

Red Hat與AMD攜手合作 打造企業AI應用平台

Joe Yu

首席架構師經理

Red Hat, Taiwan



世界發生了什麼事？



Chat GPT



GitHub Copilot

生成式 AI

一組新的模型可以從自然語言輸入中產生一些東西

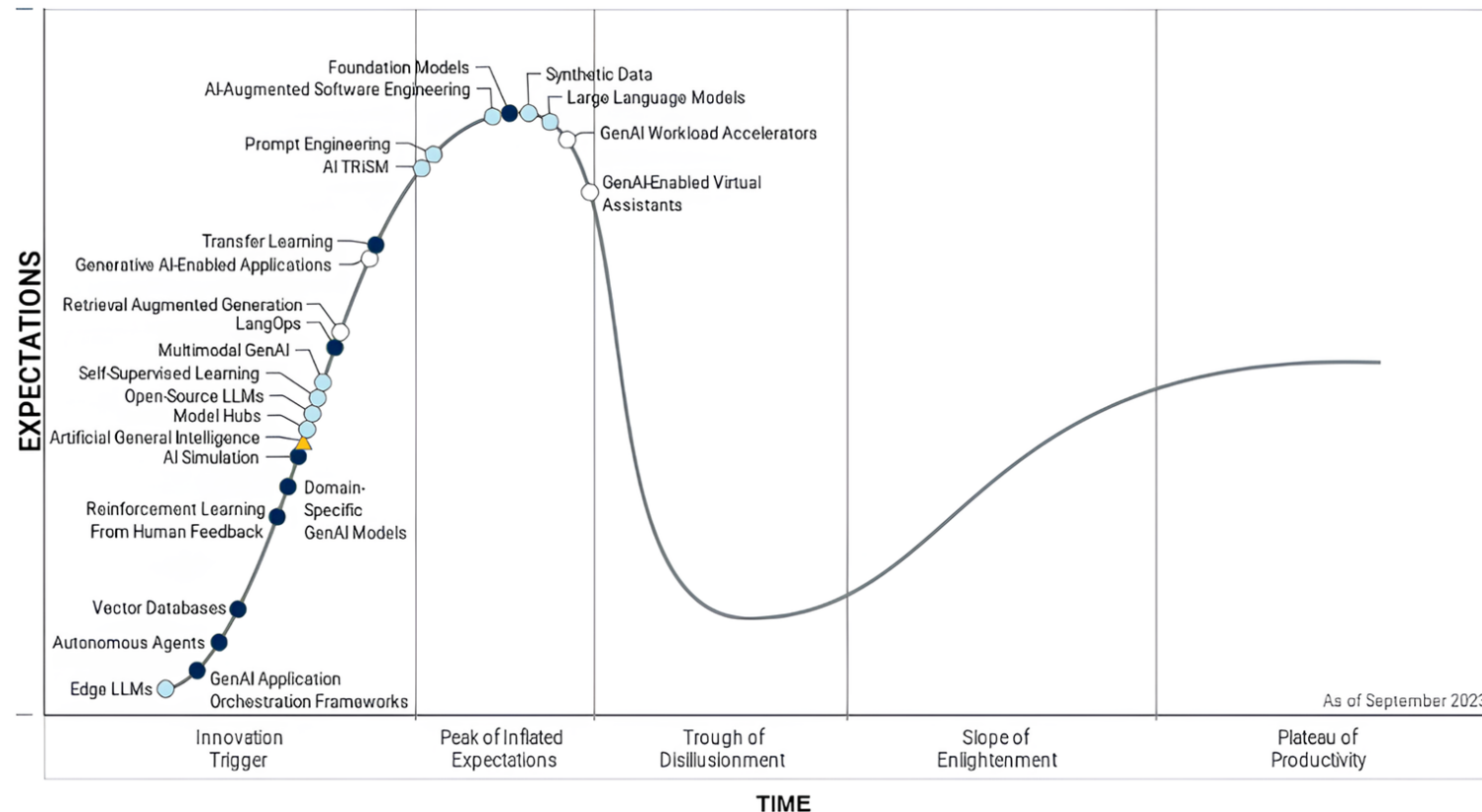
特別是LLMs

收集網路規模的文字資料集來訓練能夠令人信服地預測成文字的東西

Red Hat OpenShift AI (RHOAI)

