



# Sangfor AICP

## 輕鬆釋放GPU算力, 軟件定義AI政務創新

Steven Lo | 方案顧問

[Steven.lo@sangfor.com](mailto:Steven.lo@sangfor.com) | 深信服科技

# DeepSeek趨勢下AI進入普惠階段，但依然成本挑戰巨大



## 專業人才成本高

AI能力  
要求高

供需  
失衡

行業競  
爭激烈

培養  
成本高



## 基礎設施成本高

硬體資  
源成本

供電成  
本高

平臺工  
具成本

開發運  
維成本



## 運維成本提高

資源分  
享  
管理難

版本管  
理難

視覺化  
監控難

故障定  
位難

# 從擁抱AI的一開始到大範圍使用都要考慮ROI

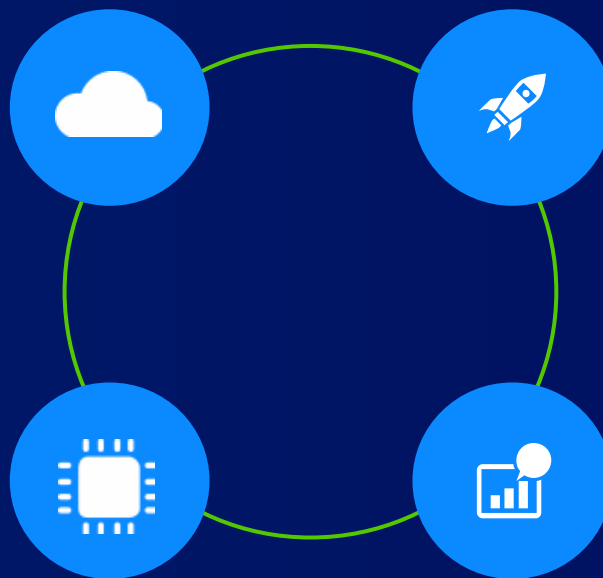
基於場景選擇模型，基於模型選擇顯卡，同時關注平臺性能優化用確定的顯卡發揮更大的效能

## 平臺性能優化

算力平臺不僅僅是一個資源管理器，更應該通過平臺自身的性能優化提升客戶推理性能

## 基於訓練/推理以及信創要求選擇顯卡

結合自身的建設需求以及需要的模型大小，在4090D、H20、沐曦、昇騰、天數、等顯卡中進行選擇



## AI優化新技術

比如異構調度優化、昇騰優化、語義Cache等新技術的發展，大幅降低客戶的推理成本

## 根據場景選擇模型大小

以Deepseek為例，要基於使用場景選擇671B、70B、32B、14B、7B、1.5B

# 深信服長期佈局AI產業，團隊和技術能力積累扎實



明確AI First戰略  
各產品融入AI能力

2021年

正式組建AIC  
產品開發團隊

2023年

中國首發  
安全GPT大模型

2023年5月

行業領先發佈  
AICP產品

2024年

AI應用創新平臺  
發佈

2025年2月

AICP2.1創新支持  
4090D商用671B方案

2025年4月

業界領先、獨具特色的  
AI Infra

未來持續構築

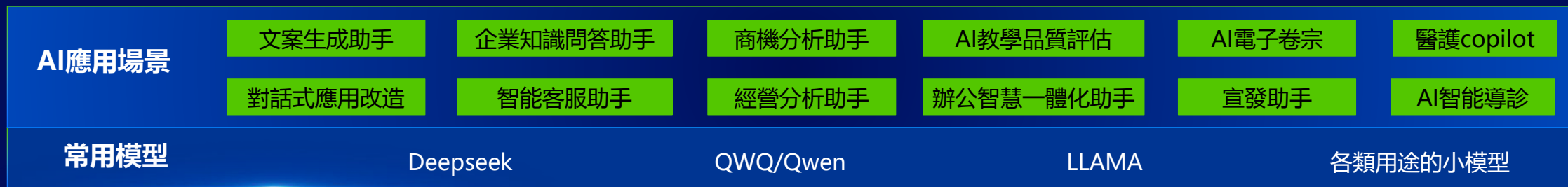


# 深信服AI升級的解決方案：AICP 及 SF-FastGPT



深信服AI平台用戶

# 深信服新一代AICP平臺架構



向上承載各類大小模型，服務好客戶的業務應用



向下開放相容，廣泛適配並調優、調度各硬體廠商硬體基礎設施



# 智慧融合架構，降低複雜性，服務化的提供基礎設施能力



# 算力基礎設施服務化，多模型、多卡相容不綁定硬體



Qwen-VL Embedding Rerank OCR .....

