型

vLLM 提供了更有效率的資源利用，這相當於需要更少的 GPU 來處理 LLM，並且對硬體有廣泛的支援。

透過統一的 AI 平台，以一套技能即可在本地端、任何雲端及邊緣環境，運用數據並部署模型。

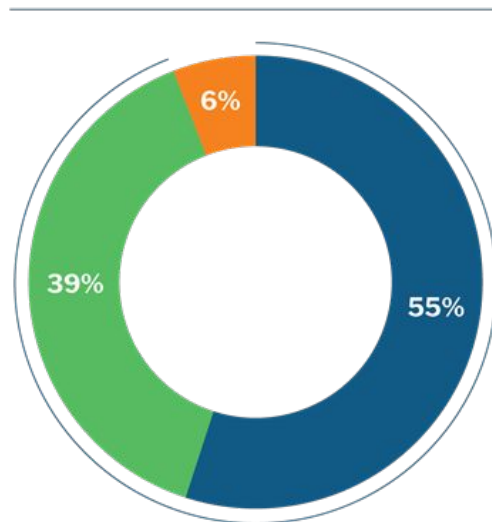
共享資源以最佳化 GPU 等硬體的使用率，並讓資料科學家能透過自助服務快速佈建實例。



AI 採用者熱衷於使用容器來改善 AI 工作負載

在生產環境中 使用容器部署機器學習模型

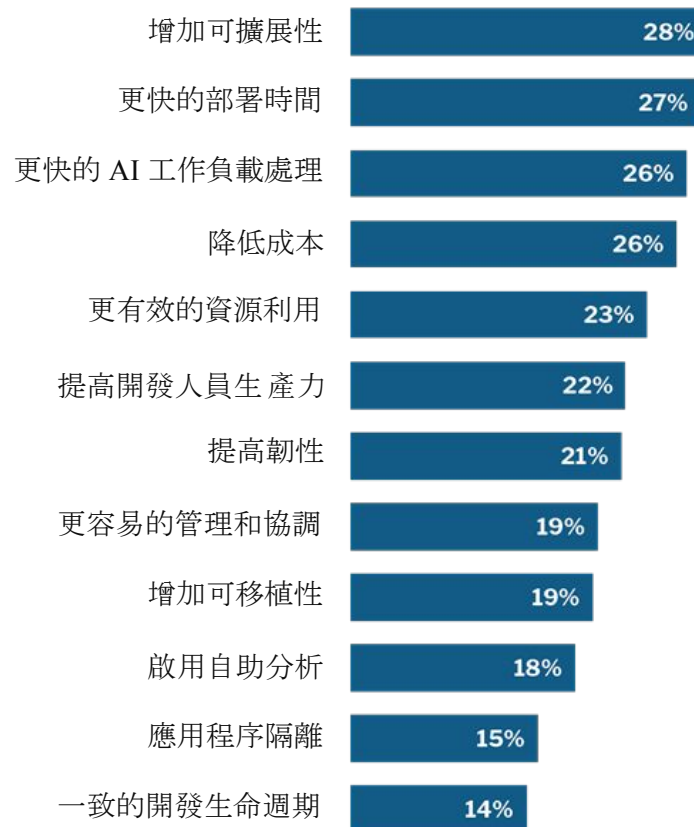
94% of AI adopters are using or plan to use containers within one year'



- Currently using
- Not currently using but plan to within the next year.
- Not currently using containerized environments and no plans in next year

