



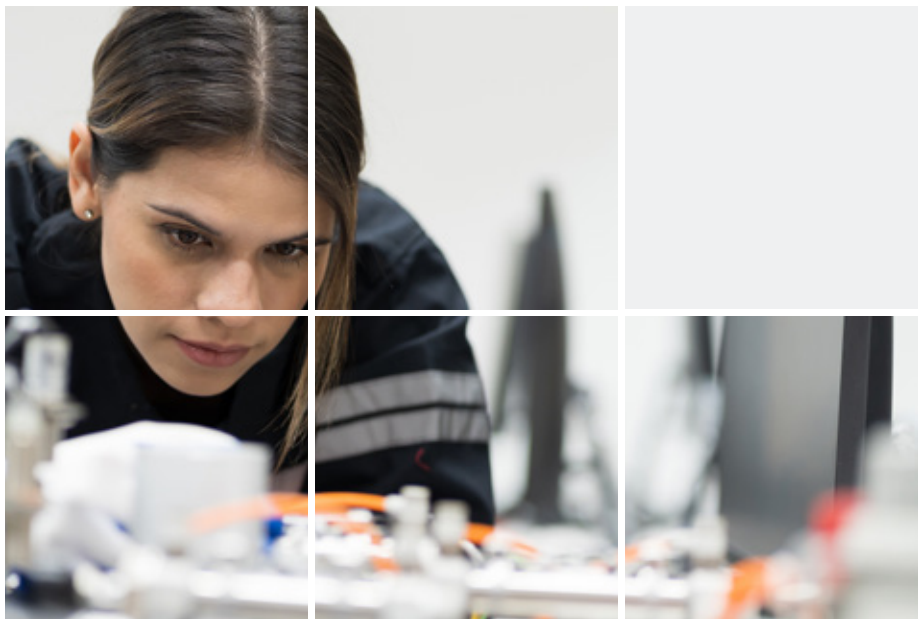
*NEXT Frontier: NEXTDC's AI in Education Series*

# Building Your University's AI Research Powerhouse

Why Next-Gen Infrastructure is Essential



**N E X T D C**  
where AI lives™



# Defining the AI Research Powerhouse: **A Strategic Imperative for Universities**

In the rapidly accelerating world of Artificial Intelligence, universities face a critical choice: merely observe the revolution or actively lead it. The ambition to push the boundaries of machine learning, train sophisticated large language models, and conduct groundbreaking data-intensive research hinges entirely on one foundational element:

**cutting-edge AI data centre infrastructure.**

Think of it not just as a collection of powerful machines, but as your **university's AI research powerhouse**, a dedicated engine designed to generate unprecedented discovery and cultivate the next generation of innovators.

This powerhouse is where AI tokens are processed at scale, where complex algorithms are refined in hours instead of weeks, and where data transforms into breakthrough insights. It's the engine room for academic excellence in the age of AI. But what does it take to build and sustain such a powerhouse, and why is doing so a strategic imperative, not just an IT project?





# The Vision: What Defines a University's AI Research Powerhouse?

An AI research powerhouse isn't simply a server room; it's a strategically designed ecosystem built for the unique demands of advanced AI. It encompasses:

## Immense Compute Power:

Primarily driven by high-density GPU clusters, capable of handling the parallel processing needs of deep learning, generative AI, and complex simulations. This means supporting the latest NVIDIA Hopper-generation chips (H100, H200) and preparing for future advancements, often in configurations that scale to NVIDIA DGX SuperPOD™ architectures.

## Vast, Accessible Data Storage:

Petabytes of high-speed storage, optimised for AI workloads, ensuring researchers have immediate access to the enormous datasets required for modern models.

## High-Bandwidth, Low-Latency Networking:

An internal network architecture designed to move massive amounts of data quickly between GPUs and storage, eliminating bottlenecks that can cripple research. NVIDIA's NVLink technology, for example, allows H100 and H200 GPUs to communicate at 900 GB/s within a server, and InfiniBand provides high-speed connectivity between servers in a cluster, **forming the backbone of scalable AI architectures like DGX SuperPOD.**

## Specialised Software & Tools:

A curated suite of AI frameworks, libraries, and orchestration tools that empower researchers to maximise the hardware's potential.