緊跟趨勢，不斷相容適配更多卡商、各類開源模型、多模態模型與微調模型

NVIDIA			META 沐曦	天数智芯	HYGON	Ascend
Nvidia A100	Nvidia 4090D	Nvidia L40S	N260	智鎧MR V100	K100-AI	910B3
Nvidia A800	Nvidia 5090D	Nvidia A6000	C500	天垓BI 150	.....	910B4
Nvidia 4090	Nvidia L40	Nvidia H100	C550	.....		300I DUO PRO
Nvidia H20	Nvidia L20	.....	.....			.....

不斷適配調優更多型號

## 深信服 AICP 算力平台

英偉達顯卡  
快速相容

新品牌型號顯卡  
1個月完成適配

擁抱開源  
最新推理框架即時支持

異構算力及模型資源  
統一管理

異構算力資源  
統一調度與性能優化

OS、AICP軟體、容器平臺  
層層解耦

# vGPU大幅提升顯卡資源利用率，承載更多AI Agent應用

RAG、智能客服等AI應用建設  
都需要大量小模型

vGPU

支持算力1%，  
256MB級顯存資源切分

運行Embedding+Rerank+OCR  
小模型場景

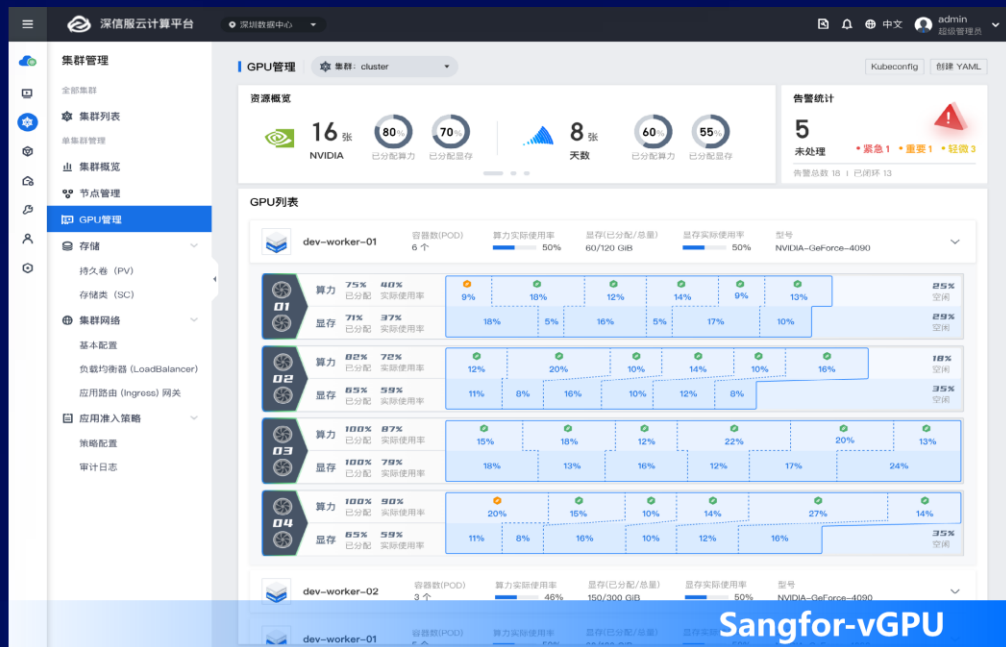
TCO分析

顯卡  
消耗

顯卡越  
高端

模型使  
用越多

越省錢



4090D

未切分

3張卡

vGPU切分

1~2張卡

H20

3張卡

1~2張卡

L20

3張卡

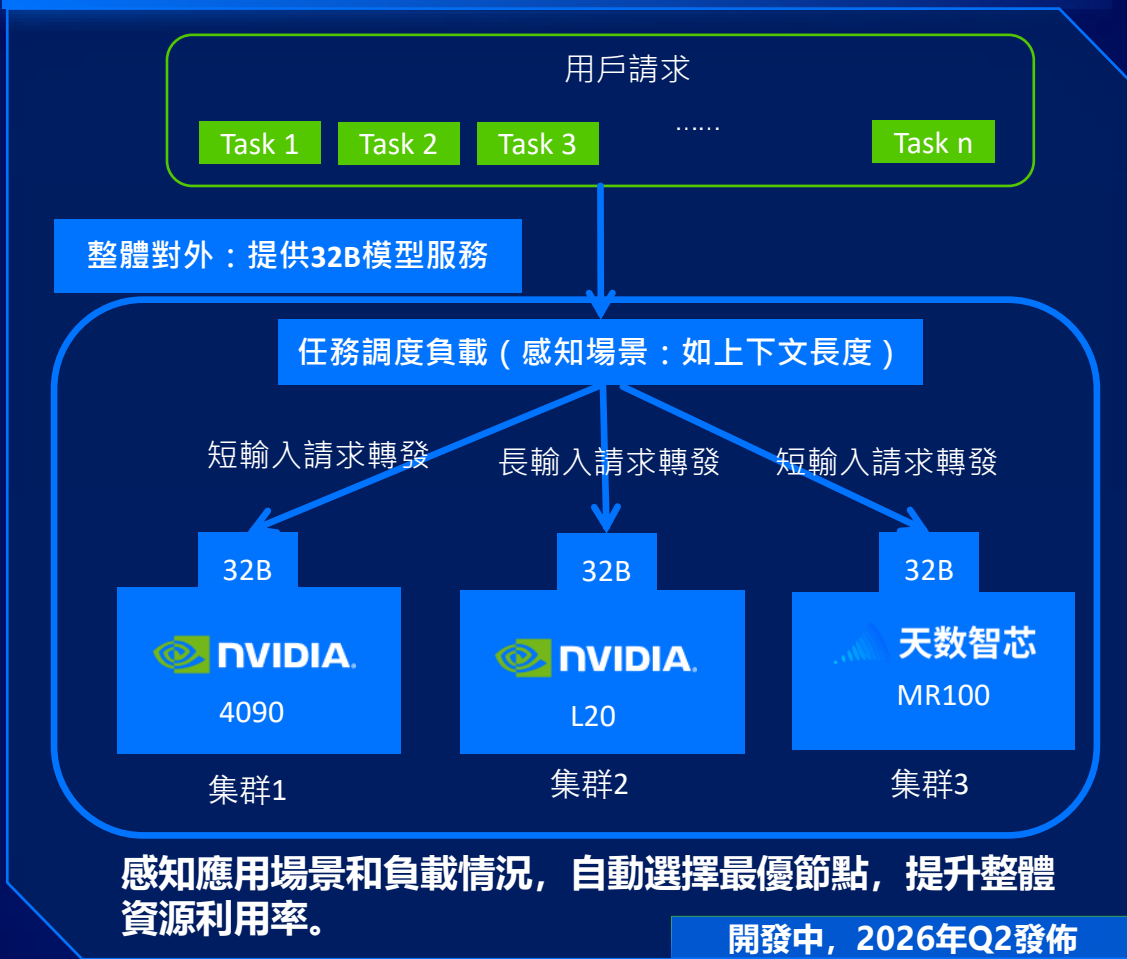
1~2張卡

單卡承  
載模型  
數量翻  
8倍+

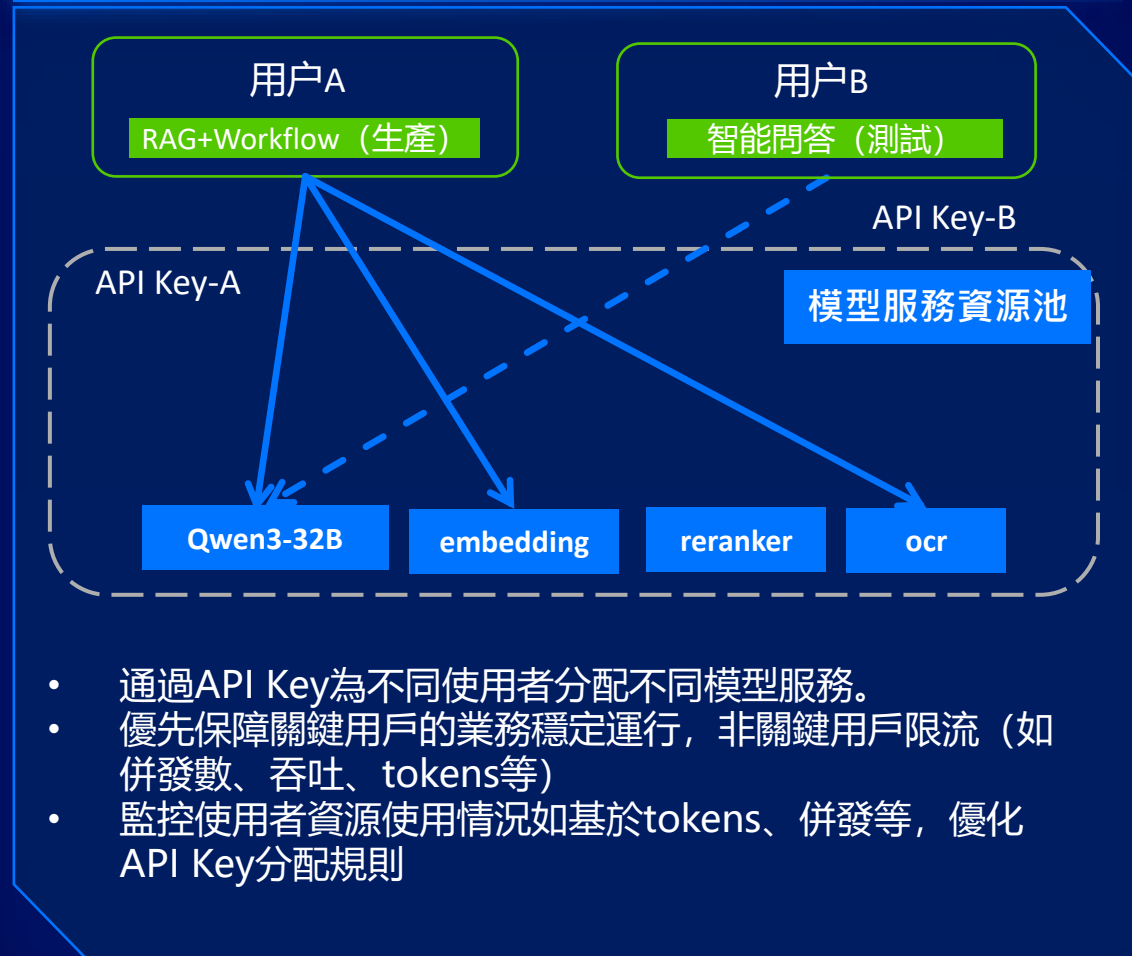
現有資  
源利用  
率同比  
3倍+

Sangfor-vGPU

## 01 自動化的異構推理集群調度



## 02 多API KEY高效利用顯卡資源，保障關鍵業務體驗



通過底層核心技術創新，發揮2~5倍智能算力設施效能



## 重載應用場景

# 創新優化啟動智慧算力效能2~5倍起步

單卡平均承載用戶數 **2x**倍提升

ROI : **2x**倍提升



**H20**

p99 TTFT **6.4x**倍提升

Decode速度 **2.0x** 倍提升

p99 E2E Latency **2.0x**倍提升



**4090D-48G**

總承載用戶數 **2.25x**倍提升

ROI : **2.0x**倍提升



**910B3**

模型：GLM4.6 (H20)、GLM4.6-Air (4090D-48G、910B3)；

SLO標準：p99 TTFT < 10s, p99端到端時延 <= 60秒；

推理框架：相比sglang最新社區版本

# 通過關鍵技術創新，不斷提升端到端ROI

AI應用

AI  
智能體

AI  
智能客服



## 自我調整架構層Smart Arc

自我調整  
模組

自我調整多版本推理引擎

自我調整多類型推理引擎

性能瓶頸  
分析

自我調整調整

回流數據  
觀測

一鍵最佳實踐配置

指標採集

開源推理框架/自研推理框架：

SIF (Sangfor Inference Framework) --加密、可靠優化、性能優化

自我調整  
原子能力

AI 開道  
智能負載調度、限流、計量

基於檢索的  
投機解碼優化

算力資源自我調整  
MOE運算元優化

基於場景化的分塊調度優化

GPU間、節點間資料傳輸優化

基於場景語義感知的Cache智慧負載

長輸入性能優化

昇騰優化

.....

