# define tech

# LMX AI

## AN AI REFERENCE ARCHITECTURE

*Best-practice suggestions for compute, networking, storage, power, cooling, and more, in an integrated AI infrastructure that serves to progress your AI strategy.*

## ABSTRACT

If 2020 has taught us anything it is that the future of our world is deeply dependant on advancing technology. From fields such as life sciences, manufacturing, materials science, medicine, energy and beyond – technology is not only helping to improve our quality of life but helping us to solve critical, global problems.

Advancing time to insight in these fields is the holy grail that all research facilities are seeking, and as a result organisations are looking to Artificial Intelligence (AI) to accelerate their research initiatives.

Of course with this AI-fueled innovation comes the enormous pressure on underlying infrastructure and the costs associated with continually upgrading legacy hardware to run AI workloads.

The LMX AI reference architecture features our best-practice suggestions for compute, networking, storage, power, cooling, and more, in an integrated AI infrastructure that serves to progress your AI strategy without blowing the budget.

The Age of AI is here. You may be ready, but is your infrastructure? With LMX Cloud, no matter the industry you work in - you can future-proof your organisations' AI strategy.

## INTRODUCTION

We live in a data-driven world. Modern workloads are adapting AI tools and techniques to shift through enormous data volumes to extract insight or enhance applications and solutions. With fast growing datasets and broadening application requirements, the dependence of our technology on AI has never been greater and will only continue on that trajectory.

Any successful AI strategy should include a platform that provides the agility to respond and scale to any workload demands ranging from compute intensive simulations to the processing of large data sets.  See fig.1.

## AI WORKFLOW

In order to design a best-fit AI reference architecture it is imperative to understand the full AI workflow life-cycle. See fig. 2 overleaf.

The AI life-cycle begins at the data capture phase. Depending on the organisations function, the AI life-cycle can include the ability to capture vastly different data formats including audio, video, images, multi-omic or sensor generated, and from multiple sources including publicly available, collaborative datasets.

This data is then ingested into the AI platforms' storage or data-lake. This ingestion can have varying requirements dependent on the type of data being ingested. For example the requirement of different storage protocols, ingestion rate and frequency – the platform should be able to accommodate all.

Preprocessing is an automated feature which can process the data as it is being ingested based on predefined policies (for example the addition of meta data based on the data-type and source).

Enrichment is the manual process of adding additional meta data to the datasets and usually requires some form of human interaction. It is this step that will enrich the data with information that will ultimately help to train the neural network.

Algorithm and model development is where researchers can run test simulations for the evaluation of their AI models and algorithms. These test runs can be deployed in virtual environments before rolling out the algorithms to production. The environment used needs to be able to handle the provisioning of multiple stacks in an automated manner.
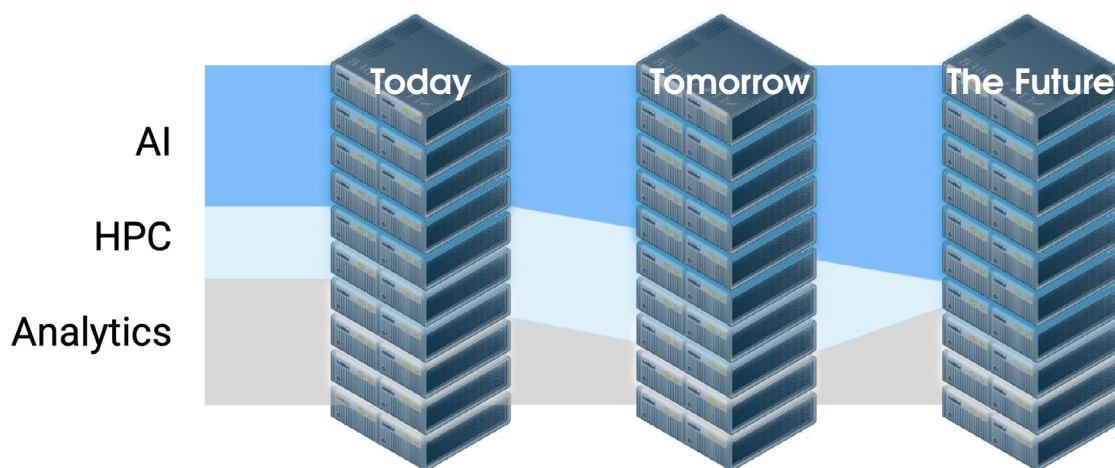