容器對 AI 採用者的好處

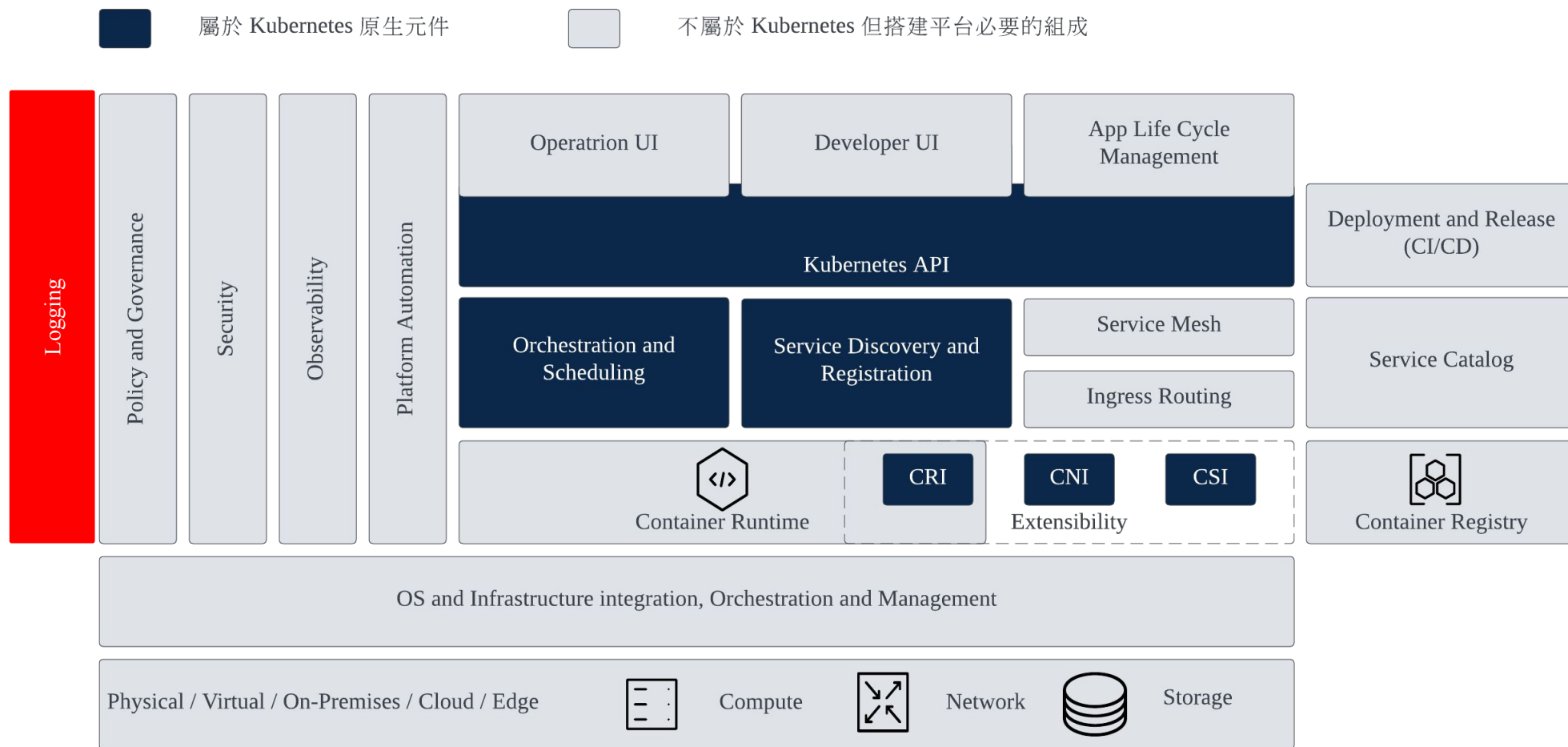
Benefits of Containers



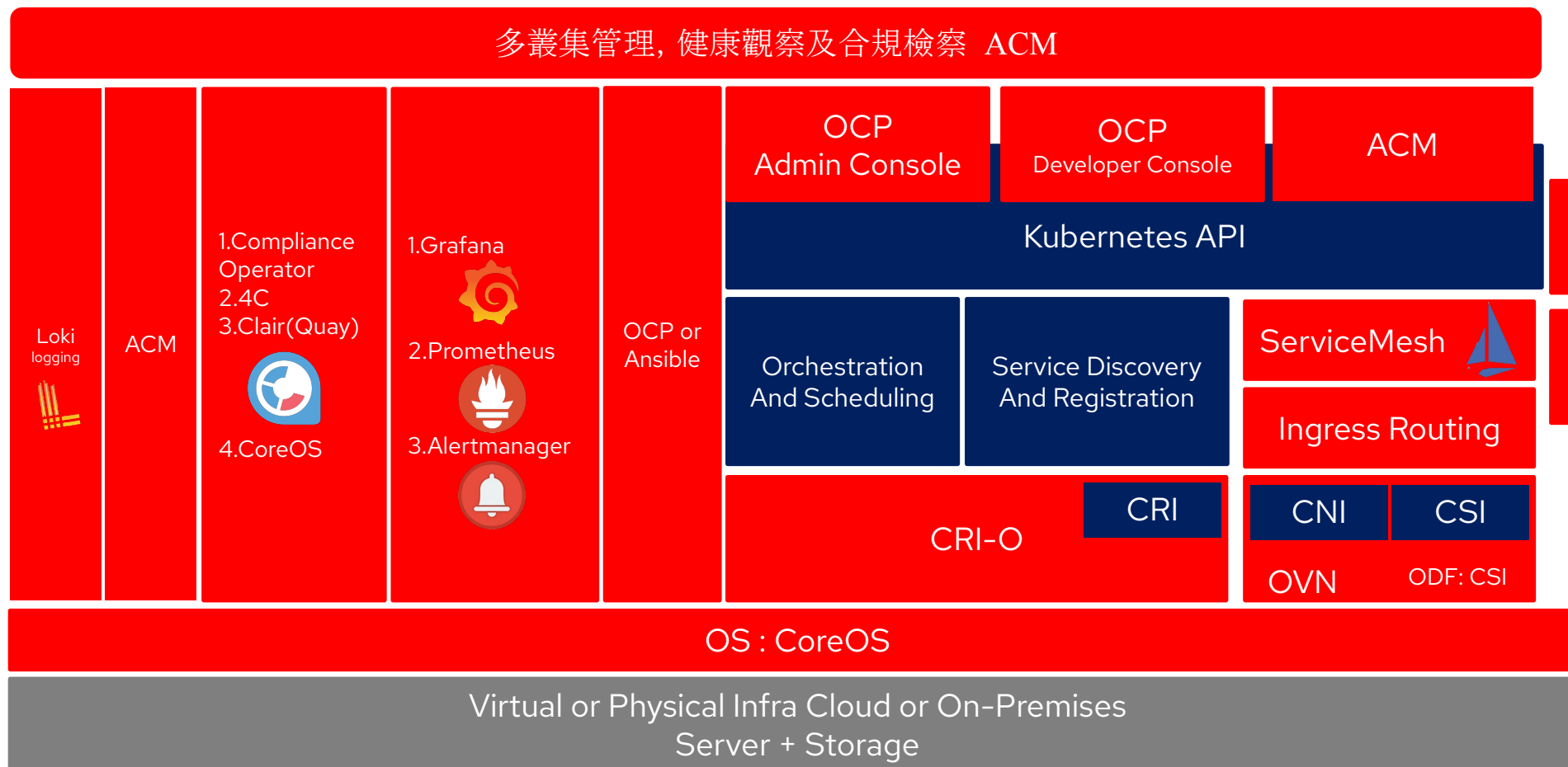
94%

的AI採用者正在用或
者計畫在一年內使用容器

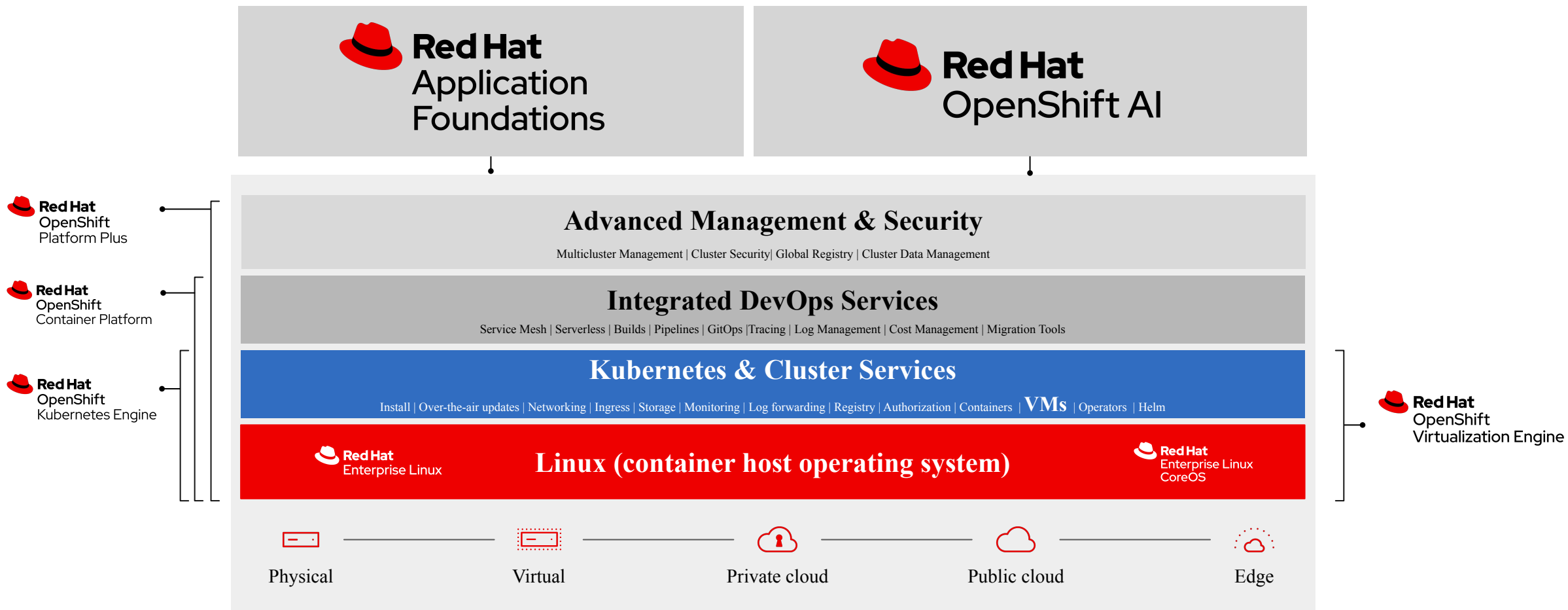
Kubernetes 元件 (Gartner 建議 k8s 平台需具備的元件)



