

Edge AI浪潮下的硬體業者發展觀察



申作昊 Zouhao Shen



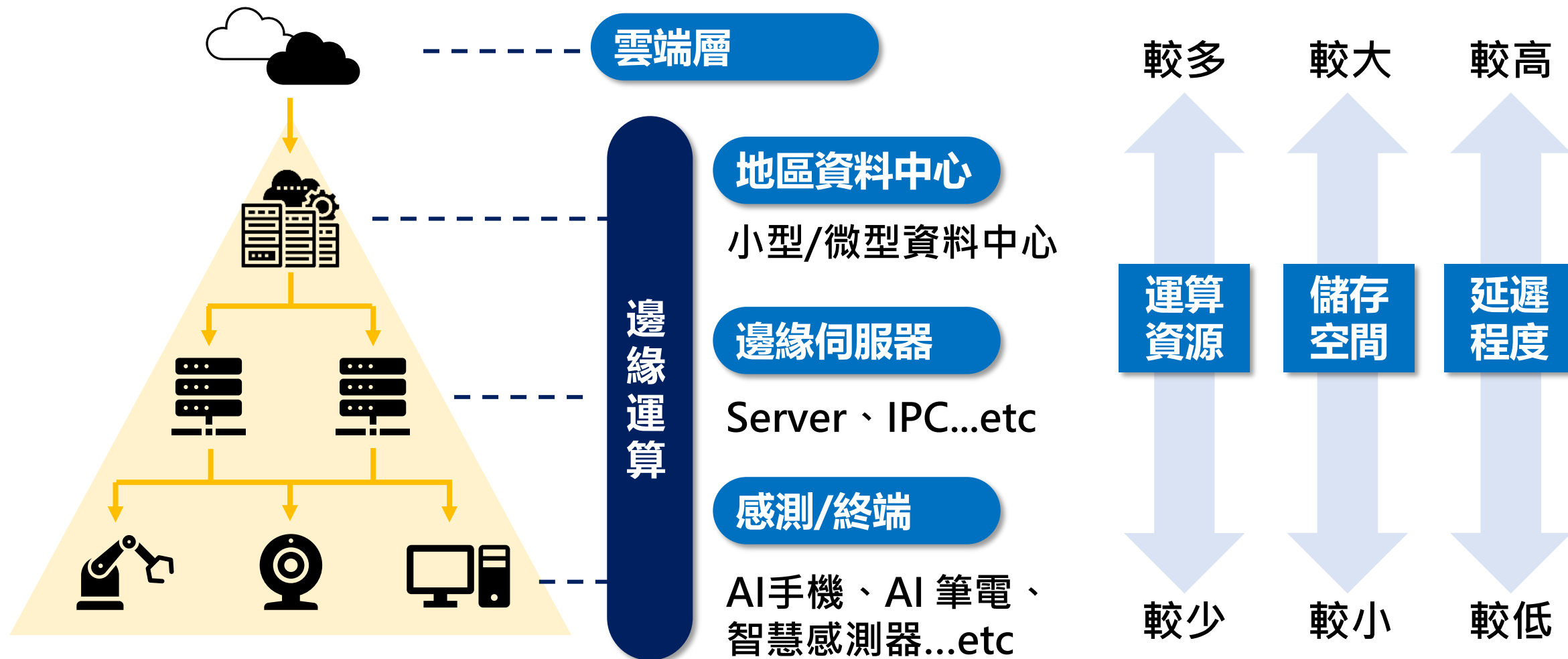
DIGITIMES研究中心



April.2024



「邊緣」運算定義：運算節點較靠近終端

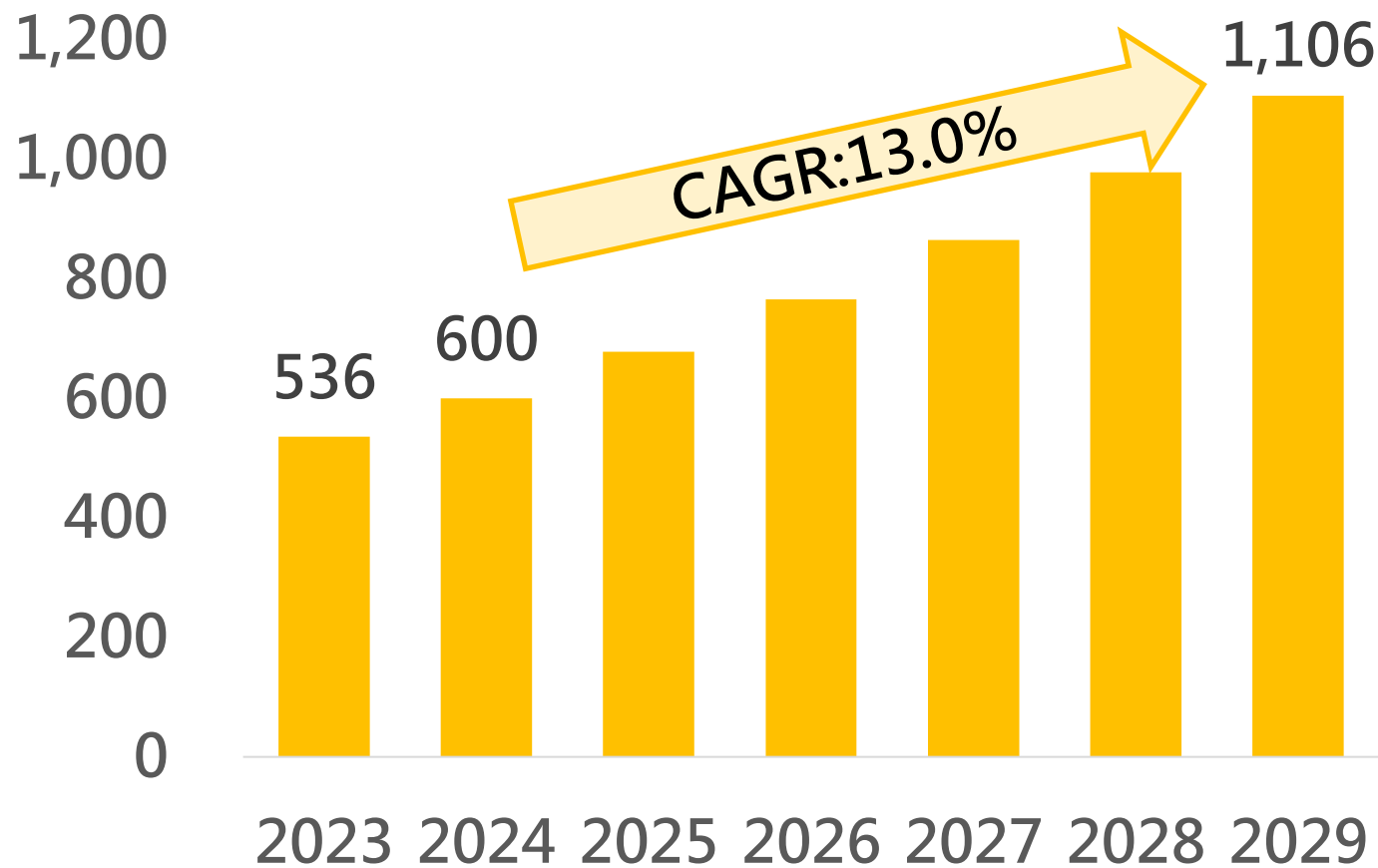