Red Hat 新解決方案名稱

隨著時代的興起，AI的導入歷程



Plateau will be reached: ○ <2 yrs. ● 2-5 yrs. ● 5-10 yrs. ▲ >10 yrs. ⊗ Obsolete before plateau

Gartner

“ More Than **80%** of Enterprises Will Have Used Generative AI APIs or Deployed Generative AI-Enabled Applications by 2026. ”

Gartner






































<https://www.gartner.com/en/newsroom/press-releases/2023-10-11-gartner-says-more-than-80-percent-of-enterprises-will-have-used-generative-ai-apis-or-deployed-generative-ai-enabled-applications-by-2026>

AI/ML技術與方案組成呈現開源多樣化

AI + MLOPS

平台

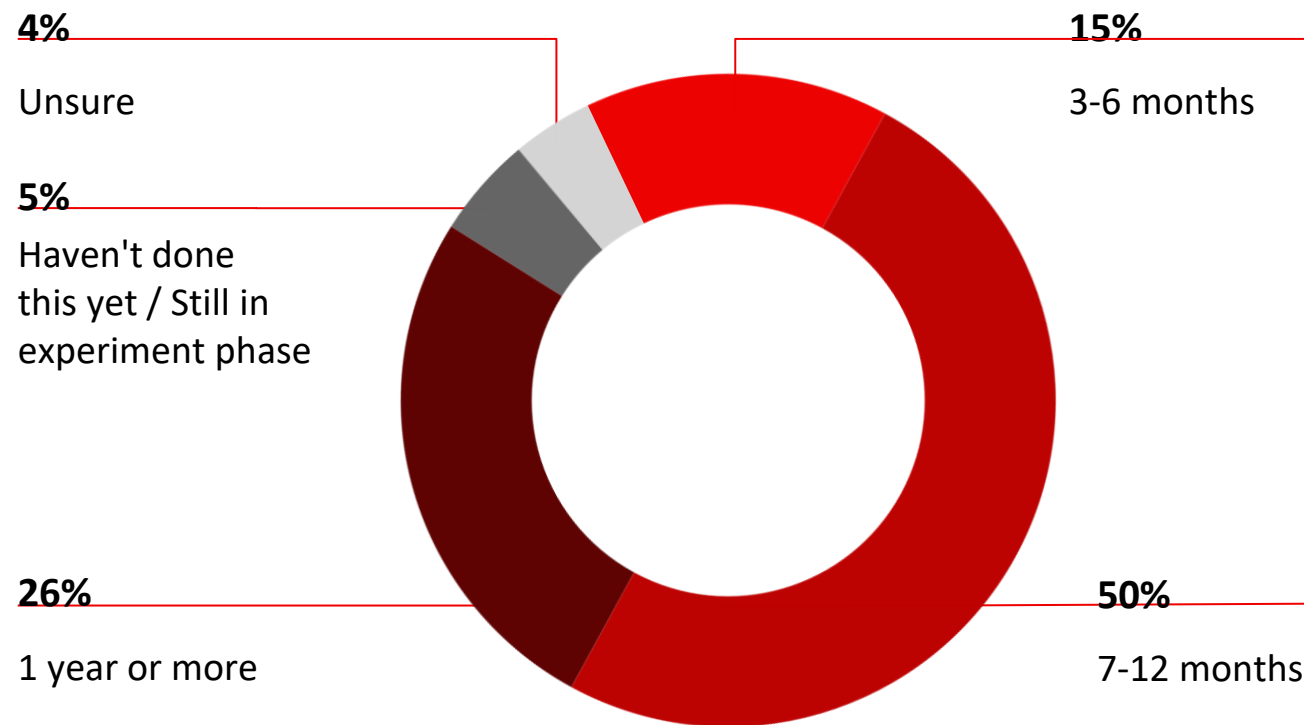
應用程式平台

機器學習函式庫	   XGBoost  NumPy  pandas  OpenCV  LangChain