OpenShift 元件說明 (紅帽提供的元件)



一站式企業級平台解決方案



vLLM: AI 推理的 Linux

Apps



Hardware

Models



Hardware

~500,000

downloads in a typical week

Source: "PyTorch+package" PyPI Stats, May 2025

生成式和預測式人工智慧的統一平台



值得信賴、一致且全面的基礎



硬體加速



Physical



Virtual



Private cloud



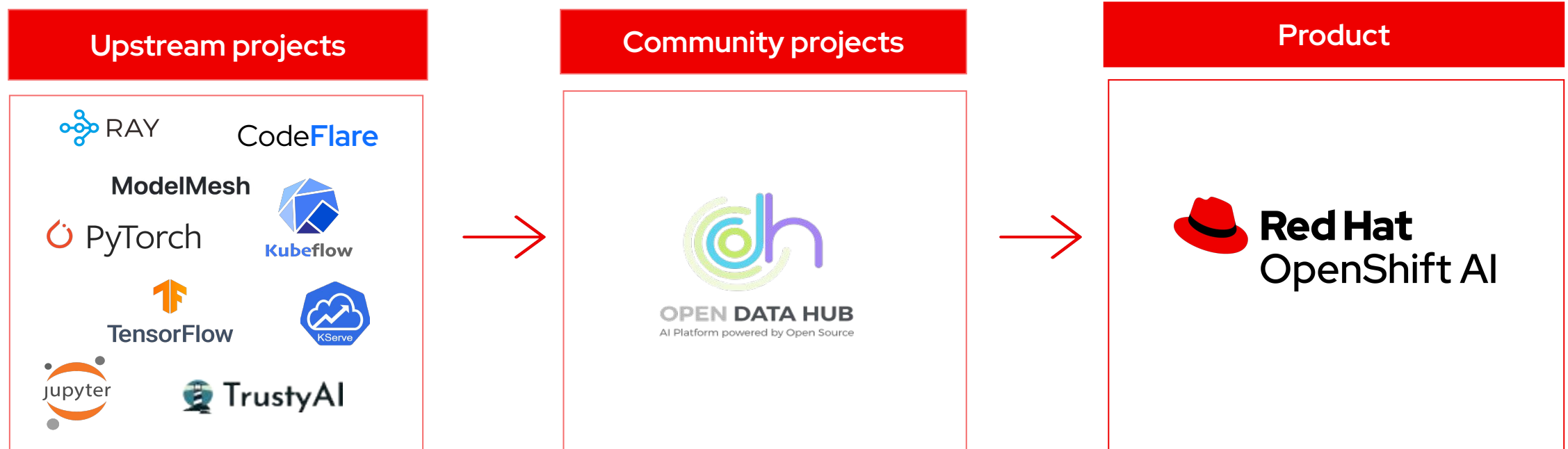
Public cloud



Edge



Red Hat's AI/ML engineering is 100% open source



靈活高效的推理

- ▶ 分散式推理 (llm-d)
- ▶ 新的經過驗證和最佳化的模型
- ▶ vLLM 增強功能
- ▶ LLM 壓縮器

將模型連接到數據

- ▶ 模組化且可擴展的方法：資料收集、合成資料產生、調優、評估。
- ▶ RAG 增強功能和合作夥伴集成
- ▶ Feature Store



**Red Hat AI
Enterprise**

Agentic AI

- ▶ AI 體驗：AI Hub 和 Gen AI Studio
- ▶ Gen AI Studio 中的 MCP 支援
- ▶ Llama Stack API 整合

AI Platform

- ▶ 模型目錄和模型倉儲
- ▶ Model as a Service 增強功能和 API 管理整合
- ▶ GPU as a Service 功能增強

單一平台可在任何加速器、任何雲端上運行任何模型

Ensuring efficient GPU sharing across users

Scenario / Use case

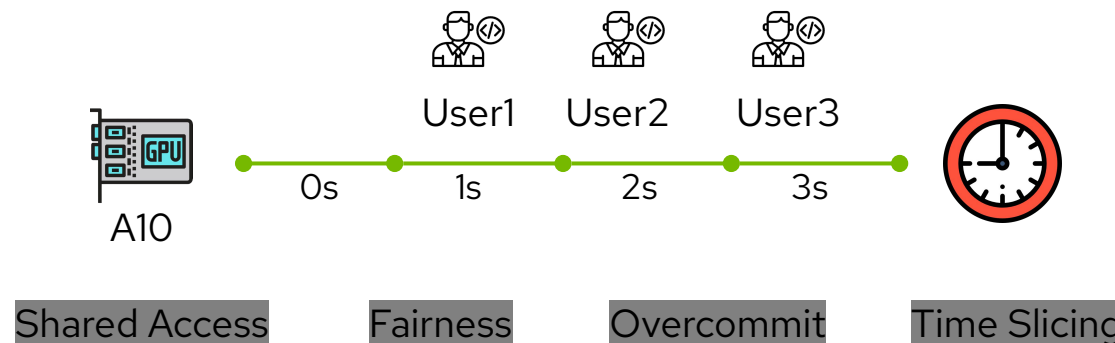
Early stage of development
(Limited computing resources)

Common Challenges

One team's workload monopolizes GPU memory, affecting others.