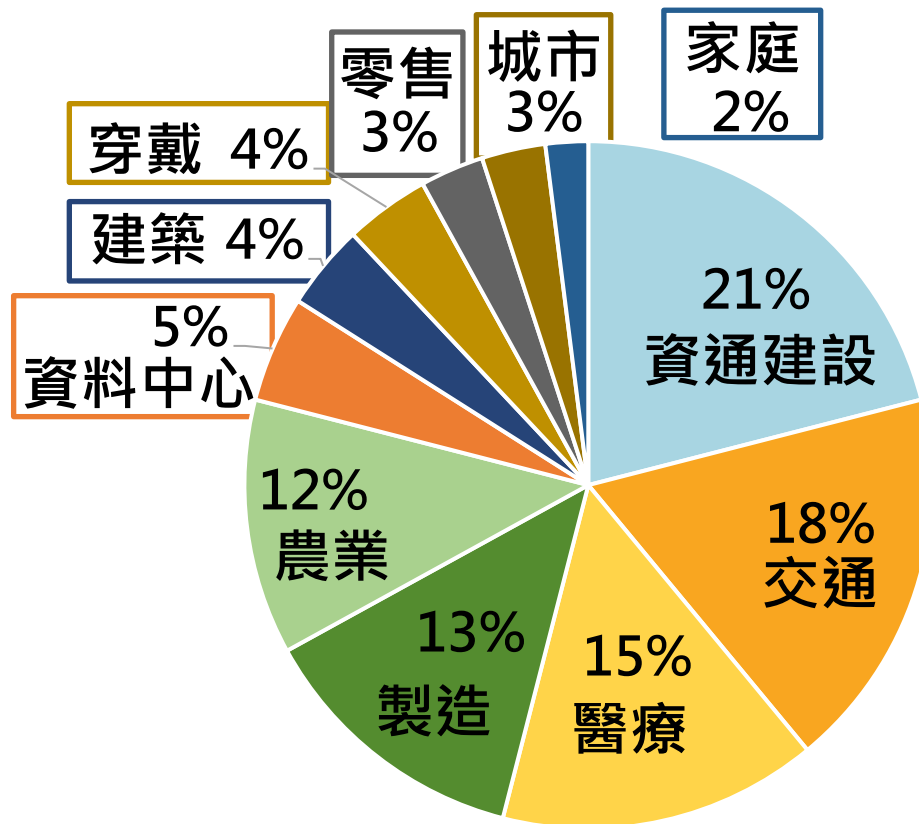
全球運算市場穩定成長 年增長率估約15.7%

單位：億美元

全球邊緣運算市場成長預測



邊緣運算各場域應用分布





邊緣市場構成與重要參與者列舉

Edge AI 軟體

整合式開發平台



預訓練模型



Edge AI晶片

通用晶片-GPGPU



MCU



NPU(XPU)