## Scalability and Agility:

The ability to rapidly expand compute resources as research demands grow and technology evolves, without prohibitive delays or costs. This includes supporting the modular scalability of AI Factory deployments.

## Expert Support:

A dedicated team of technical professionals to manage, maintain, and optimise this complex environment, ensuring maximum uptime and accessibility for researchers.

**This integrated approach enables universities to conduct research that might otherwise be impossible, from training custom foundation models tailored to specific academic disciplines to running massive climate simulations or accelerating drug discovery.** It's about empowering researchers to go beyond the confines of traditional labs and unleash exponential impact.



**Understanding Tokens:** In AI, particularly with Large Language Models (LLMs), a "token" is the fundamental unit of text (like a word, part of a word, or punctuation mark) that the AI model processes. When you input text, it's broken down into tokens, and the model then works with these tokens to understand context, generate responses, or perform analyses. The sheer volume of tokens processed during both training and inference directly correlates with the computational power required. For example, a single prompt for a generative AI model might involve hundreds or thousands of input tokens, and the output could generate many more.

**Powering Training & Inference:** Training large AI models like LLMs involves feeding them billions or trillions of tokens. This process is incredibly compute-intensive and demands sustained, maximum GPU performance for weeks or months. For instance, an NVIDIA H100 GPU can consume up to 700 Watts (W) of power, and the newer H200 GPU operates within a similar power envelope while offering significantly increased memory (141GB HBM3e vs. 80GB HBM3 on H100) and bandwidth (4.8 TB/s vs. 3.35 TB/s). This enhanced memory and bandwidth specifically doubles inference speed on models like Llama 2 70B compared to the H100 and provides substantial accelerations for training. While training demands peak power for extended periods, inference workloads, even though typically less intense per query, still require significant power, especially when serving many concurrent users or processing large batches of data rapidly.



# The Peril of Going It Alone: Why Building On-Premise is Often a Roadblock

The aspiration to build such an AI research powerhouse is strong, but the practicalities of doing so entirely on-premise present formidable, often prohibitive, challenges for universities:

## Massive Capital Expenditure (CAPEX):

The upfront cost of land, construction, specialised building modifications (for power, cooling, fire suppression), and the acquisition of high-end GPUs, servers, and networking gear for SuperPOD-scale deployments can run into tens or hundreds of millions. This diverts crucial funds from other academic priorities.

## Operational Complexity and Cost (OPEX):

Managing a 24/7, mission-critical data centre is a specialised undertaking. Universities often face:

- **Extreme Power Requirements:** As noted, an H100 or H200 GPU can draw up to 700W. A single rack housing 8-10 such GPUs can easily demand 7-10 kW of power, with dense AI racks now pushing well beyond 50 kW. An NVIDIA DGX SuperPOD can consume hundreds of kilowatts, or even megawatts, for its compute footprint alone. Providing this reliable, redundant power supply is a significant infrastructure challenge.
- **Advanced Cooling Challenges:** The concentrated heat generated by these powerful chips is immense. Traditional air cooling often becomes insufficient for high-density AI clusters. Modern solutions like direct liquid cooling (cold plates) or rear-door heat exchangers are often required to efficiently dissipate heat and maintain optimal operating temperatures. Implementing and maintaining these complex cooling systems adds significant cost and technical complexity, especially at the scale of an AI Factory.
- **Physical Security:** Protecting valuable hardware and sensitive research data requires multi-layered security protocols that most university campuses aren't equipped to provide at data centre scale.
- **Specialised Staffing:** Recruiting and retaining skilled engineers proficient in HPC, AI infrastructure, power, cooling, and networking is a constant challenge, as these roles are in high demand across all industries.
- **Maintenance & Upgrades:** The continuous cycle of hardware failure, upgrades, and patching is a relentless burden on internal IT teams.