Solution

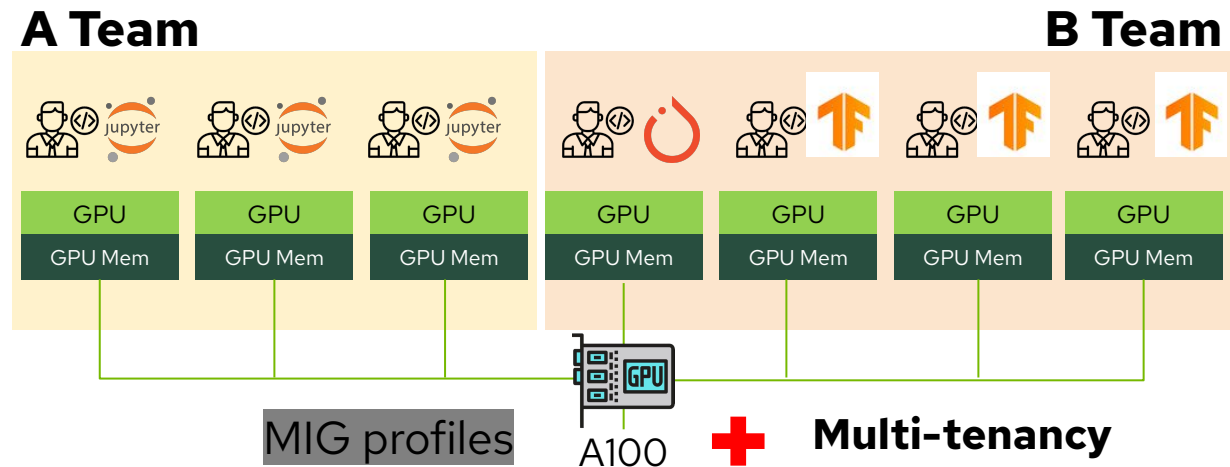
Time Slicing & Multi Tenancy



Ensuring GPU resource isolation in a multi-team environment

Scenario / Use case	Common Challenges	Solution
Multiple internal teams within a company (e.g., AI, Cloud, MLOps)	One team's workload monopolizes GPU memory, affecting others.	Multi Instance GPU & Multi Tenancy

GPU Isolation

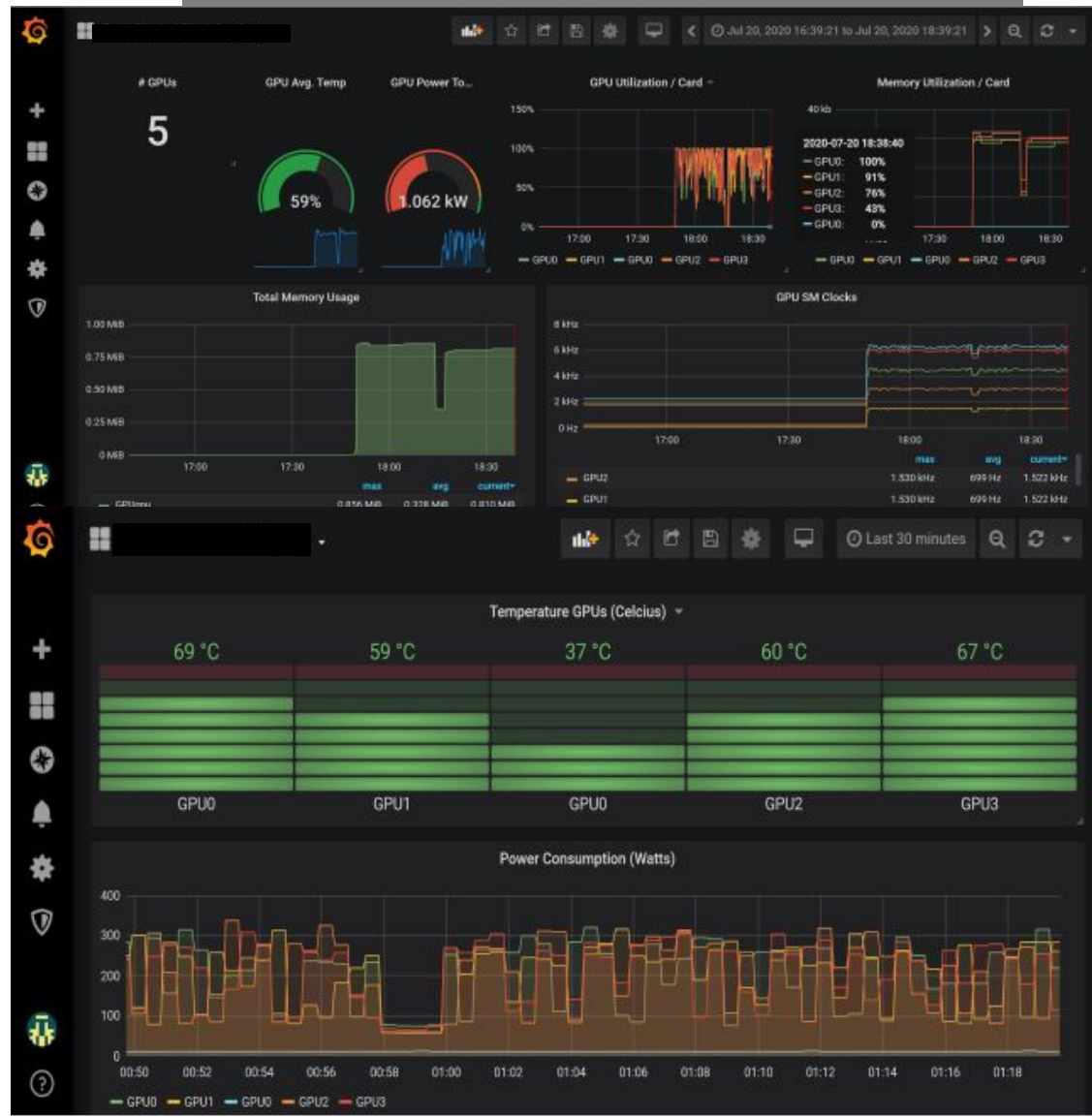
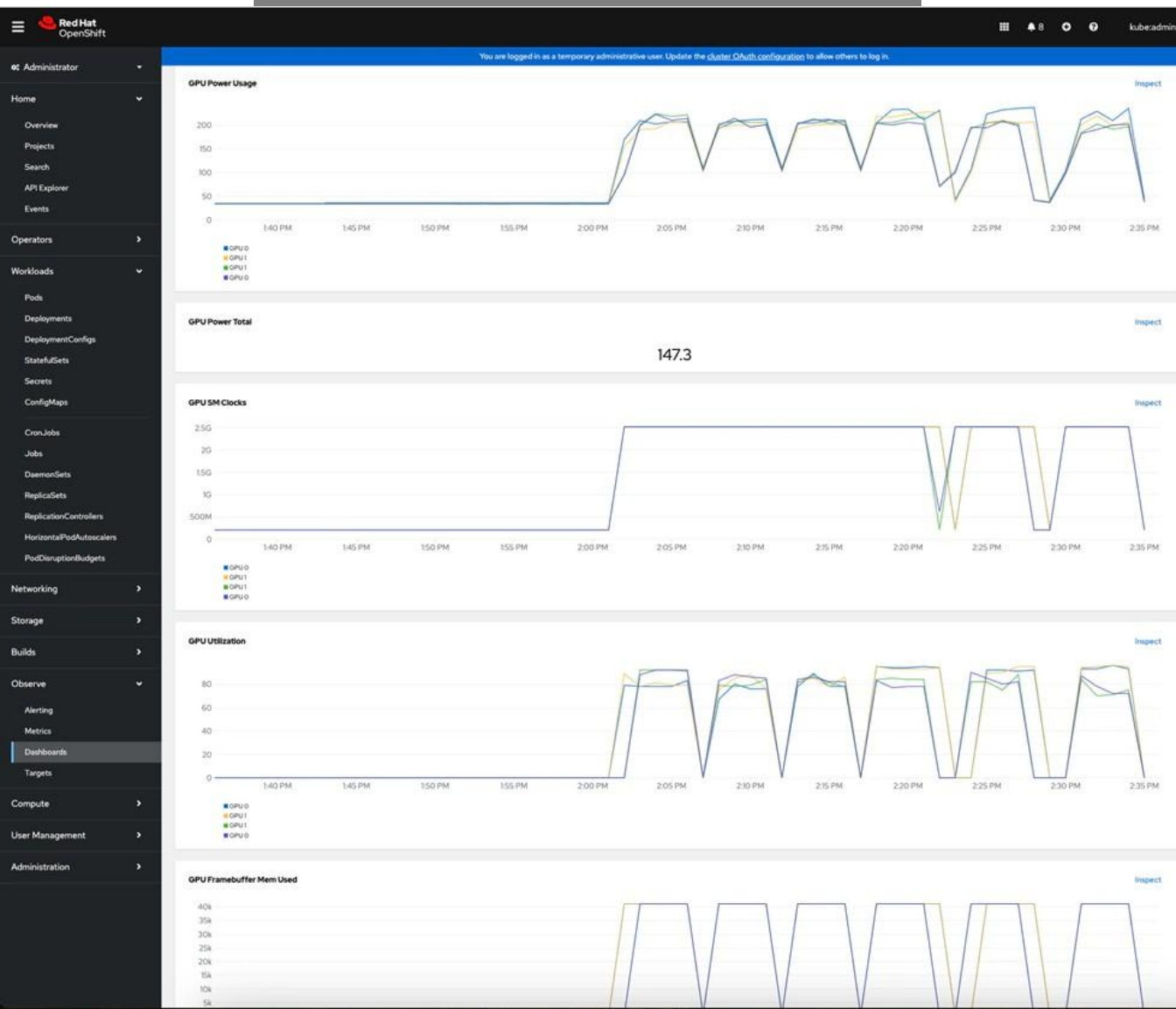


