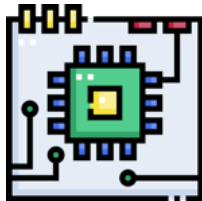# Empowering AI Transformation with AiHPC

Welcome to the AiHPC Innovation, a transformative solution designed to revolutionize the way enterprises tackle the challenges of large-scale data processing, analysis, and workflow management. Our cutting-edge platform and services can help you unlock new possibilities in your business endeavors.



AiHPC
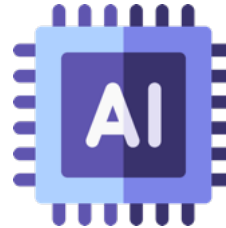ENABLER OF FUTURE INNO

# Executive Summary

Comprehensive GPU Platform for Secure, Governed Enterprise LLM Deployments

## Infrastructure Management

- Centralized GPU control with container-level allocation
- On demand HPC application deployment

## AI Capabilities

- Supporting models from LLM, NLP to OCR
- Fine-tuning and optimization pipelines

## Enterprise Framework

- Authentication integration
- Role-based access controls
- Resources monitoring

✓ **Maximize GPU Utilization** ✓ **Deploy Multiple LLM Instances**
✓ **Ensure Data Governance** ✓ **Scale with Demand**

# Navigating the AI Revolution in Enterprise Environment

**AiHPC** ENABLER OF FUTURE INNO

**Potential we have now and the challenges for the industries**

## The trend of AI cannot be reversed - Enterprises Embracing AI for Transformation

### $154.9B          25.7%

To 2030          CAGR

The enterprise AI market is growing rapidly.

**67%**
professionals regularly use AI tools

Consumers using AI have built a strong foundation for enterprise AI adoption

**96%**
Reduction in doc processing time[2]

AI-driven automation delivers dramatic substantial cost reductions and efficiency gains.

**70.0%**
Healthcare providers plan to invest in AI technologies[3]

AI enhances productivity, decision-making, and competitive advantage.

## Enterprise AI Adoption Barriers: Critical Requirements Gap

**An On-Premises solution is needed:**
Critical data and information must be stored in-house, driven by regulatory requirements and data sovereignty needs..

**Data Security & Privacy Concerns:**
Proper AuthN and AuthZ are critical in enterprise environments, requiring seamless integration with existing security systems.

**Scalability and resource allocation:**
Enterprises need efficient scalability and resource allocation for multiple AI models within limited infrastructure.