*Fig 1. Modern Infrastructure need the agility to quickly respond to changing workloads*

**Deployment**
Take models and run inference (including at the edge)

**Validation**
Measure the accuracy of the model against unknown datasets

**Simulation & Training**
Data and Compute Intensive phase

**Algo Development**
Test framework, lots of small concurrent simulations

**Data Capture**
Audio, Video, Images or sensor data

**Data Ingest**
Multi-protocol access to Data Lake

**Pre-processing**
Automated meta data extraction / file conversion or manipulation

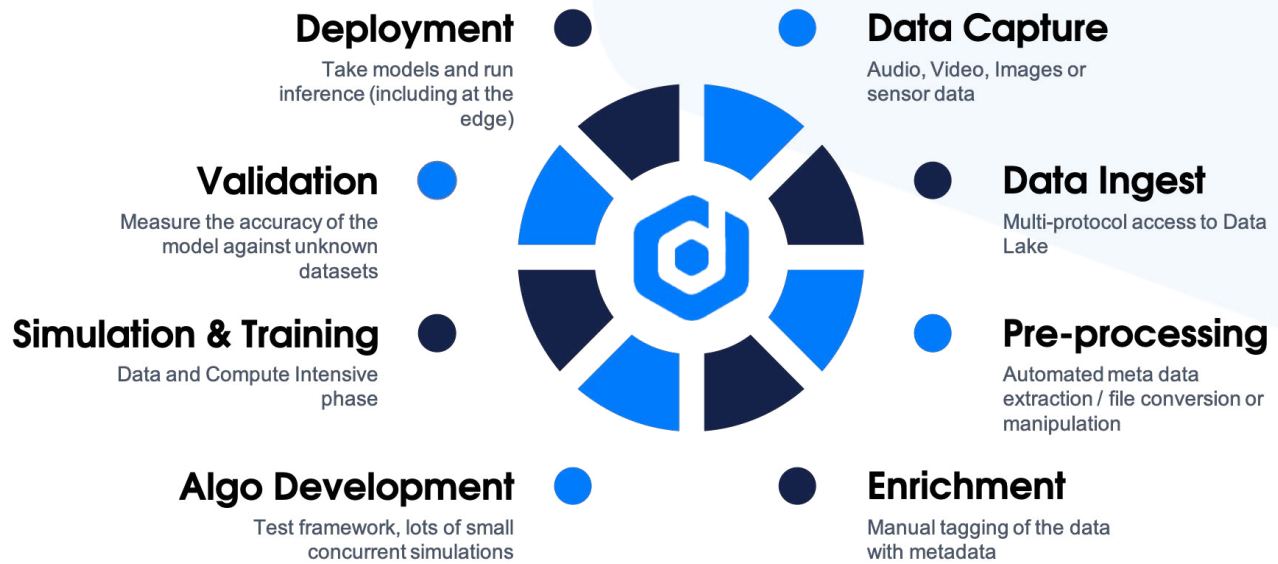**Enrichment**
Manual tagging of the data with metadata

*Fig 2. The AI workload Life-cycle*

As the name suggests, Simulation and Training is when the models finally get trained on the neural network. Depending on the size of the dataset and the neural network depth, this training can take up time and resource, which is why it is important to have a platform that can scale on-demand to meet workload requirements. Performance is critical here and a balanced architecture which provides fast data access to labeled data sets and also provides high performance access to compute resources such as GPUs, is essential.

Validation is the verification of the training model against an unknown dataset. And finally, deployment is the rolling out of the entire AI solution.