Fair sharing

OpenShift GPU 資源監控

整合 NVIDIA DCGM Operator

整合客戶自有 Grafana Dashboard
(預設支援 NVIDIA DCGM Exporter)



自助式的模型開發環境

透過 OpenShift AI 的平台，數據科學家自己可以申請和配置所需要的模型開發環境。

Red Hat OpenShift AI

Applications

Enabled

Explore

Data Science Projects

Data Science Pipelines

Distributed Workload Metrics

Model Serving

Resources

Settings

Notebook Server Administration

Start a notebook server

Select options for your notebook server.

Notebook image

- Minimal Python 2024.1
Python v3.9
[Versions](#)
- CUDA 2024.1
CUDA v12.1, Python v3.9
[Versions](#)
- TensorFlow 2024.1
CUDA v12.1, Python v3.9, TensorFlow v2.15
[Versions](#)
- HabanaAI 2024.1
Python v3.8, Habana v1.13
[Versions](#)
- Standard Data Science 2024.1
Python v3.9
[Versions](#)
- PyTorch 2024.1
CUDA v12.1, Python v3.9, PyTorch v2.2
[Versions](#)
- TrustyAI 2024.1
Python v3.9
[Versions](#)
- code-server 2024.1
Python v3.9
[Versions](#)

Deployment size

Container Size

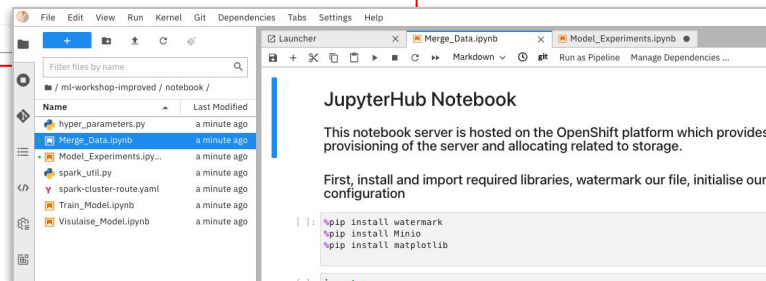
Small

Environment variables

Add more variables

Start server Cancel

Start server in current tab



Notebook image

預設支援常用ML 開發框架(Framework) 等容器映像檔。可透過創建自訂義映像檔*, 您也可以新增所需的函式庫(Library), 或使用Jupyter Notebook 以外的編輯器 (e.g. VSCode)。