671B低門檻起步：當前業界唯一商業化，4090D起步運行含671B的端到端AI應用場景



【方案1】全量8卡：8併發80tokens

【方案2】4+4：DeepSeek：4卡8併發80tokens；QWQ 32B：4卡100+併發1000+tokens

【方案3】4+2+X：DeepSeek：4卡8併發80tokens；QWQ 32B：2卡；X小模型2卡

# 其他價值優勢

廣泛ISV生態  
適配調優

中科聞歌

衡石

行云创新  
CLOUD·TOTO

鼎捷软件  
DigiwinSoft

私有雲  
AI升級

HCI加1台GPU節點即可輕鬆獲取AI能力

AI能力與線上體驗一致

構建  
立體AI  
安全能力

模型訓練一  
鍵微調

嚮導化配置一鍵微調

模型  
加密

模型倉庫

支持DeepSeek等開源大模型統一運  
維管理上傳管理

支援從外部上傳自訂大小模型運  
維管理

通用場景、長輸入場景最佳實踐，  
一鍵發佈模型服務

數據  
內容安全

# 模型倉庫讓各類大小模型運維管理更簡單省心

開源  
模型

支持DeepSeek等開源大模型統一運維管理上傳管理

自訂大小  
模型

支持從外部上傳最新、自訂大小模型運維管理

一鍵發佈  
模型服務

無需配置算力，一鍵發佈均衡、長文本、自訂模型服務

The screenshot shows the 'Model Warehouse' (模型倉庫) interface. On the left, a sidebar lists navigation options: 總覽, 數據集管理, 模型倉庫, 模型訓練, 模型優化, 模型服務, 在線體驗, and 平台管理. The main area displays a list of models:

- DeepSeek-R1-671B-sangfor (公共)
- Qwen2.5-72B-Instruct-GPTQ-Int4 (公共)
- Qwen2.5-14B-Instruct (公共)
- Qwen2.5-72B-Instruct (公共)
- Qwen2.5-32B-Instruct (公共)
- Qwen-2.5-vl-7b-Instruct (公共) - qwen2.5 multi modal vision language model.
- DeepSeek-R1-Q4 K M (公共)

Below the list, a table shows details for two models:

Qwen2.5VL-7B	1	版本數
DeepseekQ4-671B	1	版本數

The 'Create Model' (創建模型) form is also visible, with a red box highlighting the 'Import Method' (導入方式) section, where 'Import from external platform' (從平台外部導入) is selected.

The screenshot shows the 'Create Service' (創建服務) configuration page. A red box highlights the 'Basic Configuration' (基本配置) section:

- Service Name: DeepSeek-R1-Distill-Qwen-32B
- Model: DeepSeek-R1-Distill-Qwen-32B
- Model Version: V1

The 'Resource Configuration' (資源配置) section shows:

- Compute Power: NVIDIA A100 推理池
- Deployment Mode: 均衡模式 (Selected), 長文本優化, 自定義模式
- Single Instance Compute Power: 算力: A100-80GB(卡數:1, 顯存:80GB), CPU 內存: 8 核 16 GB
- Instance Quantity: 2

The 'Service Configuration' (服務配置) section shows:

- Network: tgw-0802
- Service Port: HTTP, 10.10.10.10, 請輸入端口, 即: 4430 (已用端口)
- Encryption: 已開啟