## Scalability Headaches:

Predicting future AI compute needs is notoriously difficult. Building too small means rapid obsolescence; building too large means wasted capital and stranded assets. Expanding on-premise, particularly to SuperPOD or AI Factory scale, is often a disruptive, multi-year construction project.

## Supply Chain Volatility:

The global chip shortage and ongoing supply chain disruptions mean universities often wait months or even years for critical GPU hardware, delaying research and missing crucial funding opportunities.

## Sustainability Challenges:

Powering a high-density AI cluster on campus without the efficiencies of a purpose-built facility can lead to a significant environmental footprint and make it difficult to meet institutional sustainability goals.



**These challenges illustrate why attempting to “go it alone” in building an AI research powerhouse can quickly turn into a financial drain, an operational nightmare, and a bottleneck for research, rather than an enabler.**





# The Smart Path Forward: Partnering for Your AI Research Powerhouse

For universities committed to leading in AI, a strategic partnership with a specialised data centre provider offers a compelling solution, transforming these formidable challenges into manageable opportunities. By leveraging colocation, universities can:

## Shift from CAPEX to OPEX:

Reduce massive upfront capital investment by utilising a provider's existing, purpose-built infrastructure.

## Access World-Class Facilities:

Benefit from redundant power, advanced cooling (including the capability to support high-density racks with liquid cooling systems like rear-door heat exchangers or direct-to-chip solutions, critical for handling the 700W+ thermal loads of H100 and H200 architectures), and robust physical security that few universities can afford to build or operate in-house. These facilities are designed to host NVIDIA DGX SuperPODs and AI Factories.

## Ensure Operational Excellence:

Offload the burden of 24/7 monitoring, maintenance, and facility management to expert teams, freeing up university IT staff to focus on supporting researchers directly.

## Achieve Unprecedented Scalability:

Rapidly scale compute resources up or down on demand, aligning AI data centre infrastructure with evolving research needs without disruptive construction or lengthy procurement cycles. Providers like NEXTDC are specifically engineered to support the ultra-high-density power and cooling demands of NVIDIA Hopper (H100, H200) architectures, including DGX SuperPODs, ensuring seamless integration and future-proofing for next-generation systems like Blackwell.influence.



## Enhance Connectivity:

Gain direct, high-speed access to major cloud providers, research networks, and other institutions within a vibrant interconnection ecosystem, fostering collaboration and data sharing.

## Boost Sustainability:

Leverage a provider's energy-efficient designs and renewable energy commitments to meet institutional environmental goals, often achieving superior Power Usage Effectiveness (PUE) compared to custom on-premise builds.

**This strategic approach allows universities to establish their AI research powerhouse rapidly and efficiently, providing the power, resilience, and scalability needed to drive breakthroughs and attract top talent. It's about choosing a foundation that empowers your academics, rather than burdening your administration.**



# Ready to build **Your University's AI Powerhouse?**

In today's AI-driven world, infrastructure isn't just a technical consideration—it's a strategic differentiator. If your institution is serious about leading in breakthrough discovery, attracting global talent, and winning major research grants, the foundation must be built now.

**"A true AI research powerhouse empowers your academics, not burdens your administration."**

As La Trobe University has demonstrated, pairing NVIDIA DGX H200 systems with purpose-built environments accelerates discovery in ways once thought impossible. But building and operating this infrastructure alone is costly, complex, and slow. That's why forward-thinking universities are turning to strategic colocation with trusted partners.

Leveraging a specialised provider like NEXTDC for your H100/H200 deployment allows your university to maximise its research ambitions without the complexities of building and managing bespoke infrastructure. **Consider these advantages:**

- **NEXTDC is NVIDIA DGX Data Centre Ready certified**

This certification confirms our facilities are purpose-built to support the most demanding AI workloads, ensuring optimal performance and reliability for your NVIDIA DGX systems, including NVIDIA DGX SuperPOD™ deployments and AI Factory architectures.