Deployment size

指定所需的資源(CPU/Memory) 來啟動容器

透過設定預設的停止時間, 可以自動停止非活躍的環境並回收資源

Number of GPUs

加入指定 GPU 數量資源

Number of GPUs

0

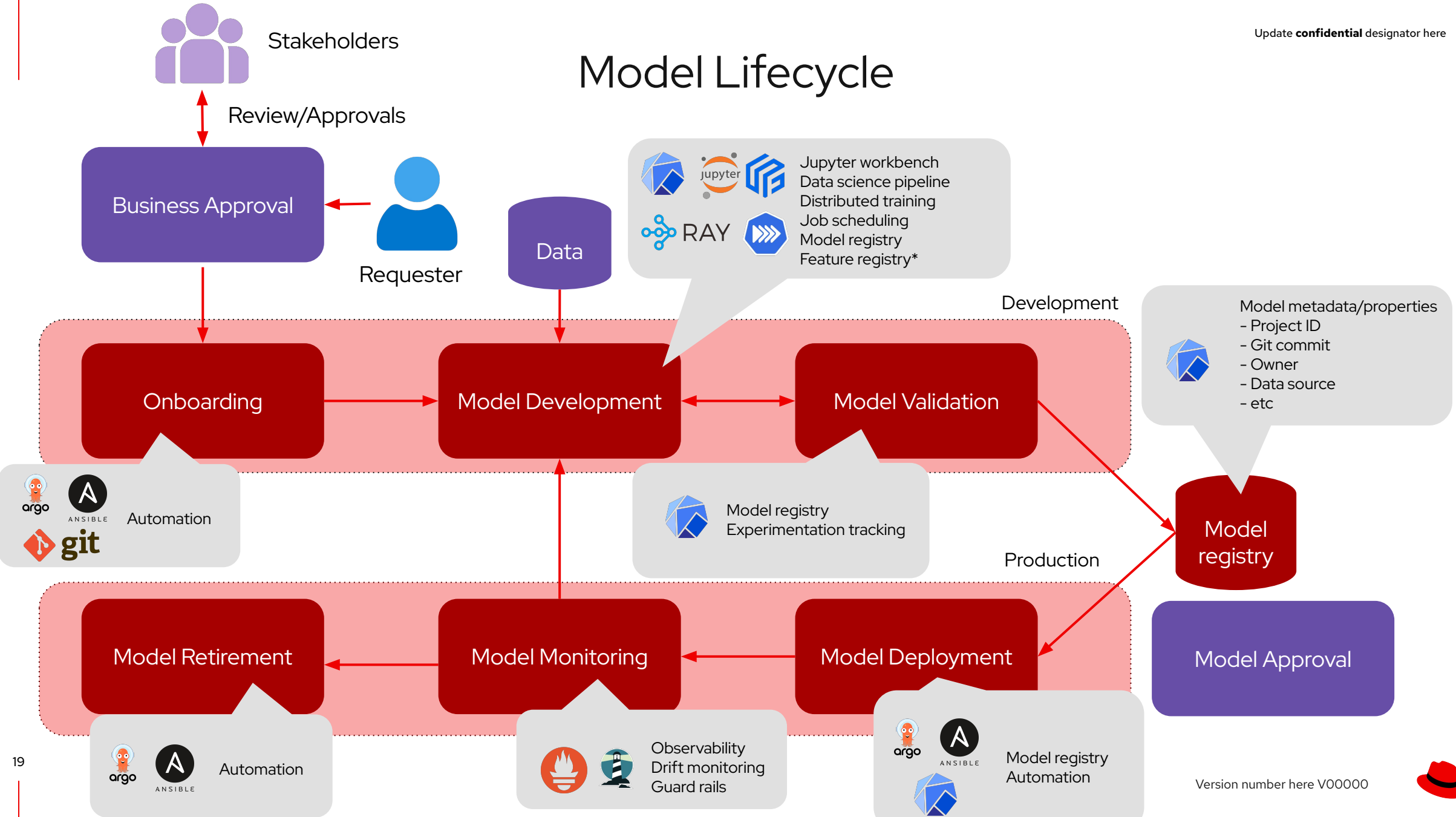
1

2

3

4

Model Lifecycle



將傳統 VM 遷移進 OpenShift Virtualization

虛擬機遷移

只需幾個簡單的步驟建立轉移計畫(Migration Plan)
，即可將虛擬機器大規模遷移到 OpenShift
Virtualization。

目前支援來源端的虛擬平台有：

- Red Hat Virtualization
- KVM
- OpenStack
- VMware vSphere
- OVA File
- OpenShift Virtualization

The image displays two screenshots from the OpenShift OperatorHub interface. The top screenshot shows the 'OperatorHub' page with a search for 'Migration Toolkit for Virtualization' and a red box highlighting the operator card. The bottom screenshot shows the 'Create migration plan' configuration screen for the Migration Toolkit for Virtualization operator, with a red box highlighting the 'Review' step in the progress indicator.

OperatorHub

Project: vmexamples

Discover Operators from the Kubernetes community and Red Hat partners, curated by Red Hat. You can purchase commercial software through Red Hat Marketplace developers. After installation, the Operator capabilities will appear in the Developer Catalog providing a self-service experience.

All Items

Migration Toolkit for Virtualization Operator provided by Red Hat

Facilitates migration of VM workloads to OpenShift Virtualization

Create migration plan

1 General
2 VM selection
3 Network mapping
4 Storage mapping
5 Type
6 Hooks
7 Review

Review the migration plan

Review the information below and click Finish to create your migration plan. Use the Back button to make changes.

Plan name: move-webtriv

Source provider: rhvcnv

Target provider: host

Target namespace: vmexamples

This is a new namespace that will be created when the plan is started.

Migration transfer network: Pod network

Selected VMs: 1

Network mapping: Source networks

Target namespaces / networks: Pod network

Migration Toolkit for Virtualization

Migration plans > my_second_migration_plan

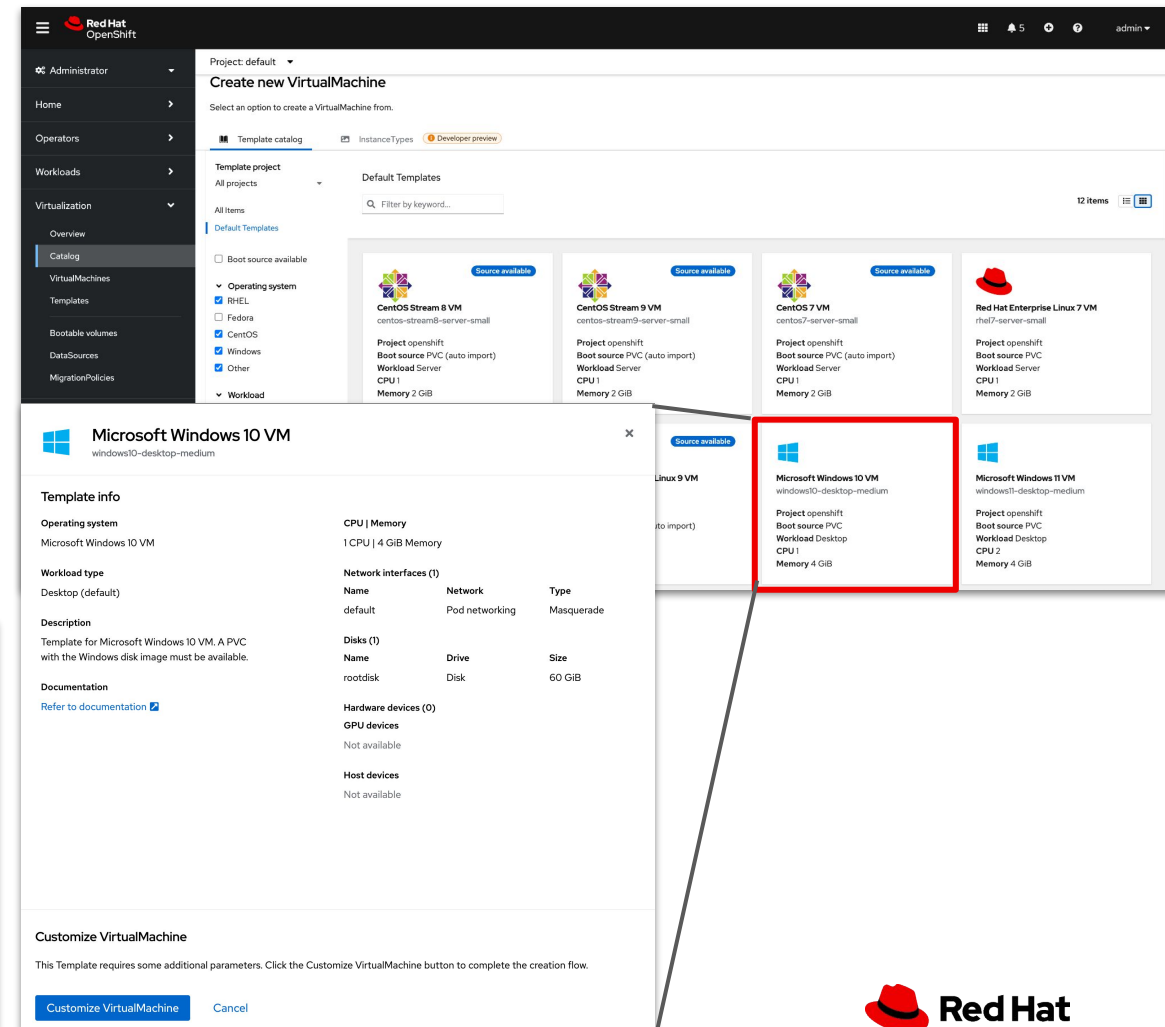
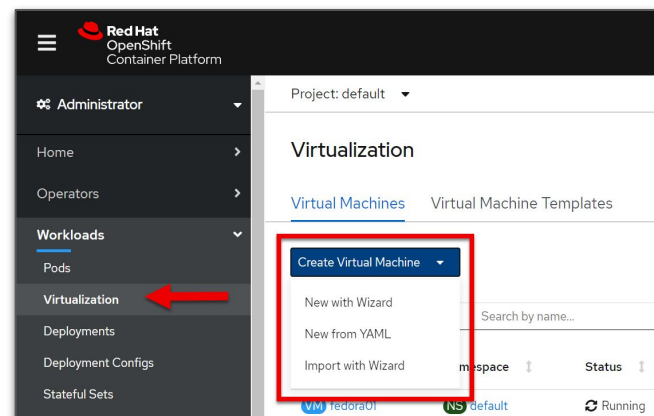
Migrations