邊緣設備/邊緣終端

邊緣伺服器



邊緣IPC



消費型終端

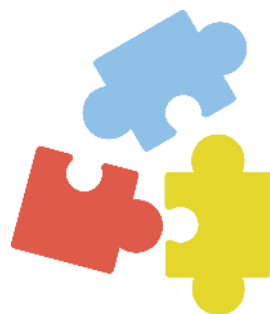


智慧終端





邊緣端產業生態兩大發展重點



市場需求

應用專業化



供應商策略

強調軟硬整合



Edge AI市場往專業分工領域發展

DIGITIMES



邊緣端市場趨向應用專業化



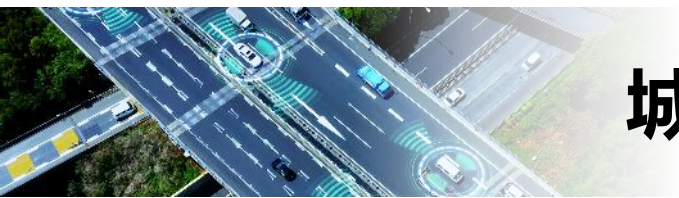
車載電腦



工業應用



安防監控

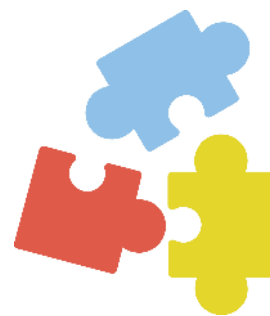


城市交通



智慧家居

...etc



市場需求

應用專業化

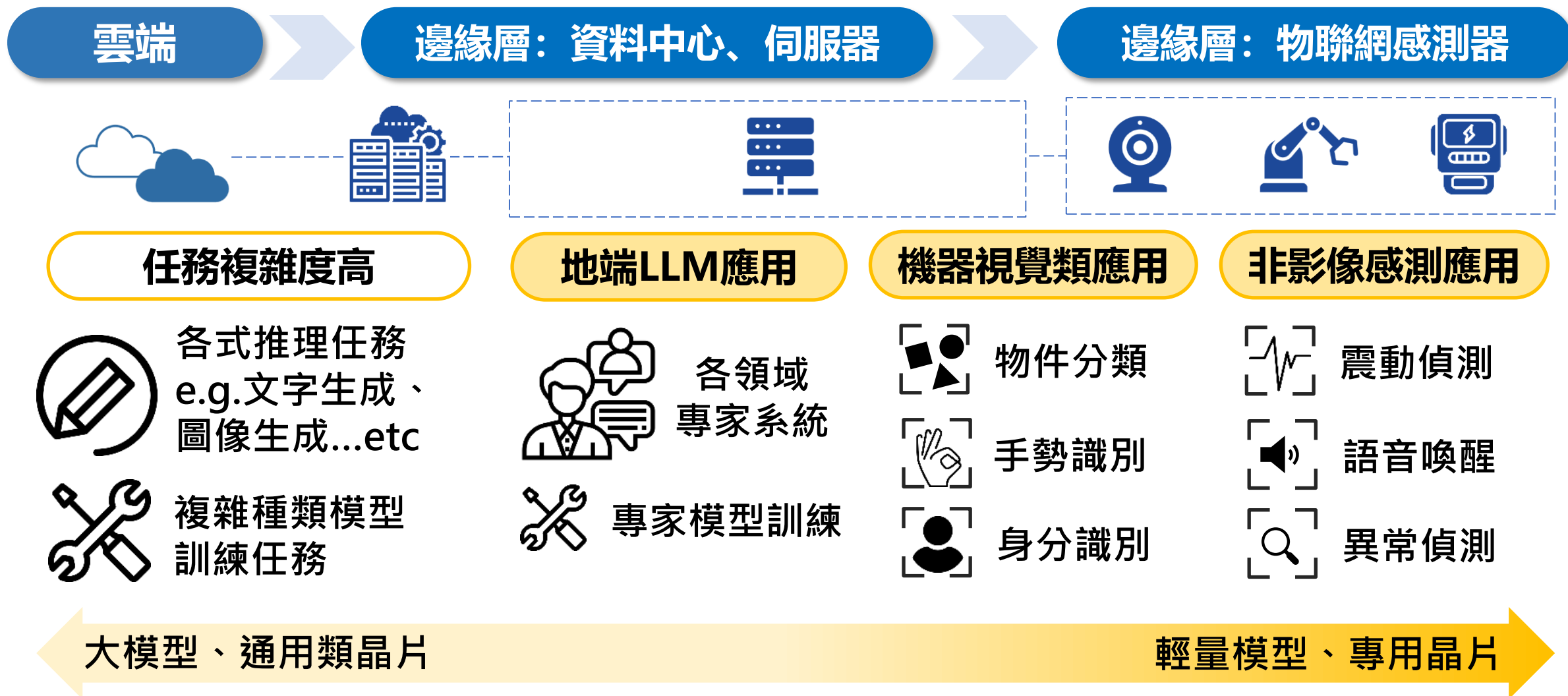


供應商策略

強調軟硬整合

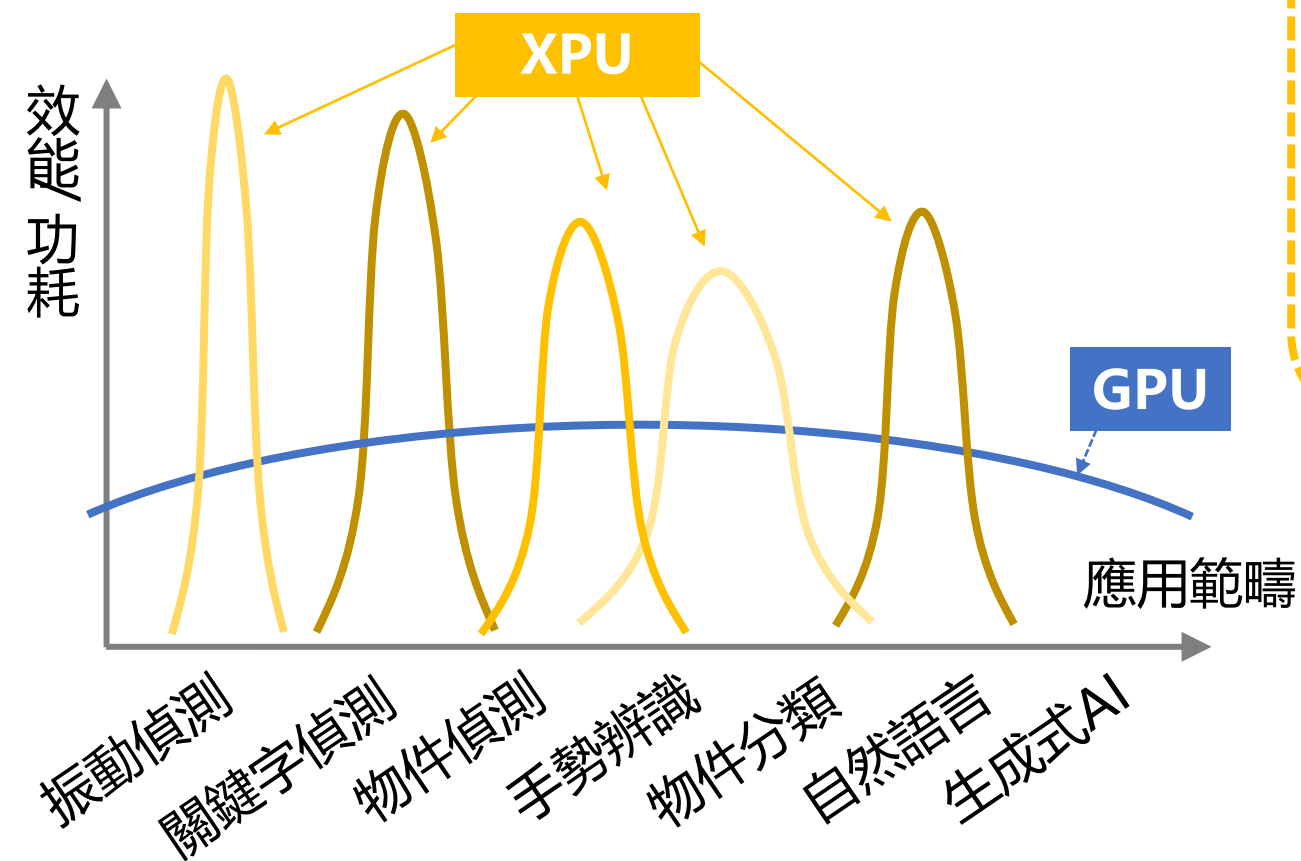


邊緣架構依情境處理不同類型AI任務





邊緣端AI晶片類別(GPU、XPU)與差異



圖片來源: Hailo

NPU(XPU)