# 立體安全能力，構築AI時代新安全體系



## 大模型自身安全

業界首推模型動態加密，確保模型與私有訓練內容不被竊取使用



## 環境安全

數據當地語系化，線上專屬雲AI環境，保障內網隱私資料安全



## 資料與參數內容安全

- 資料安全平臺：大模型應用安全護欄，實現提示注入攻擊防護以及有害或敏感資訊輸出的安全防護



開發中，2026年Q2發佈

# 模型訓練嚮導化配置，一鍵微調提升AI生產效率

## 01 选择基础模型

**Yi-34B**  
Yi系列是由01.AI的开发者从头开始训练的大型语言模型  
[查看建议](#)

**Qwen-72B**  
通义千问超大规模语言模型，支持中英文等不同语言输入。  
[查看建议](#)

**LLaMA2-7B V3.0**  
由Meta AI研发并开源的7B参数大语言模型，在编码、推理及知识应用...  
[查看建议](#)

**Baichuan2-7B**  
由百川智能推出的新一代开源大语言模型，采用2.6万亿Toke，由百川...  
[查看建议](#)

更多参数配置

开始训练 取消

## 02 选择数据集

安全知识量化对齐数据集	准备的安全量化对齐数据集	<input type="checkbox"/>
安全知识问答数据集	准备的安全相关微调数据集	<input type="checkbox"/>
医学知识量化对齐数据集	v3 (最新)	准备的医学知识量化对齐数据集 <input checked="" type="checkbox"/>
计算机相关数据集	准备的计算机相关垂直问题数据集	<input type="checkbox"/>
计算机知识量化对齐数据集	准备的计算机知识量化对齐数据集	<input type="checkbox"/>

数据配比: ● 自动 ○ 指定 1:

ChineseMedicalDialogue Data	医疗语料	中文医疗对话数据集由792099个问答对组成，包括男科、内科、妇产科、肿瘤...	<input checked="" type="checkbox"/>
FinCUGE_FinNL	金融语料	金融新闻分类数据集。对于给出的金融新闻，需要模型将其多标签分类到可能的...	<input type="checkbox"/>
BELLE_SchoolMath_data	中文语料	包含约25万条由BELLE项目生成的中文数学题数据，包含解题过程。	<input type="checkbox"/>
法律相关知识问答数据集	法律语料	法律知识问答数据集	<input type="checkbox"/>
QIZhen_data	医学语料	针对药品知识问答发布了评测数据集，后续计划优化疾病、手术、检验等方面的...	<input checked="" type="checkbox"/>