開發語言與工具	 python  R  jupyter  
資料視覺化、標籤與處理	 Superset  Label Studio  trino  Pachyderm  Apache Spark  Apache Airflow  RAY  FEAST  Katib  Kubeflow  mlflow  Flask  KServe 實驗 & 模型生命週期
軟體定義存儲	 ceph  MINIO  kafka  Istio  3scale <small>by Red Hat</small> 整合
大規模容器管理	 etcd  K  Prometheus  HELM  Grafana  TEKTON  argo 自動化軟體測試交付
容器化與資源編配	 docker  podman  CoreOS  kubernetes
程序調配與硬體加速	 LINUX + AI Accelerator

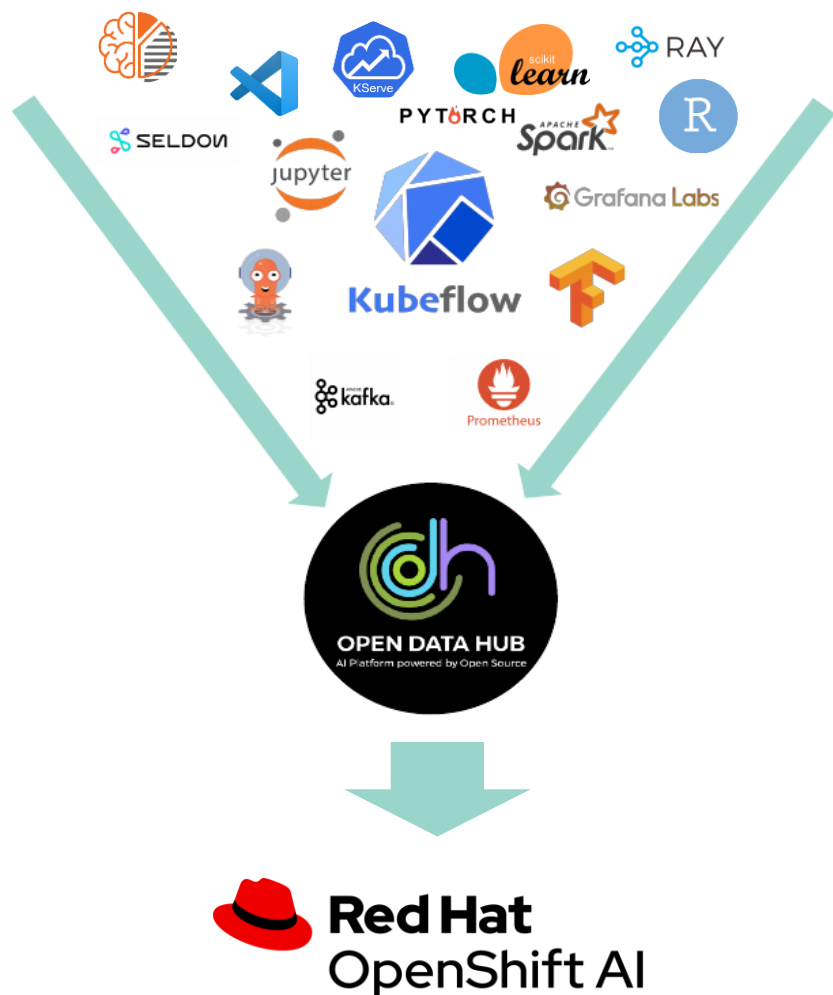
實施AI是一個充滿挑戰的過程

AI/ML 從想法到模型實施的時間平均需要多少？

一半的受訪者 (50%) 表示，他們從構思到實施模型的平均 AI/ML 時間為 **7-12 個月**。



讓企業可以專注 - Red Hat將 AI/ML的開源創新賦予企業支持



模型開發

使用Red Hat提供或自有的image，在 JupyterLab 中進行探索性實驗，並存取核心 AI/ML 程式庫和框架，包括 PyTorch 和 Tensorflow