專用度較高



- ✓ 多針對特定AI任務開發
- ✓ 特定情境下效能、功耗等較佳
- ✓ 須一併考量AI演算法開發
- ✓ 較適合邊緣端單一應用情境

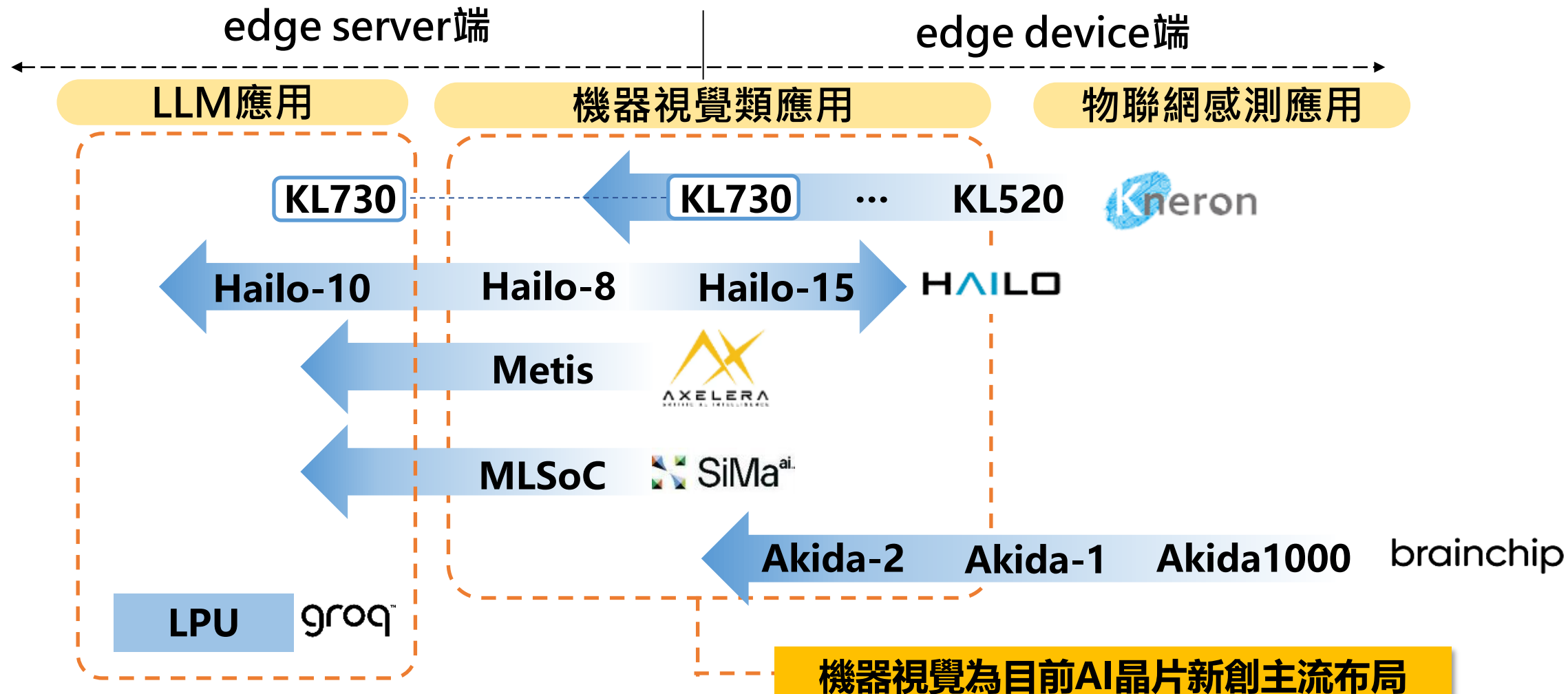
GPU

通用度較高



- ✓ 可處理多種、複雜AI任務
- ✓ 較適合用於AI模型訓練
- ✓ 軟體開發生態較開放、完善
- ✓ 較適合雲端多樣任務情境

機器視覺為當前AI晶片業者主流應用方向



部分業者近年開始發展LLM解決方案

機器視覺為目前AI晶片新創主流布局

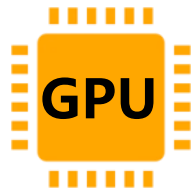


邊緣運算各環節AI加速硬體需求不同 解決方案各異



雲端/資料中心

- AI任務種類多樣
- 可擴充性要求高
- 訓練模型需求高



通用性高
處理各種任務

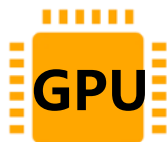


雲端業者自研
配合自家算法

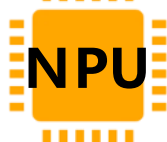


Edge server

- 可擴充性要求高
- 需視應用調整規格



通用解決方案



專用解決方案
保留可擴充性



整合運算單元
具性能優勢



Edge device

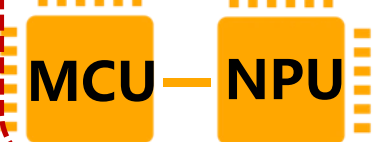
- 應用單一、專用度高
- 成本、能耗敏感



外接運算能力
保留設計彈性



整合運算單元
具性能優勢



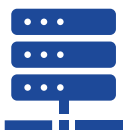
微控制器
內嵌運算能力



NPU業者多布局機器視覺技術 應用領域各不相同

機器視覺晶片各應用領域產品

業者逐漸從外掛NPU轉向發展整合度較高產品



資料中心/伺服器

- 運算力需求高
- 大量記憶體空間
- 可擴充性高