## LMX SOLUTION OVERVIEW
With LMX, we have architected a high-performance, scalable AI solution that's optimised, affordable and easily scalable. Powered by NVIDIA DGX™ A100 GPUs, our own AI optimised Cloud software and Lightbits LightOS™ NVMe storage. The platform includes NVIDIA® Mellanox Spectrum® ethernet and NVIDIA Mellanox Quantum™ InfiniBand switches,

delivering excellent linear scaling for training and inference workloads. The solution combines the best of breed hardware with software that puts the user in control of how the environment is configured and tuned.

## SOFTWARE - LMX CLOUD
LMX Cloud is a comprehensive Cloud Software stack, optimised for AI and Deep Learning that supports a broad range of workloads and software environments, enabling organisations with an agile and scalable AI-driven infrastructure.

With LMX Cloud, you get the flexibility of a software-defined architecture with accelerated hardware, optimised storage and pre-integrated application stacks so that you can focus on harnessing the most value from your data.

LMX is designed to remove the administrative burden of running an infrastructure for next generation workloads and empowers users to create appropriate environments for conducting AI exploration, investigation, creation and deployment.

## KEY FEATURES

- Remote Visualisation Support to help with data set labeling or visualisation operations.
- Support for vGPU, MIG (Multi-instance GPU), PCI pass-through GPU and bare metal provisioning
- Full support for containerised workloads and applications from NGC (NVIDIA GPU Cloud)
- Complete AI user environment (Kubernetes (K8s) support, SLURM support, monitoring, scientific libraries, compilers, profilers, debuggers)
- Control infrastructure via open cloud APIs
- Manage entire infrastructure including compute, storage and networking from a single interface.
- Comprehensive monitoring and alerting
- Support for virtual machines (for training and POC) as well as bare metal provisioning
- Web UI Portal with support for file transfers, workload management, and on demand VNC, RStudio and Jupyter support.
- No single point of Failure / Zero touch provisioning / Rolling upgrades / Zero downtime



*Fig 3. LMX Cloud Software Interface*

## STORAGE - LIGHTBITS LIGHTOS™

GPUs are computationally very capably devices and to ensure they are kept busy and running efficiently, you need to ensure the storage back-end can serve up data sets as fast as the GPUs can process them.

LMX Cloud has integrated Lightbits LightOS™ - a software-defined block storage solution that delivers composable, high-performance, scale-out and redundant NVMe/TCP storage that performs like local flash.

## ACCELERATION - NVIDIA DGX-2 A100

Developed to meet the demands of AI and analytics, NVIDIA DGX™ Systems are built on the revolutionary NVIDIA Volta™ GPU platform. Combined with innovative GPU-optimized software and simplified management tools, these fully-integrated solutions deliver ground-breaking performance and results. NVIDIA DGX Systems are designed to give data scientists the most powerful tools for AI exploration—from your desktop to the data center to the cloud.

NVIDIA DGX™ A100 is the universal system for all AI workloads, offering unprecedented compute density, performance, and flexibility in the world's first 5 petaFLOPS AI system. NVIDIA DGX A100 features the world's most advanced accelerator, the NVIDIA A100 Tensor Core GPU, enabling enterprises to consolidate training, inference, and analytics into a unified, easy-to-deploy AI infrastructure that includes direct access to NVIDIA AI experts.