**Accountability and monitoring:**
Enterprise-level audit mechanisms and performance monitoring are essential in AI governance.
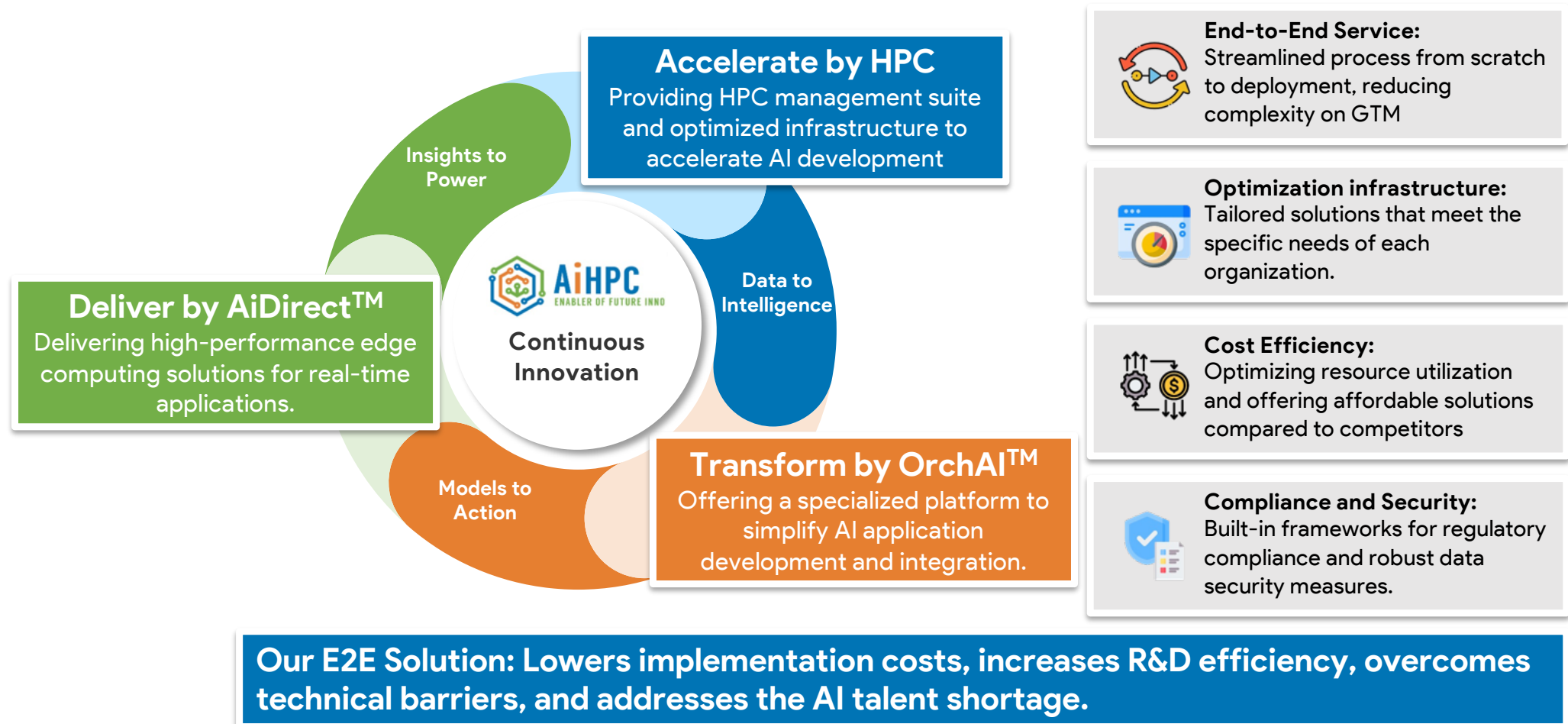
**Limited Solutions and talents:**
Shortage of mature, turnkey AI solutions for private infrastructure, compounded by the scarcity of specialized AI talent.

> **How can organizations effectively overcome the technological, regulatory, and resource barriers to adopt and scale AI securely and efficiently within their own on-premises or hybrid environments?**

# AiHPC Delivers E2E Solutions to Accelerate AI Adoption

**AiHPC**
ENABLER OF FUTURE INNO

## Our Integrated Approach Overcomes Barriers and Drives Outcomes

**Insights to Power**

**Accelerate by HPC**
Providing HPC management suite and optimized infrastructure to accelerate AI development

**Deliver by AiDirect™**
Delivering high-performance edge computing solutions for real-time applications.

**AiHPC**
ENABLER OF FUTURE INNO
**Continuous Innovation**

**Data to Intelligence**

**Models to Action**

**Transform by OrchAI™**
Offering a specialized platform to simplify AI application development and integration.

**End-to-End Service:**
Streamlined process from scratch to deployment, reducing complexity on GTM

**Optimization infrastructure:**
Tailored solutions that meet the specific needs of each organization.

**Cost Efficiency:**
Optimizing resource utilization and offering affordable solutions compared to competitors

**Compliance and Security:**
Built-in frameworks for regulatory compliance and robust data security measures.

**Our E2E Solution: Lowers implementation costs, increases R&D efficiency, overcomes technical barriers, and addresses the AI talent shortage.**

# End-to-end AI Infrastructure Solution: Simplifying Complex HPC for Enterprise AI

**Supercharge Data Processing Capabilities with Our HPC Solutions**

## AiHPC Portal
Managing AI workloads is complex → Simplified management

- Enterprise-Grade Security
- Interactive Job Control
- Dynamic Resource Allocation
- Smart Usage Analytics

## AiHPC Meter
Performance optimization is challenging → Automated optimization

- Framework-specific metrics
- Multi-aspect monitoring
- Real-time performance insights
- Optimization recommendations

## AiHPC Service
Support needs are diverse → Complete support

- Infrastructure provisioning
- Framework deployment
- Technical consultation
- Ongoing maintenance & support