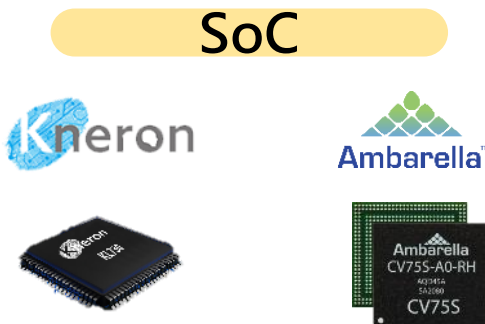
車載系統

- 運算力需求高
- 高可靠性
- 延遲程度低



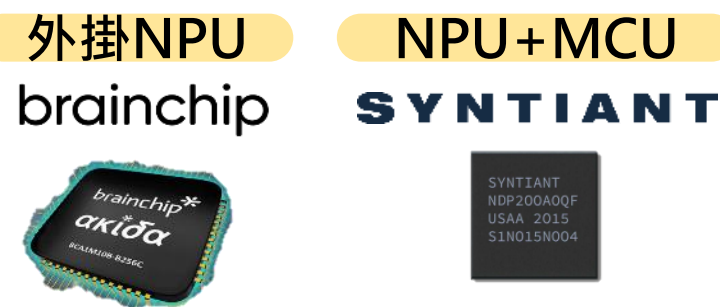
智慧鏡頭

- 延遲程度低
- 低耗能要求
- 成本敏感度高



物聯網終端

- 低耗能要求
- 成本敏感度高
- 保留設計彈性





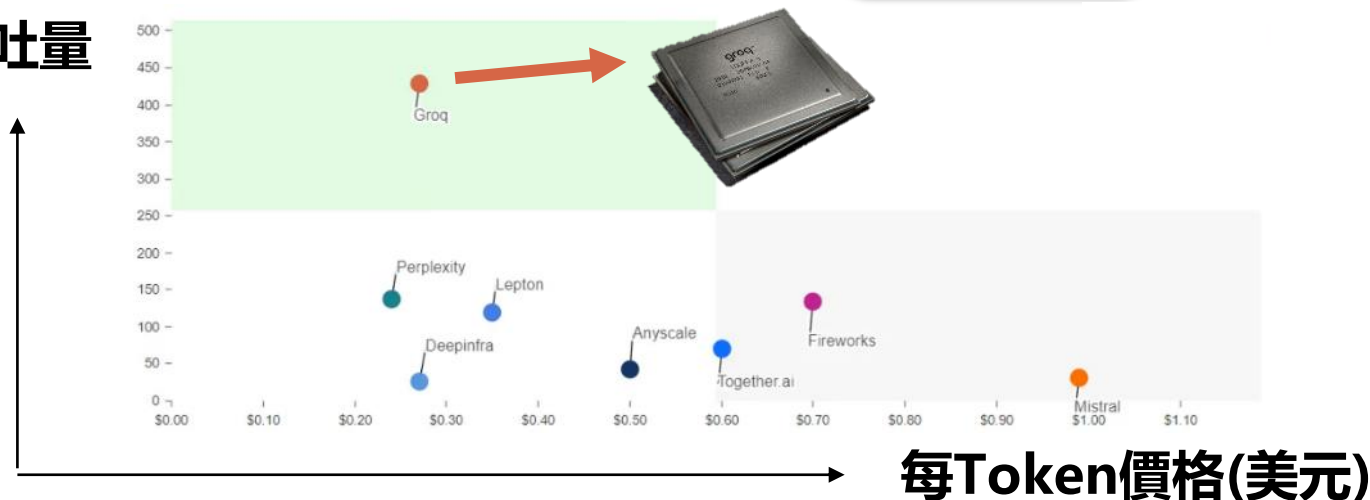
部分NPU新創業者近期推出LLM專用AI加速晶片

LLM專用晶片

groq™

LPU

吞吐量



推廣Groq cloud服務

- 購併 AI 解決方案公司 Definitive Intelligence
- 以高吞吐量、低收費為優勢

推廣Groq 晶片至資料中心

- 成立Groq Systems推銷晶片
- Edge LLM尚無特殊用例
- 建置成本為一大劣勢

HAILO

Hailo-10



- 以M.2加速卡形式推出，支援邊緣server或PC
- Hailo-10H提供40TOPS運算力
- 功耗低於 3.5W



- 表達近年將推出LLM專用開發軟體或硬體



處理器、MCU大廠新產品趨向內嵌NPU增加運算力



處理器、MCU大廠展出新品內嵌NPU



Versal AI Edge Gen 2

- 整合NPU的系統單晶片，共6款產品
- 標榜低功耗與低延遲適於嵌入式應用
- 以車電、工業、醫療等為目標市場

AMD邊緣AI軟硬體解決方案

軟體



Vitis AI開發環境
(具模型壓縮優化工具)

硬體



Versal AI Edge SoC
運算力20~200 TOPS



PSoC Edge MCU系列

- 以物聯網應用為目標市場
- E83與E84 MCU內建Ethos-U55 NPU
- 搭配ModusToolbox軟體開發工具
- 2025年量產，預估單價6~8美元