The NVIDIA A100 Tensor Core GPU delivers unprecedented acceleration for AI, data analytics, and high-performance computing (HPC) to tackle the world's toughest computing challenges. With third-generation NVIDIA Tensor Cores providing a huge performance boost, the A100 GPU can efficiently scale up to the thousands or, with Multi-Instance GPU, be allocated as seven smaller, dedicated instances to accelerate workloads of all sizes.
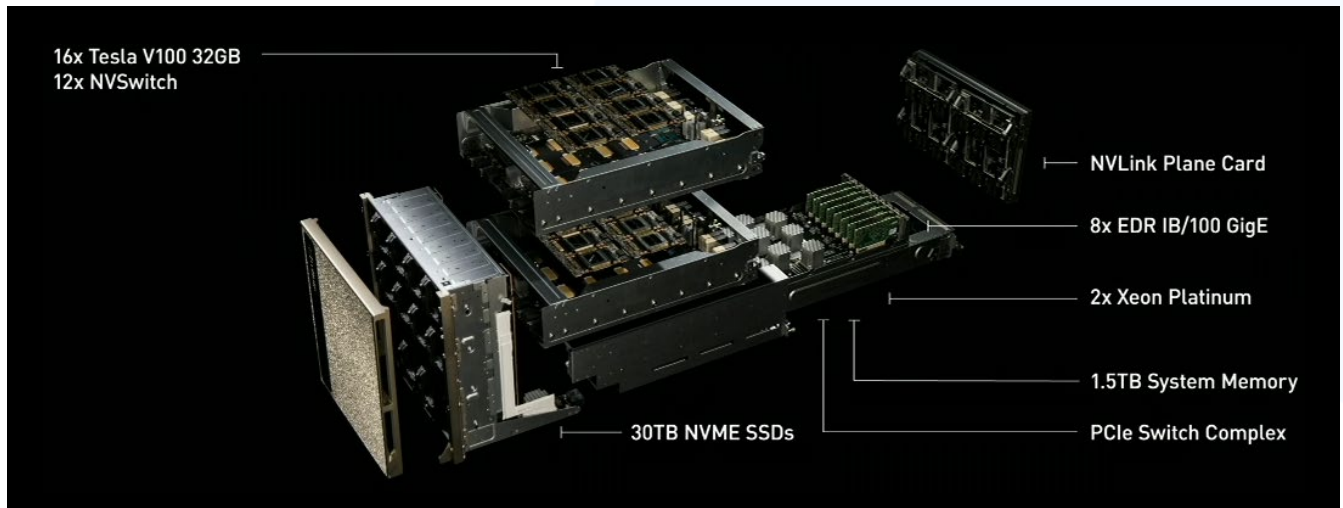
*Fig 4. NVIDIA DGX2*

## FABRIC - NVIDIA MELLANOX RDMA AI

GPU accelerated computing and ever-scaling Deep Learning/Machine Learning workloads are posing a unique challenge to network architects looking to design the perfect interconnect fabric. Efficient and sustainable scaling requires expanding the role of interconnect beyond standard message-passing agent to a more intelligent entity that can accelerate the overall compute process.

Designed specifically for the needs of GPU acceleration, NVIDIA Mellanox GPUDirect® RDMA provides direct communication between NVIDIA GPUs in remote systems. This eliminates the system CPUs and required buffer copies of data via the system memory, resulting in 10X better performance. See fig. 5.

## AI FRAMEWORKS  - TENSORFLOW

The LMX Cloud platform comes pre-integrated with the latest versions of industry standard AI frameworks such as Tensorflow, Caffe and Theano among others. These frameworks can be deployed quickly and easily from our built-in container registry so users waste no time in setting up their AI environments.

## APPLICATION REPOSITORY - NGC

LMX Cloud comes pre-integrated with NVIDIA GPU Cloud (NGC™). The NGC catalogue is the hub for GPU-optimized software for deep learning (DL), machine learning (ML), that accelerates deployment to development workflows so data scientists, developers, and researchers can focus on building solutions, gathering insights, and delivering business value.
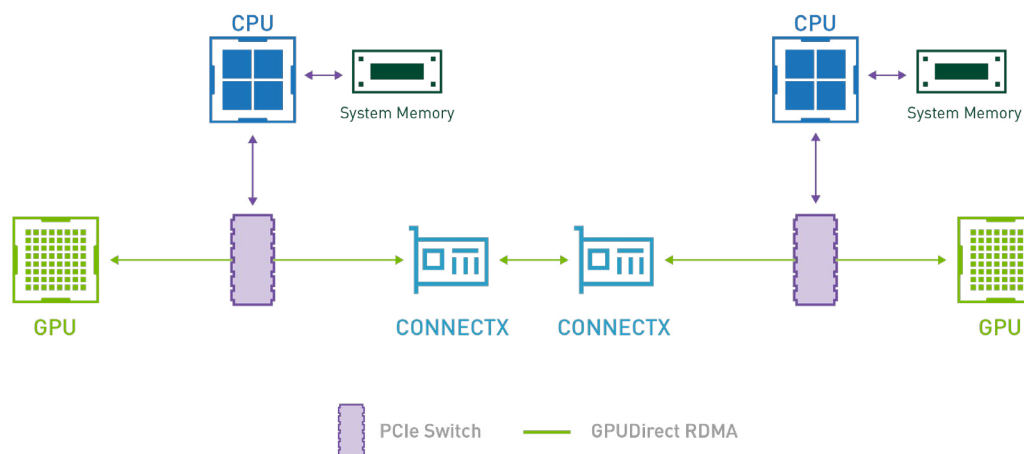