**Powered by:**

jupyter · mxnet · ONNX · mlflow · R Studio · slurm workload manager · kubernetes · S · Docker · hadoop · BeeGFS · DVC · DELL · Hewlett Packard Enterprise · NVIDIA · intel · HAILO · AMD

**Selected stories:**

- Deployment of AiHPC Portal at **100+ nodes** HPC environment for interactive jobs and billing system.
- Designed and implemented HPC architecture for **Genomic Research**
- Provided AiHPC Portal for enhanced resource and workflow management.
- Designed and developed **AI software-orientated GPU card performance evaluation** under the liquid cooling environment with AiHPC Meter suite.

CityU

**We empower AI availability with highly flexible, cost-efficient, user-friendly enterprise HPC solutions.**

# Technology: Enterprise-Grade Infrastructure Management

**AiHPC** ENABLER OF FUTURE INNO

## What We Believe is Maximizing Hardware Investment with Granular Control

### On-Demand Resource Allocation

Dynamic allocation of GPU resources based on workload demands

Efficient containerization enables a 300% improvement in GPU utilization rates[1]

### Fine-grained GPU Partitioning

Centralized administration with API integration

Dynamic GPU sharing shows 4.5x improvement in specific research environments.[2]

### Enterprise Integration

Role-based access control and seamless Integration with existing identity management systems

- Supports LDAP, OAuth, and other enterprise authentication standards
- Automated user provisioning and de-provisioning workflows

### Monitoring and Tracking

Comprehensive usage analytics and reporting Cost allocation and chargeback capabilities

- Real-time resource utilization visualizations and historical performance trending, and predictive analytics

NGC APPLICATION CONTAINER
GPU-ACCELERATED APPLICATION
REQUIRED LIBRARIES
CUDA-X AI
CUDA TOOLKIT

MAPPED NVIDIA DRIVER
CONTAINER BASE IMAGE (OS, LIBRARIES)

KUBERNETES

NVIDIA GPU OPERATOR
NVIDIA GPU DEVICE PLUG-IN FOR KUBERNETES
NVIDIA CONTAINER RUNTIME
NVIDIA DRIVER

CONTAINER ENGINE
(DOCKER, CRI-O, SINGULARITY)
LINUX HOST OS

GPU-ACCELERATED CLUSTER
GPU-ACCELERATED SERVERS
NVIDIA GPUs (T4, V100)

**A well-architected HPC infrastructure is the foundation determining the success of the entire AI journey.**

[1] NVIDIA. (2023). Maximizing GPU utilization for AI workloads with MIG technology.
[2] Zhang, Y., et al. (2021). Dynamic GPU sharing for deep learning in Kubernetes.

# AiHPC Portal to Manage AI Data Flow and Computing for Research Institute

**Bridging Researchers, Data, and Computing Power Through Intelligent Orchestration**



**Distributed DBMS**

**Business Workflow Development**

**Menu Systems**

**Tools set**

GUI

Slurm Interactive Session

Terminal

Slurm Script/Job Script

**Presentations Layer**

**Storage Framework**

**Computing Pool**

By unifying access, computation, and storage, the AiHPC Portal creates the one-stop solution optimizing institutional infrastructure investments.

# AiHPC Portal in Action: Streamlining Job Management and Workflow Management

**Demonstration of the HPC management with a GUI platform**



Demonstrates comprehensive HPC job monitoring with real-time tracking of computing resources, partition utilization, and historical job retrieval capabilities.

Showcases script-free launching of scientific applications, including Jupyter Lab, Notebook, and RStudio, enabling researchers to instantly access computing resources without technical barriers.

# AiHPC Portal in Action: Streamlining Job Management and Workflow Management

**Demonstration of the HPC management with a GUI platform**



Example usage: The Portal provided comprehensive workflow monitoring with real-time status tracking, resource utilization metrics, and process-level performance analysis for computational biology applications.



Comprehensive hardware monitoring dashboards display real-time GPU performance metrics including temperature, power consumption, and utilization trends for optimizing computational resource efficiency.