FRDM開發板內嵌N947 MCU

- 具兩顆Cortex-M33核心
- 內建自研eIQ Neutron NPU
- 搭配MCUXpresso開發工具

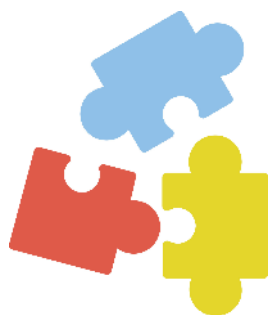


Edge AI硬體業者走向軟硬整合

DIGITIMES



邊緣端產業生態兩大發展重點



市場需求

應用專業化



供應商策略

強調軟硬整合

客製化程度高
開發導入時間拉長



雲端通用AI模型
不適用邊緣硬體



客戶獨立開發軟體解決方案面臨挑戰



產業應用軟體合作



軟體開發工具

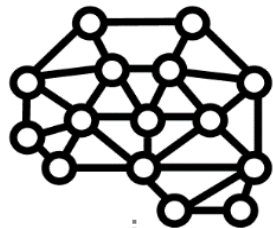


邊緣整合式開發業者將雲端開發工具串聯至邊緣端

Edge AI整合式開發平台業者



壓縮



編譯



部署



串聯雲端、邊緣開發生態



預訓練模型
與開發工具



模型壓縮
優化工具



串聯MCU業者開發工具



生成數據
模擬驗證









導入邊緣
硬體模組



AI晶片業者提供軟體開發套件支援客戶開發相容模型

AI晶片業者 提供軟體開發套件

-  ▶ MCUXpresso
-  ▶ ModusToolbox
-  ▶ ST Edge AI Suite
-  ▶ Palette
-  ▶ AI Software Suite
-  ▶ VOYAGER SDK