模型服務和監控

跨任何雲端、代管與自我維運的 OpenShift 以部署和擴展模型，並集中監控其效能。



生命週期管理

建立用於模型訓練和驗證的可重複資料科學pipeline，並將其與 DevOps機制整合，以便在整個企業中交付模型。



提昇能力/協作

建立專案並在團隊之間共用。結合Red Hat套件、開源以及 ISV 認證軟體。

RHOAI: 提供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。

The screenshot displays the 'Start a notebook server' configuration page in the Red Hat OpenShift AI console. It features a sidebar with navigation options like 'Applications', 'Data Science Projects', and 'Settings'. The main content area is titled 'Start a notebook server' and includes sections for 'Notebook image' and 'Deployment size'. The 'Notebook image' section lists various pre-configured environments, each with a 'Versions' link. The 'Deployment size' section shows a 'Small' container size. Below these sections, there are 'Start server' and 'Cancel' buttons. A Jupyter logo is positioned to the right of the configuration options. At the bottom, a preview of a Jupyter Notebook interface is shown, displaying code for installing dependencies and importing libraries.

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

可以加入指定 GPU 數量資源

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

*: 客製化的映像檔需要由客戶自行管理。

RHOAI: 充分使用OpenShift平台上的GPU資源

Start a notebook server

Select options for your notebook server.

Notebook image

Minimal Python ⓘ
Python v3.8

PyTorch ⓘ
Python v3.8, PyTorch v1.8, CUDA v11.4

Standard Data Science ⓘ
Python v3.8

TensorFlow ⓘ
Python v3.8, TensorFlow v2.7, CUDA v11.4

CUDA ⓘ
Python v3.8, CUDA v11.4

Deployment size

Container size

Small ▼

Number of GPUs

0 ▼

Number of GPUs

1 ▼

0

1 ✓

[+ Add more variables](#)