# Transform AI Integration with OrchAI™ Platform

Unified Orchestration for Enterprise AI Operations

**Data Integration → Model Selection → Monitoring → Optimization**

## Quality Local Data Hub
Connect & Transform

- Data Connectors
- ETL Pipeline
- Data Validation
- Real-time Sync
- Security Assurance

## One System, multiple AI
Store & Utilize

- Multi-model Registry
- Version Control
- RAG Integration
- Model Marketplace

**OrchAI**
**One Stop Enterprise AI Platform**

## Enterprise Security and Monitoring
Review & Improve

- Performance Metrics
- Cost Analysis
- Usage Tracking
- Real-time Alerts

## Interaction with your AI
Continues & Reproduce

- Workflow Engine
- Task Scheduling
- Auto-scaling
- Pipeline Management

✓ **Rapid AI Deployment -** From concept to production in days ✓ **Data Security:** Your data stays within your control
✓ **Scalable Architecture:** Grows with your enterprise needs ✓ **Full Transparency:** Track every AI decision and interaction

One System, Multiple AI Models

**Performance and Usage Monitoring**

# Local Data Hub Integration

# OrchAI Hybrid Architecture: Intelligent Workload Distribution

**AiHPC**
ENABLER OF FUTURE INNO

**Utilizing the resources to gain the greatest user experience in AI adoption**

## Sever Solution

### Centralize Management

- Knowledge hub
- AuthN, AuthZ
- AI workflow orchestration
- Resources management

### Intensive AI Workflow

- Complex training
- Large parameter models
- Multiple GPU inference
- Model repository
- Version control

**OrchAI**

- Workload classification
- Resource availability
- On demand routing
- Task priority management.

**Sending AI Tools to Your PC** →

← **Sending Big Jobs to the Server**

## PC Solution

### Edge Computing

- Fast model access
- Quick responses
- Simple deployment
- Local processing

### Specific AI Agent

- AI Writing & Email Assistants
- Smart Document Preparation
- Interactive Data Queries
- Real-time Code Completion

### Scenario

- Classroom
- Office
- Laboratory
- Workstation

**OrchAI: The right AI task, on the right device, at the right time.**

# The AiDirect™ – Deliver AI Models to Real-World Applications

## The Professional AI Channel for Enterprise Innovation

AiHPC
ENABLER OF FUTURE INNO

**Model Producer** → Contribute → **AiDirect** → Delivery (API, Container, Edge) → **Real World Impact**

**Model Producer** ← Review ← **AiDirect** ← Collaboration ← **Real World Impact**

AiDirect is a **domain specific AI model distribution platform** designed to deliver AI models into real-world applications.

Similar to how **LinkedIn** caters to **professional networking** compared to **Facebook**, **AiDirect** focuses on professional, secure, and compliant AI deployment, whereas platforms like Poe cater to general AI usage.

vs AiDirect

**Professional-Grade AI Model Library:**
- Access a **curated selection** of AI models developed for professional applications in different industries.
- Models are vetted for **compliance, performance, and reliability** in operational settings.

**Seamless Integration into Real-World Applications:**
- Efficient deployment AI models with **AI Agent** available.
- Supports various integration methods, including **APIs, containers, and edge deployment**.

**Compliance and Security:**
- Data Protection: Comply with relevant industry regulations (e.g., GDPR, CCPA, and others) and data privacy mandates**.**
- Security Protocols: Implements robust security measures to protect sensitive enterprise data and system integrity.

**Collaborative Platform:**
Enables collaboration between AI developers, domain experts, and organizations to foster innovation and address real-world business challenges**.**

**Intelligence Decision Support Systems:**
Augment expert decision-making with AI-driven insights for improved accuracy and efficiency across various functions..

**Medical Imaging Analysis:**
Enhances imaging technologies, providing faster and more precise interpretations of scans.

**Streamlined Processes:**
Automation scheduling, billing, and other administrative tasks, increasing efficiency.

17

# Accelerating AI Deployment with Purpose-Built Edge Computing

**Bringing AI to the Point of Care with Enterprise-Grade Performance**

- **Customized Edge Devices** – e.g. optimized for medical imaging workloads
- **Real-Time Analytics** - 95% faster than traditional cloud processing
- **Secure Deployment** – on-premises and e2e encryption
- **Remote Management** - for enterprise-wide monitoring