主流模型框架



TensorFlow



PyTorch



Keras



ONNX



mxnet



Caffe2

...etc

Edge模型訓練

模型框架

使用者
模型

模型庫
示範

編譯器

處理器應用平台

使用者
應用

應用
示範庫

模型部署/更新工具

設備監控工具

合作Edge平台整合業者



EDGE
IMPULSE



Deeplite



deci



LatentAI



EDGE
NEURAL.ai



alwaysAI



Plumerai

...etc

合作產業應用軟體業者



Micro.ai



RealityAI



fogsphere



CVEDIA

...etc

AI晶片業者提供預訓練模型加速客戶開發

預訓練邊緣模型

外部合作取得資源

deci.
Break the AI Barrier

物件偵測、物件分類、
語意分割

Deeplite

物件偵測、物件分類



LatentAI

物件偵測、物件分類



imagimob

聲音偵測

EDGE

NEURAL.ai

物件偵測

內部研發

NXP

RENESAS

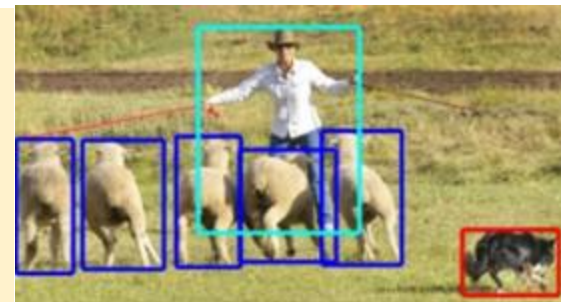
ST
life.augmented

HAILO

AXELERA
ARTIFICIAL INTELLIGENCE

晶片業者本身
開發預訓練模型

物件偵測



物件分類



語意分割





AI晶片業者與場域軟體業者合作提供應用解決方案

製造場域應用



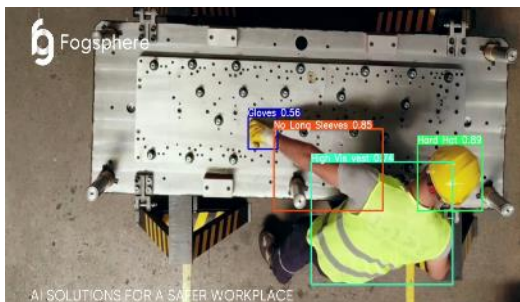
訂單管理系統



- 訂單管理系統
- 機器視覺：機器手臂抓取點、拾取物體識別
- 連線倉儲管理系統，達成自動化訂單管理



工安偵測系統

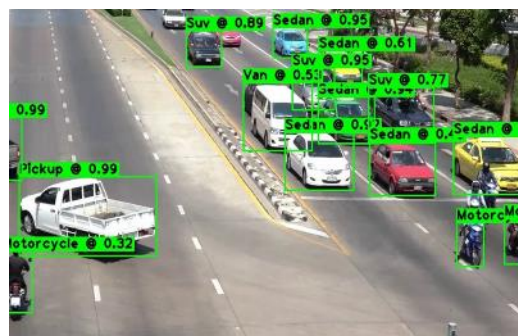


- 機器視覺：安全裝備配備狀況
- 行為安全、緊急通報、個人防護設備檢測

交通場域應用



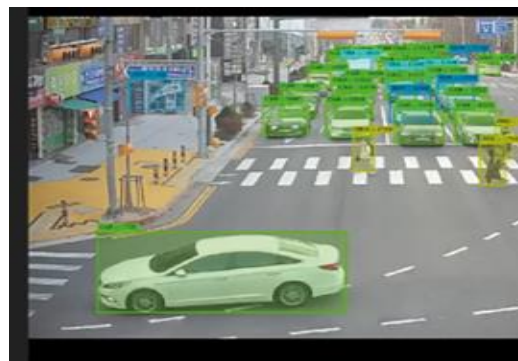
智慧交通系統



- 智慧交通系統
- 機器視覺：車輛識別、行人識別
- 車流計算、分類、人流分析等應用



路況偵測系統

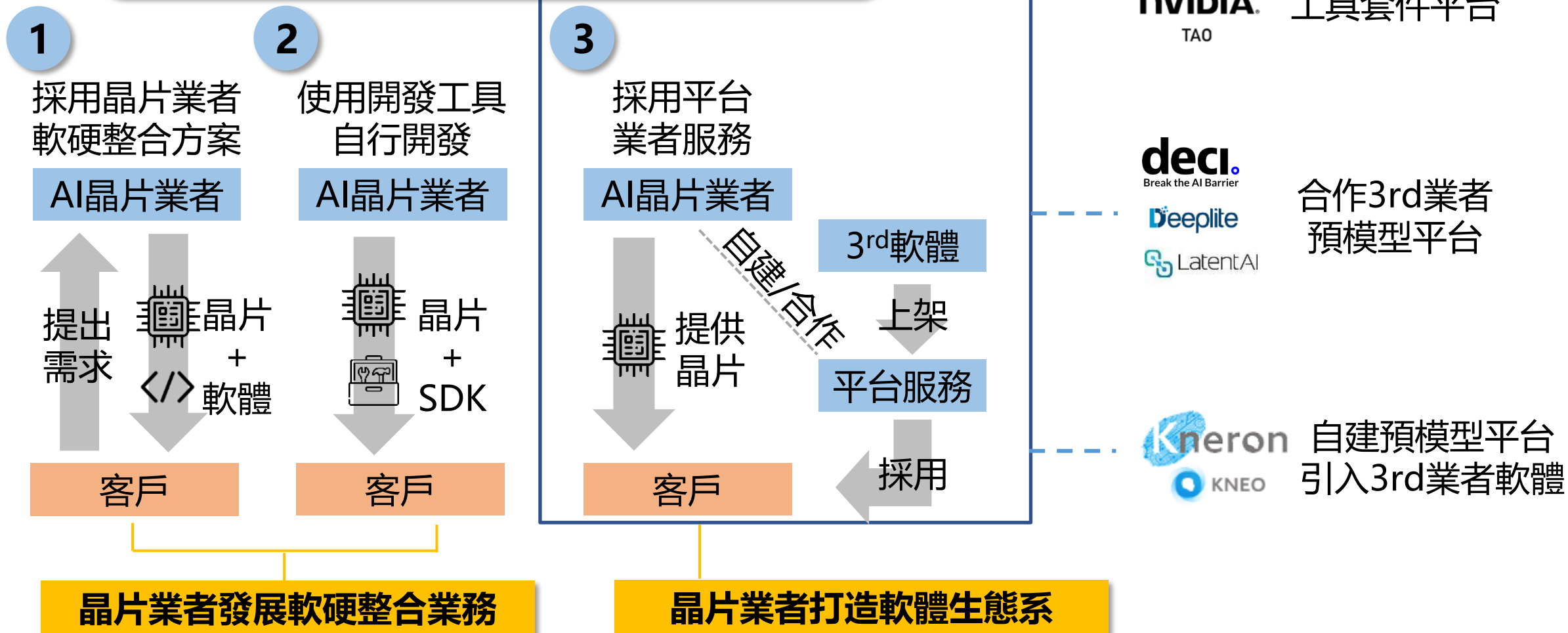


- 機器視覺：車輛識別、行人識別
- 路口監控、事故偵測、智慧停車場等應用



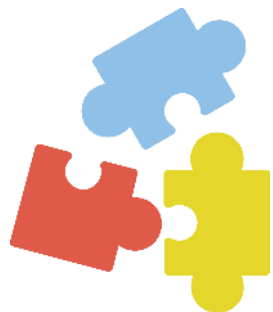
AI晶片硬體業者逐步打造自家軟體生態體系

邊緣AI晶片導入模式





結論



市場需求

應用專業化

穩織合度

適得其所

- 硬體求取成本與效能平衡
- 通用性與專用性的取捨



供應商策略

軟硬整合方案

串聯生態系

降低採用門檻

- 不再固守硬體業務
- 發展微垂直領域解決方案

DIGITIMES

THANK YOU



申作昊 Zouhao Shen



DIGITIMES研究中心



April.2024