Number of GPUs

0 ▼

0 ✓

1

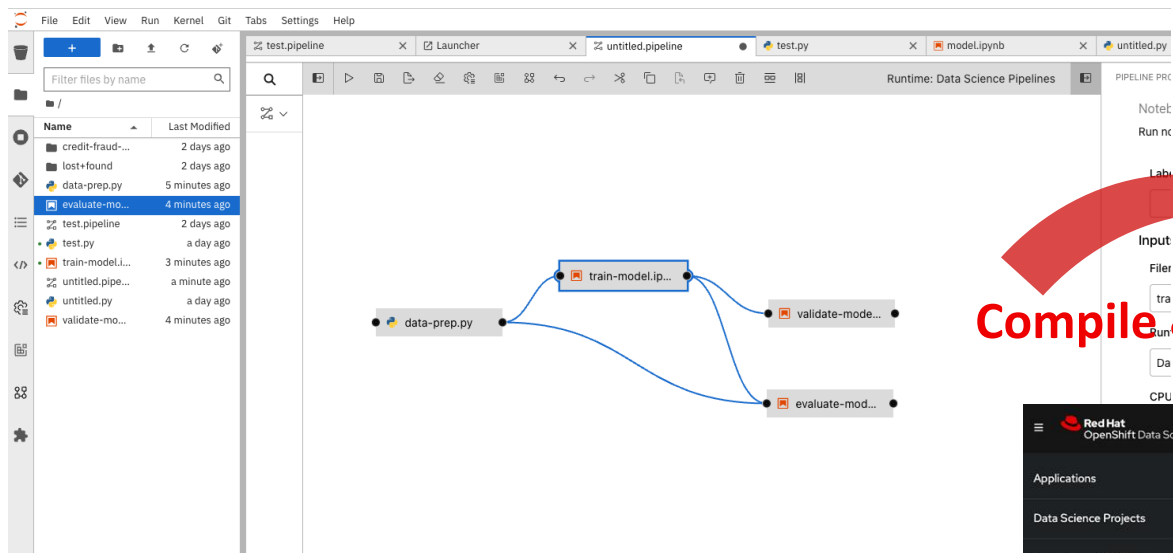
2

3

4

RHOAI: 支援透過流水線 (Pipeline) 實現流程自動化

把寫好的模型程式碼組合起來，建立一條流程管線 (Pipeline)，從而自動化多個和模型重新訓練或部署相關的步驟。



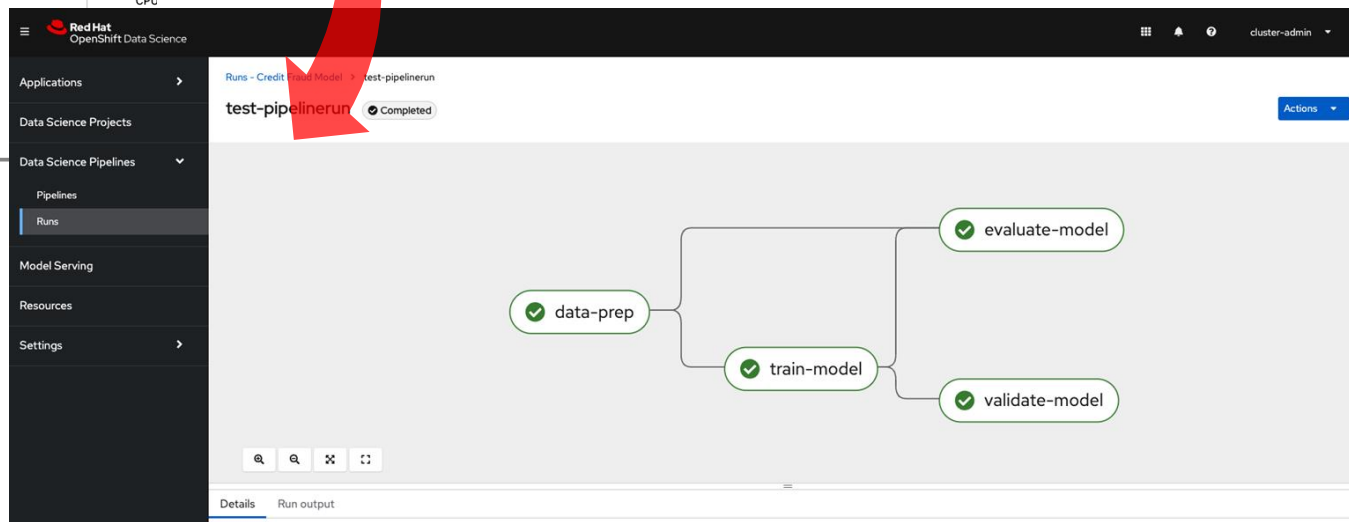
Compile & Run

Process Automation

透過自動化流程管線的執行，可以極大地減少 AI/ML 模型重新訓練和部署所需的手動操作。

Pipeline Creation

你可以基於已有的 Python 腳本或 Notebook 檔案，透過圖形使用者介面 (GUI) 來組建流程管線。



RHOAI: 簡化模型服務 (Mode Serving) 提供

透過 OpenShift AI Portal 網站部署已訓練的模型，我們提供能夠以API形式呼叫的終端點。



Deploy model

將儲存在物件儲存空間的已訓練模型，進行容器化後部署對應後續服務提供。



Model serving
Manage and view the health and performance of your deployed models.

Name Find by name Deploy model

Model name	Project	Inference endpoint	Status
credit card fraud	Credit Fraud Model	https://credit-card-fraud-credit-fraud-model.apps.mycluster2.88o5.p1.openshiftapps.com/v2/models/credit-card-fraud/infer	✓

Inference endpoint

將模型轉換成 API 形式，並提供外部應用程式調用的推論(inference)端點。

Red Hat OpenShift提供最佳的應用程式平台支持並完整整合RHOAI

Figure 1: Magic Quadrant for Container Management



Red Hat OpenShift 在 2024 Gartner® Magic Quadrant™ for Container Management研究報告中

- 延續2023年，持續處於**領導象限**中
- 最具有願景(Vision)