## 03 选择资源池

**NVIDIA训练资源池**  
训练

显卡使用	8张 / 16张
GPU使用	124 GB / 192 GB
存储使用	591 GB / 7 T

**NVIDIA推理资源池**  
推理

显卡使用	4张 / 4张
GPU使用	124 GB / 192 GB
存储使用	591 GB / 1.7 T

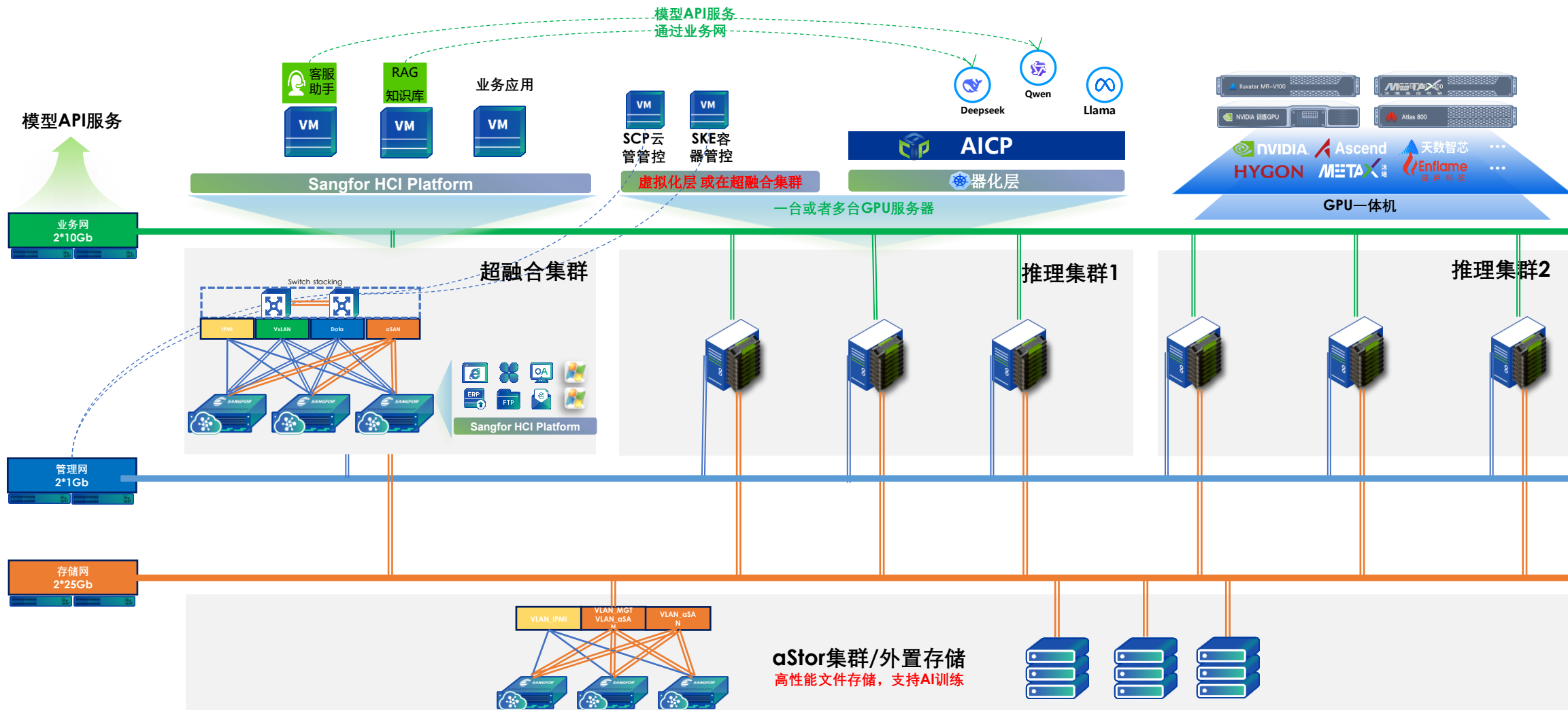
**华为昇腾推理资源池**  
推理

显卡使用	2张 / 8张
GPU使用	124 GB / 192 GB
存储使用	390 GB / 1 T

**华为昇腾推理资源池**  
推理

显卡使用	4张 / 8张
GPU使用	124 GB / 192 GB
存储使用	390 GB / 1 T

# AICP + HCI 架構圖



## AICP 能力要點

- 異構算力抽象與池化：多模型、多卡相容且不綁定硬體，持續適配 NVIDIA 與多家國產加速卡並統一調度與優化。
- 服務化接口與治理：每模型提供 API Gateway 與 API Key，支援按 Key 的限流與配額，保障多部門公平與穩定。
- 私有雲部署與全鏈路支持：主流/自有模型託管、資料版本化、單介面微調訓練、推理解耦優化（量化、蒸餾、KV cache）、模型加密。

## 成本與效能

- 藉由調度、混部與推理優化，在多款 GPU 上可實現約 2-5 倍效能提升，支撐高併發與長文本等重載場景，同時提升資源利用率與能效。
- 小模型（Embedding、Rerank、OCR 等）可使用 vGPU 資源分割能力，避免整卡佔用造成浪費。

## 完善生態協作

- 與資深 ISV 夥伴深度整合，提供從諮詢到運維的全流程支援

## 彈性起步，平滑擴展

- 最低僅需 1 個 HCI 節點 + 1 個 GPU 節點即可啟動
- 支援橫向擴展，達到公有雲級別的服務體驗

## 一站式開發平台

- 一鍵微調功能，快速適配業務場景
- 模型倉庫支援主流開源模型與自訂模型導入
- 覆蓋數據接入、訓練、管理、壓縮、加密到推理的完整生命週期

## 企業級安全合規

- 內建模型加密與數據安全機制
- 私有化部署，符合國際標準與行業合規要求



# 謝謝聆聽

---

讓每個用戶的數位化更簡單、更安全！  
[www.sangfor.com](http://www.sangfor.com)