- **Shift from CapEx to OpEx**

Reduce massive upfront capital investment by utilising a provider's existing, purpose-built infrastructure. This delivers a **lower Total Cost of Ownership (TCO)** over the long term, avoiding the escalating costs of on-campus expansion.

- **Certified, AI-Optimised Facilities**

Built to support high-density NVIDIA H100/H200 racks with advanced cooling and power efficiency. This offers immediate access to scalable power and cooling without delays or upfront capital expenditure.

- **Scalable On-Demand Compute**

Grow capacity as your research expands without waiting on construction or hardware availability. NEXTDC's modular, rapidly scalable infrastructure is designed for AI growth, unlike limited on-campus options.

- **Expert Uptime & Operational Management**

Let specialists handle power, cooling, security, and compliance, offloading the operational burden from your internal IT teams and ensuring maximum uptime. This effectively **transfers the risk of planning, construction, and delivery of the data centre to NEXTDC, not the university.**

- **Global AI Routing & Interconnection**

NEXTDC's proximity to international subsea cable landing stations enables ultra-low-latency cross-border AI workloads, including federated learning, multi-region training, and edge inferencing. This provides a **strategic edge through next-generation connectivity**, including direct, high-speed links to research networks like AARNet.

- **Sustainability Alignment**

Hit your ESG targets faster with energy-efficient, renewably powered data centre operations. This helps overcome the **sustainability challenges** of energy-intensive on-campus data centres and contributes directly to Net Zero targets.

- **Free Up Campus Space**

By shifting backend IT off-site, universities can **reallocate valuable campus real estate** to core academic priorities such as teaching, research labs, or student services.



# Accelerate Your University's AI Leadership

NEXTDC's NVIDIA DGX-certified data centres are purpose-built to support the most demanding AI workloads, from foundational model training to real-time inference. With GPU-optimised designs, super high-density power, and advanced cooling architectures, we deliver the infrastructure that accelerates every stage of the AI lifecycle.

Strategically located near Australia's major international subsea cable landing stations, NEXTDC enables ultra-low-latency pathways for cross-border GPU workloads, including federated learning, multi-region training, and edge inferencing. For Australian universities building global AI platforms or deploying GPU-as-a-Service at scale, this means your AI workloads are just milliseconds away from key markets across Asia and the Pacific. This strategic connectivity is the hidden foundation for academic success in an AI-driven world.

The institutions that ask better questions today will lead with better outcomes tomorrow.

Whether you're advancing medical research, shaping the next generation of AI talent, or building a new innovation precinct—your infrastructure must be as ambitious as your mission. milliseconds away from key markets across Asia and the Pacific.

**Partner with NEXTDC to unlock secure, scalable, and sustainable access to NVIDIA H100/H200-powered AI environments, certified, sovereign, and ready for immediate impact, ensuring your AI infrastructure acts as a catalyst for discovery, not a bottleneck, supported by our unique combination of technical infrastructure and campus strategy expertise.**

Download our AI Ready Checklist



Speak with a NEXTDC specialist



**N E X T D C**  
where the cloud lives™

136 398

sales@nextdc.com

**nextdc.com**

This document is correct at the time of printing and is for presentation purposes only. This document does not constitute an offer, inducement, representation, warranty, agreement or contract. All information contained in this document (including all measurements, photographs, pictures, artist's impressions and illustrations) is indicative only and subject to change without notice. NEXTDC Limited, its employees, representatives, consultants and agents make no representations or warranties as to the accuracy, completeness, currency or relevance of any information contained in this document and accept no responsibility or liability whatsoever for any discrepancy between the information contained in this document and the actual data centres or services provided by NEXTDC Limited or for any action taken by any person, or any loss or damage suffered by any person, in reliance upon the information contained in this document. © 2025 NEXTDC Limited ABN 35 143 582 521.

UNI05\_2025\_030825\_01