Type	Start time	End time	Data copied	Status
Migrate	09 Aug 2019, 11:34:56		91 / 123 GB	In progress 42%
Stage	08 Aug 2019, 8:19:11	08 Aug 2019, 22:33:44	87 / 87 GB	Complete 100%

OpenShift 次世代容器平台提供運行虛擬機的能力

OpenShift Virtualization 虛擬化功能

- 容器平台內建功能
- 可透過使用者介面(GUI)簡化並創建 VM 虛擬機
- 介面整合系統資源、網路和儲存配置選項
- 支持 Linux 和 Windows 虛擬機器作業系統
- 可從 VMware vSphere 遷移虛擬機器



FSI

台灣規模 Top 3 金控集團之一

客戶挑戰

2024 年被公認為人工智慧元年。該集團已確認利用生成式 AI 技術為客戶提供更全面的金融服務，包括金融詐騙防範、AI 智能客服、貸款/承保智能協助等。但因人員技能尚未完備、故不願意花時間建構和維護 AI 相關的基礎設施和平台工具，而是期待一個開箱即用的企業級 AI/ML 解決方案，提供穩定且有效率的 GPU 資源管理和跨團隊 AI 專案協作平台。

解決方案

提供 Red Hat OpenShift AI 和 GPS 團隊實施的諮詢服務，以協助該集團標準化 AI 專案的基礎執行環境，並集中現有 OpenShift 叢集上的 AI 周邊應用程式管理。

Why Red Hat

該集團的痛點是缺乏完整的 AI/ML 解決方案、和與 AI 相關的技能知識，這意味著開發人員無法專注於服務創新，從長遠來看將減緩業務成長。

這兩個問題都可以透過 Red Hat RHOAI 和專業服務來解決，與競爭對手相比，它具備企業級支持、且運行於經市場驗證過穩定性和安全性的 OpenShift 平台。

效益

Red Hat 協助該集團建構了集團內第一個完整的 AI/ML 平台，並提供教育和培訓，使 IT 團隊掌握維運技能。這使得開發團隊從維護 AI 工具中解放出來，並使他們能夠專注於應用服務開發。

1. AI 項目必要環境建置時間減少 70%。
2. 將 AI 服務開發、測試和發佈時間縮短 30%。

軟體與服務

Red Hat OpenShift Platform Plus
Red Hat OpenShift AI

Red Hat Consulting





簡化 AI 模型服務部署，搭配專業
在地技術支援，加速創新應用落地
，有效降低導入成本與複雜度。

客戶挑戰

現有 AI 環境多為分散建置，缺乏統一管理，導致維運複雜、效率不彰。部分單位選擇台製品牌導入時，常出現穩定性不足或無法使用的情況，需一個高效穩定且容易管理的平台，來簡化 AI 導入流程、提升效能，確保政策順利落實。

解決方案

完整的 AI 平台和專業服務，協助單位快速建置可支援多種 AI 模型(如 TAIDE、Mistral、Gemma3)的統一應用平台，並維持企業級的技術支援。解決快速部署與穩定性的需求，使 AI 應用得以穩定且迅速上線，顯著縮短導入時程，並簡化日後維運管理，大幅提升系統整體穩定性與運行效率。

Why Red Hat

紅帽解決方案的關鍵在於其穩定、安全與可持續的開源創新能力。支援多種主流與在地 AI 模型，且具備完善的技術服務。平台設計便於管理與擴展，協助組織有效推動 AI 應用落地，降低部署與營運上的複雜度。

效益

- 單一平台整合: 提供多元且彈性擴展的 AI 模型服務，快速滿足各種應用需求。
- AI 模型服務: 透過簡易網頁介面部署與操作，降低 AI 模型服務部署技術導入門檻。
- 專業服務在地團隊: 提供完整技術支援，簡化後續管理的複雜度。