## Traditional AI Deployment vs. Edge Approach:

| | Traditional | AiDirect |
|---|---|---|
| Processing Time | Minutes | Milliseconds |
| Data Privacy | Network Risk | On-Premise |
| IT Complexity | High | Simplified |
| Deployment Cost | Per-Use | Fixed Cost |



**PCIe x16 AI Acceleration Card Hailo-8™ AI Processor x4**

✓ **Future-Proof Investment:** Scale as AI models evolve ✓ **Reduced Bandwidth Costs:** Process locally, only send results to cloud
✓ **No Vendor Lock-in:** Standard ML framework compatibility ✓ **Simplified Compliance:** Streamlined data governance and auditing

# Our Dynamic AI Ecosystem: Driving Innovation & Rapid Industry Adoption

**Making advanced AI accessible and rapidly deployable across industries.**

## HPC Consultancy → AI Implementation → System Integration → Frontend Support Services

### (The Technology & Expertise Core)

- **Core Technology & IP:** AiHPC Portal, AiHPC Meter, OrchAI, AiDirect
- **Empowerment & Advanced Services:** Deep AI/HPC/domain expertise, solution architecture, specialized technical support
- **Market-Driven R&D & Co-Development:** Actively incorporates industry feedback and engages in joint R&D projects to enhance existing solutions and pioneer new ones.

### Strategic SI Partners
#### (Market Access & Operational Excellence)

- Industry Qualification & Trust:
- Extensive Sales & Distribution Networks:
- Local Implementation & Integration:
- Frontend Support & Customer Management:

**Tech & Expert**

**Co-development**

**GTM, Delivery, Engagement**

**Target Industries: Government, Finance, Healthcare, Education ...**
(Delivering Enterprise AI solutions fit to specific industry demand)
**Accelerated Time-to-Market:** Rapidly deploy proven, enterprise-grade AI solutions across industries.
**Scalable Delivery Models (B2B / B2B2C):** Flexible engagement delivering tailored solutions.
**Target Industries Served:** Providing specialized AI solutions to meet specific vertical needs.

## AiHPC's Technical Innovation + SI's Market Expertise = Accelerated Customer Success

## AiHPC's core technology, amplified by SI partnerships and enriched by industry co-development, delivers continuously evolving AI solutions.

# Founding Members

**Pioneering Enterprise AI with Deep Technical and Market Expertise**



## Sam Chu, Ph.D.
### Chief Executive Officer, Founder

Sam Chu science and technology entrepreneur, serving as CIO for Sanomics Limited and ACT Genomics, overseeing a clinical trial on early-stage cancer detection using plasma DNA tests. Ensuring efficient data integration, storage, and analysis, they have also played a key role in developing the tele-oncology platform, Aurora Tele-Oncology, enabling remote cancer care during the pandemic.

## Sammy Tang
### Chief Technology Officer, Co-Founder

Sammy Tang is a highly skilled professional with expertise in mathematical modeling, scientific problem-solving, HPC research, and clinical solution deployment. With a track record of managing engineering support science projects in areas such as data centers, air pollution studies, and ocean modeling, they have made significant contributions to research and development in various fields

## Mark Chang, D.Phil.
### Chief Strategy Officer, Co-Founder

Mark Chang is a highly accomplished individual with a strong background in healthcare, technology, and government initiatives in Taiwan. With roles such as Taiwan Managing Partner, Assistant Professor, Committee Member, and Chief Information Officer, he contributes significantly to the advancement of healthcare and technology integration in Taiwan..

## Hong Kong • Taiwan • Singapore

## Providing local support with global standards across key Asian technology and healthcare hubs

# Summary: Total Solution of Cloud/HPC/AI



**AiHPC Portal**
**AiHPC Meter**
**Software Products**

Public Pre-trained AI Models

User Applications

User data

on-line

AI Model Training

AI Model Deployment

off-line

<u>Cloud/HPC Services:</u>
- Design
- Implement
- Optimize
- Maintenance

<u>AI Factory Services:</u>
- Data Processing
- Model Training
- Model Deployment
  - User Application
  - PCIe Card

<u>Domain:</u>
Healthcare/Manufacturing/
Retail/Security/Finance/...

**Edge Inferencing Cards**