The Forrester Wave™: Multicloud Container Platforms, Q4 2023

- 多雲容器平台**名列第一**

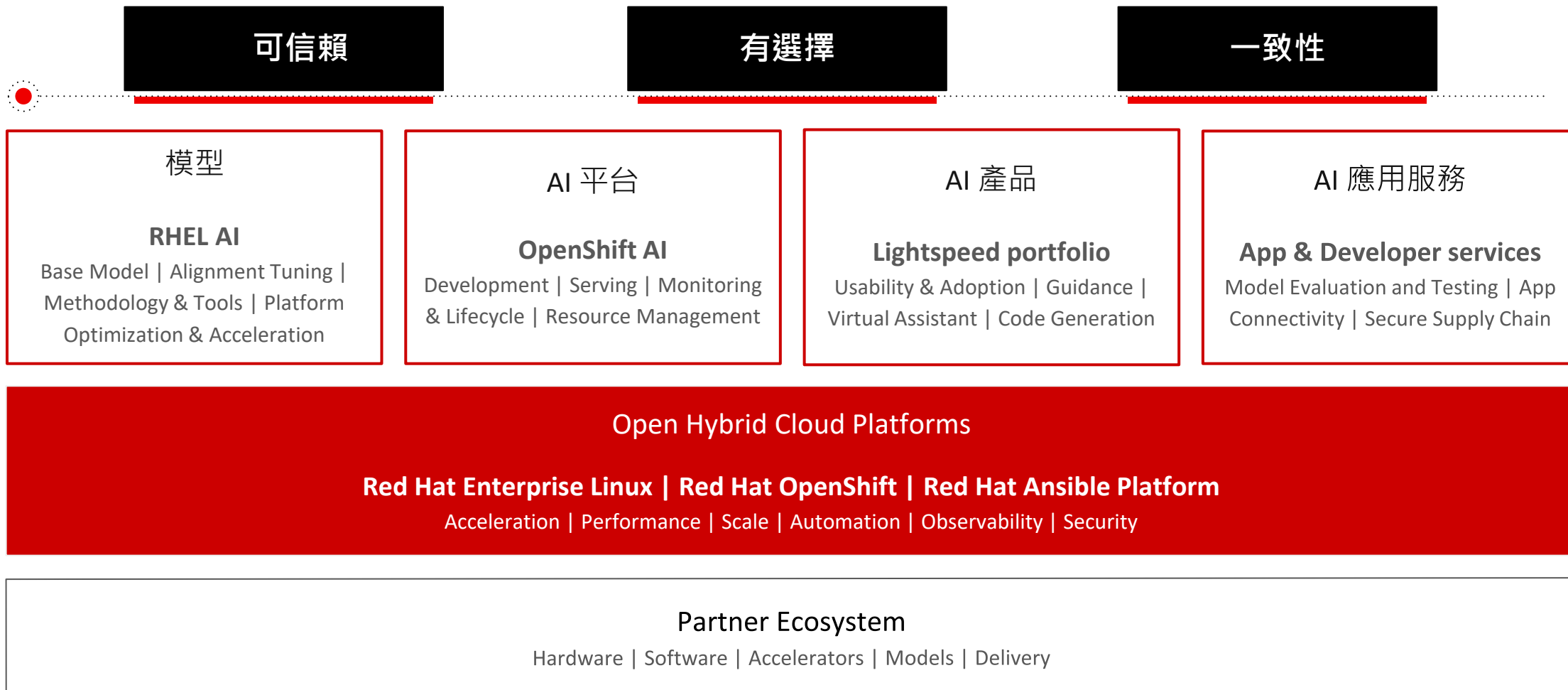


Source: Forrester Research, Inc. Unauthorized reproduction, citation, or distribution prohibited.