軟體與服務

Red Hat OpenShift Platform Plus
Red Hat OpenShift AI

Red Hat Consulting
AI Innovation LAB

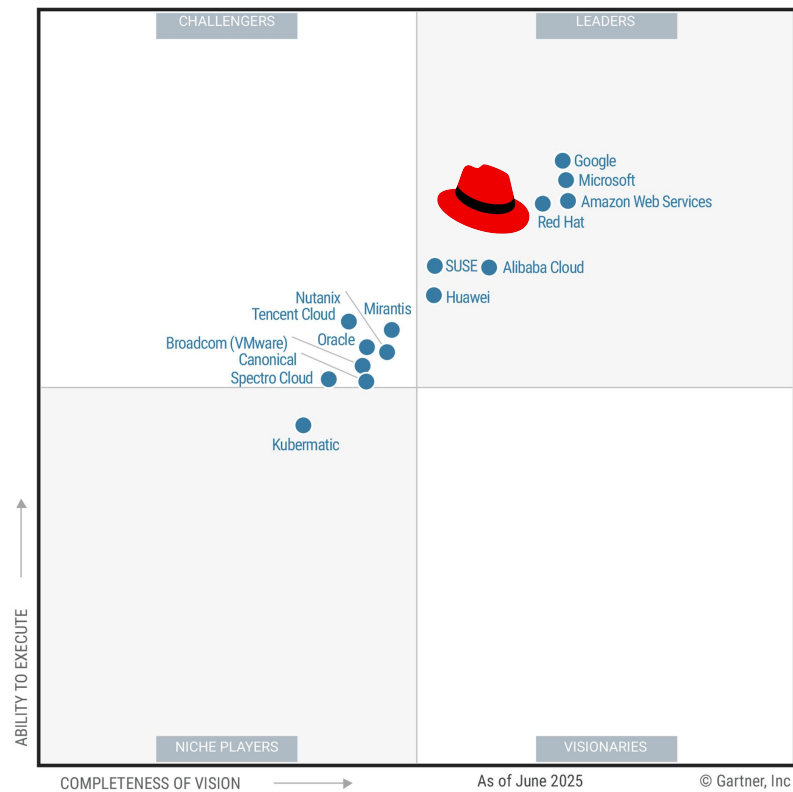


2025 Market Ranking

Gartner June 2025 Container management

Forrester Wave Q3 2025 Multicloud Container Platforms

Figure 1: Magic Quadrant for Container Management

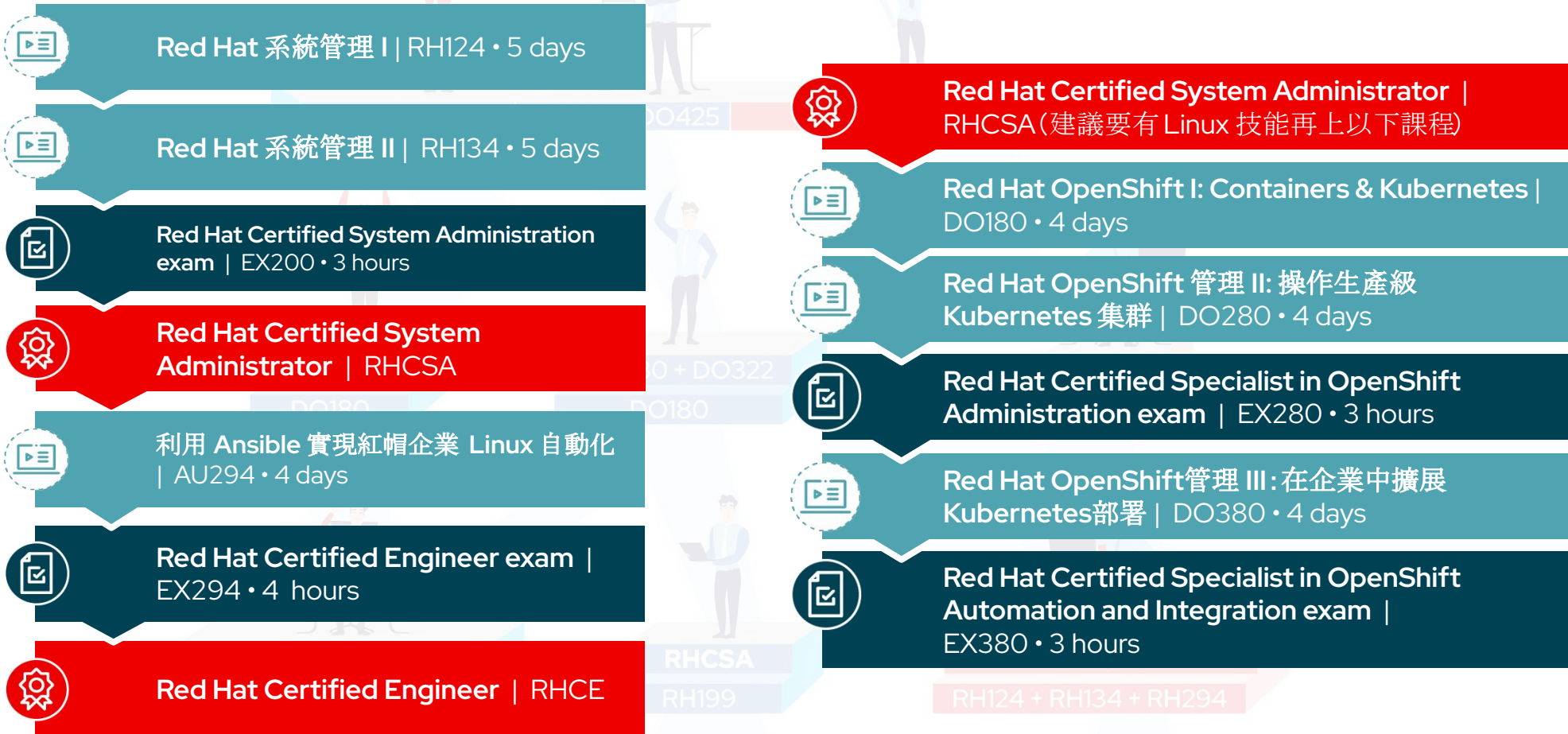


Gartner

THE FORRESTER WAVE™
Multicloud Container Platforms
Q3 2025



完整的教育訓練, 從 Linux 到 容器平台 到 AI



為學生準備就職就緒程度及留任率



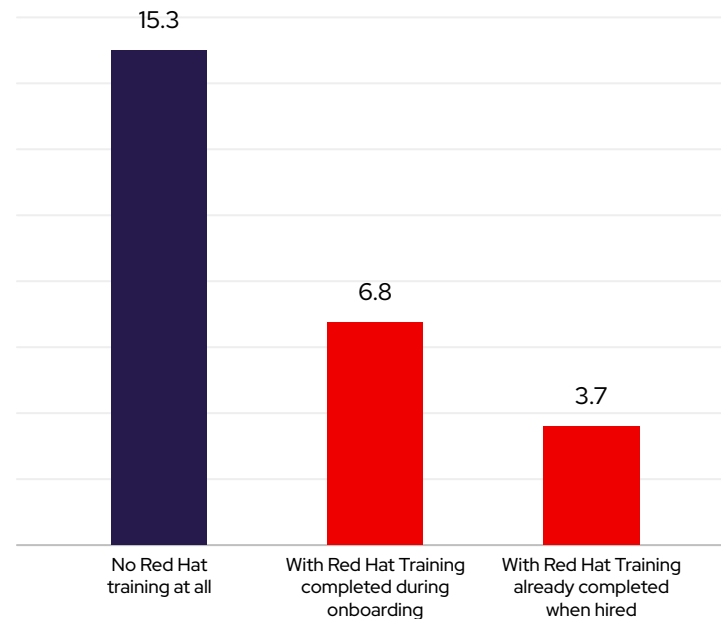
工作準備就緒速度提高 **76%**(上任之前受訓)
when a new hire had already completed Red Hat training prior to onboarding.



工作準備就緒速度提高 **55%**(新人訓受訓)
where new team members completed Red Hat training as part of the onboarding process.



任期延長 **8%**
Red Hat-trained employees have 8% longer tenure on average than untrained staff members.



達到完全生產力之時間(週數)

Demo

快速部署一個模型服務

Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.

 [linkedin.com/company/red-hat](https://www.linkedin.com/company/red-hat)

 [youtube.com/user/RedHatVideos](https://www.youtube.com/user/RedHatVideos)

 [facebook.com/redhatinc](https://www.facebook.com/redhatinc)

 twitter.com/RedHat