*Fig 5. NVIDIA Mellanox RDMA Fabric*

The NGC catalog provides a comprehensive hub of GPU-accelerated containers for AI, and machine learning that are optimized, tested and ready-to-run on supported NVIDIA GPUs on-premises and in the cloud. In addition, it provides pre-trained models, model scripts, and industry solutions that can be easily integrated in existing workflows.

Compiling and deploying Deep Learning frameworks from is time consuming and error-prone. Optimizing AI software requires expertise. Building models requires expertise, time and compute resources.

The NGC catalogue takes care of these challenges with GPU-optimized software and tools that data scientists, developers, IT and users can leverage so they can focus on building their solutions.

## CONTAINER MANAGEMENT - KUBERNETES
LMX Cloud features on-demand Kubernetes provisioning and scheduling for simplified container management. Kubernetes is a platform that automates the deployment and management of containerised applications, including

complicated workloads like AI and machine learning.

## AI AT THE EDGE
Our multi-cloud technology allows organisations to deploy smaller or edge environments that can replicate the central infrastructures functionality or enable remote inference workloads to process data at the edge of the network. LMX Cloud allows organisations to manage a geographically distributed infrastructure from a single dashboard and helps reduce the complexities of running cloud native workloads across multiple cloud environments.

## AI REFERENCE ARCHITECTURE
The LMX Reference Architecture is a suggested best practice for any organisation considering an AI strategy, and can fit any budget. Start small and scale as your organisation and workloads grow. Our platform is built using building blocks that enable you to infinitely scale your compute, storage or acceleration resources depending on requirements. The LMX suite of building block hardware includes hyperconverged servers for the clusters control plane, compute featuring the latest generation AMD EPYC™ processors, NVMe and Object storage blocks and the NVIDIA DGX2™ A100 for acceleration.

## RACK ARCHITECTURE
In the rack architecture diagrams we assume up to 30Kw of power and cooling can be provided. For datacenters not optimized for high density power/cooling we can re-architect this configuration to fit within more modest power budgets per rack.

The suggested configuration in Fig. 7 includes:
- 1x LMX Deployment Node (1U)
- 3x HCI (Control Plane) (2u each)
- 2x DGX A100 (6u each)
- 3x NVMe Storage (1U Each)
- 3x 4U Storage (4U each)



*Fig 6. NVIDIA NGC Application catalogue*

- 2 x Mellanox 100GB Switches
- 1x Supermicro Management
- 1x Supermicro IPMI Management
- Total Rack : 40U (total power 24kw)

## POWER
The solution is designed for high capacity/high density server environments where a 30KW/rack is assumed in terms of power and cooling support. The componets are very high performance and dense in terms of space utilisation so careful consideration is required to ensure appropriate cooling and power is provided in the datacentre.

## NETWORK ARCHITECTURE
The high-performance fabric is used for both storage and inter-process communication of parallel workloads. To ensure that the GPUs are kept active, its essential to deliver high bandwidth low latency access to the data sets. This fabric can also be used to scale applications across multiple nodes (using MPI or NCCL) which also provides GPU Direct / RDMA between GPUs in separate servers. This architecture also supports

NVMe/TCP for accelerated IO to the Lightbits LightOS™ NVMe-over-fabric storage. The network design is resilient so there are no single points of failure in the fabric.