The Forrester Wave™ is copyrighted by Forrester Research, Inc. Forrester and Forrester Wave™ are trademarks of Forrester Research, Inc. The Forrester Wave™ is a graphical representation of Forrester's call on a market and is plotted using a detailed spreadsheet with exposed scores, weightings, and comments. Forrester does not endorse any vendor, product, or service depicted in the Forrester Wave™. Information is based on best available resources. Opinions reflect judgment at the time and are subject to change.

11 GARTNER is a registered trademark and service mark of Gartner and Magic Quadrant is a registered trademark of Gartner, Inc. and/or its affiliates in the U.S. and internationally and are used herein with permission. All rights reserved. This graphic was published by Gartner, Inc. as part of a larger research document and should be evaluated in the context of the entire document. The Gartner document is available upon request from Red Hat. Gartner does not endorse any vendor, product or service depicted in its research publications, and does not advise technology users to select only those vendors with the highest ratings or other designation. Gartner research publications consist of the opinions of Gartner research organization and should not be construed as statements of fact. Gartner disclaims all warranties, express or implied, with respect to this research, including any warranties of merchantability or fitness for a particular purpose.

Red Hat AI的策略與服務 - 協助企業導入A的旅程



Red Hat與AMD攜手合作推動企業AI時代



Red Hat

Products

Solutions

Training & services

Resources

Partners

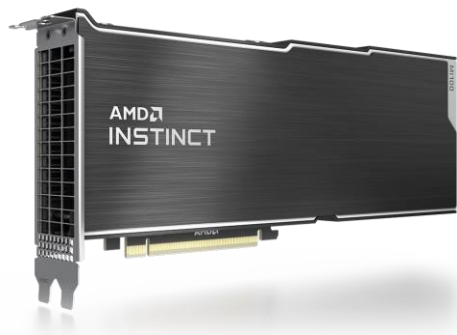
About

[Press releases](#) > Red Hat and AMD Collaborate to Advance AI Solutio...

Red Hat and AMD Collaborate to Advance AI Solutions and Empower Enterprises for the Cognitive Era

AMD GPUs on Red Hat OpenShift AI will enable greater choice and flexibility in AI architectures, lowering entry barriers for enterprises looking to embrace AI workloads

AMD Instinct MI 系列 GPU 加速器



MI100

32GB VRAM

Arcturus GPU
7680 cores
7nm CDNA1
(2020)



MI210

64GB VRAM

Aldebaran GPU
6656 cores
6nm CDNA2
(2022)



MI250

128GB VRAM

2 x Aldebaran GPUs
13312 cores
6nm CDNA2
(2022)



MI300X

192GB VRAM

MI300
14080 cores
5nm CDNA3
(2023)

RHEL AI平台支援AMD GPU

RHEL AI 1.2 現在支援 AMD accelerators



Products

Technologies

Learn

Events

Developer Sandbox

Blog

Videos

Announcing the General Availability of Red Hat Enterprise Linux AI (RHEL AI) Version 1.2!

October 15, 2024



[Yashwanth Maheshwaram](#)



Related topics: [AI/ML](#)

Related products: [Red Hat Enterprise Linux AI](#)

OpenShift平台支援AMD GPU

AMD GPU Operator Installation

The screenshot displays the OpenShift OperatorHub interface. The top navigation bar includes the Red Hat OpenShift logo, a hamburger menu, and user information (kube:admin). A blue banner indicates the user is logged in as a temporary administrative user. The left sidebar shows the navigation menu with 'OperatorHub' selected. The main content area is titled 'OperatorHub' and contains a search bar with 'amd' entered, resulting in one item: the AMD GPU Operator. The operator card shows the AMD logo, 'Community' tag, and a description: 'AMD GPU Operator provided by amd-gpu-operator. Operator responsible for deploying AMD GPU kernel drivers and device plugin.'

Red Hat OpenShift

You are logged in as a temporary administrative user. Update the [cluster OAuth configuration](#) to allow others to log in.

Project: openshift-kmm

OperatorHub

Discover Operators from the Kubernetes community and Red Hat partners, curated by Red Hat. You can purchase commercial software through [Red Hat Marketplace](#). You can install Operators on your clusters to provide optional add-ons and shared services to your developers. After installation, the Operator capabilities will appear in the [Developer Catalog](#) providing a self-service experience.

All Items

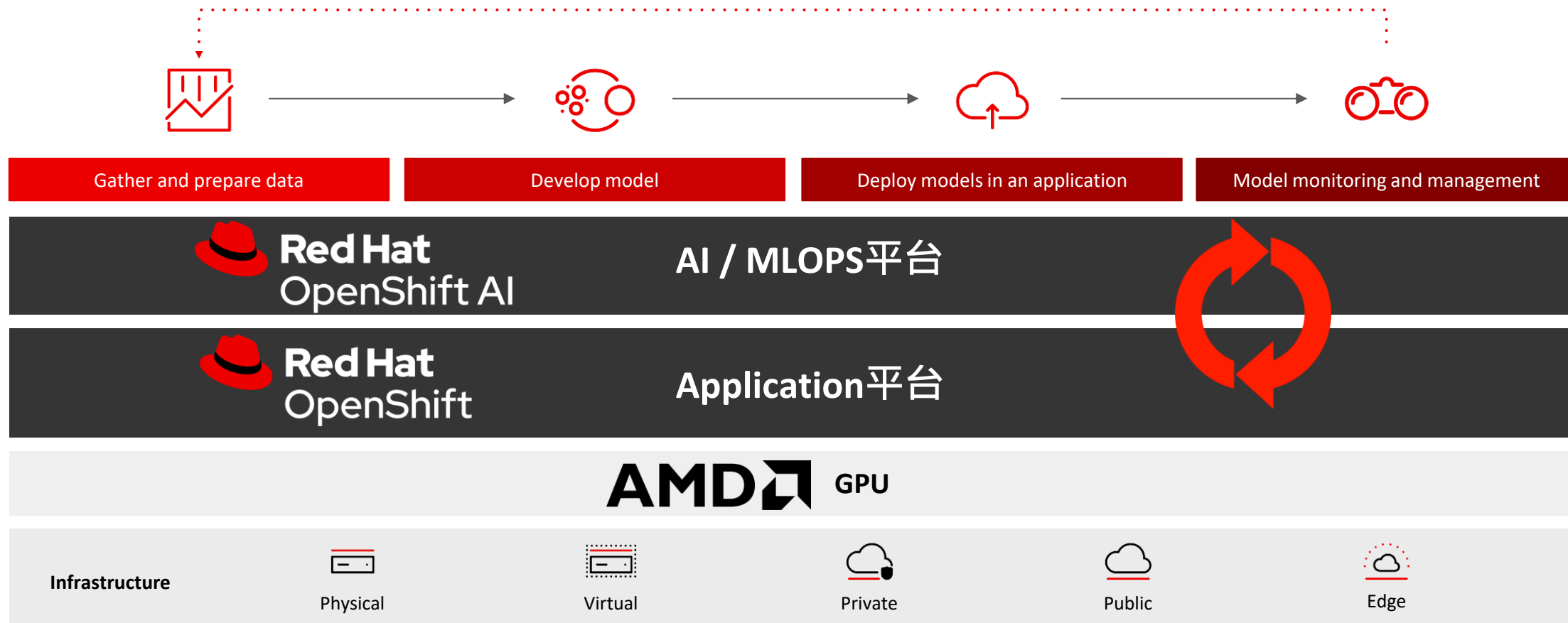
Search: amd 1 items

AMD Community

AMD GPU Operator
provided by amd-gpu-operator

Operator responsible for deploying AMD GPU kernel drivers and device plugin

Red Hat+AMD提供完整的應用程式平台+運算加速能力+AI MLOps平台




Thank You!

如果沒有全球合作夥伴生態系統的支持和積極協作，Red Hat就無法實現我們的AI願景。

AI並沒有“一套適合所有的標準（one size fits all）”，也沒有任何一家供應商能夠單獨滿足所有客戶的需求。

資料來源：[The AI opportunity is defined by a skilled ecosystem](#)

 [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