## MANAGEMENT
The management network and OOB (Out-of-band/IPMI) network are used as the control fabric within the architecture. IPMI provides full system power control, full remote console access and full KVM (Keyboard/Video/Mouse) is also available to allow teams remotely manage the infrastructure.

## AI REFERENCE ARCHITECTURE USE CASE
Recently, a government organisation focused on natural disaster detection and prevention has deployed a multi-region, GPU accelerated HPC cluster based on the LMX AI reference architecture, as a national AI cloud for cloud native workloads.

The deployment serves as a 7-Region national scale cloud, which is primarily being used for Disaster detection and prevention. Featuring a
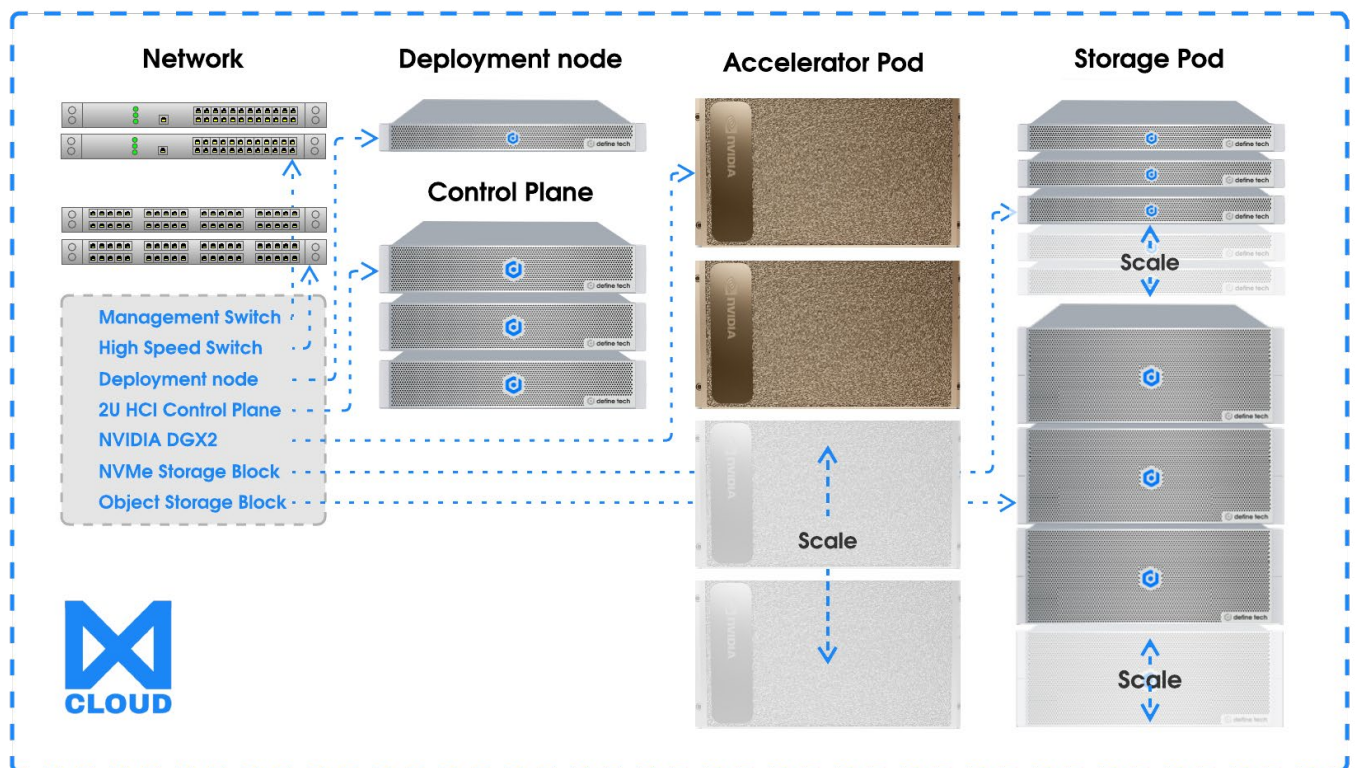


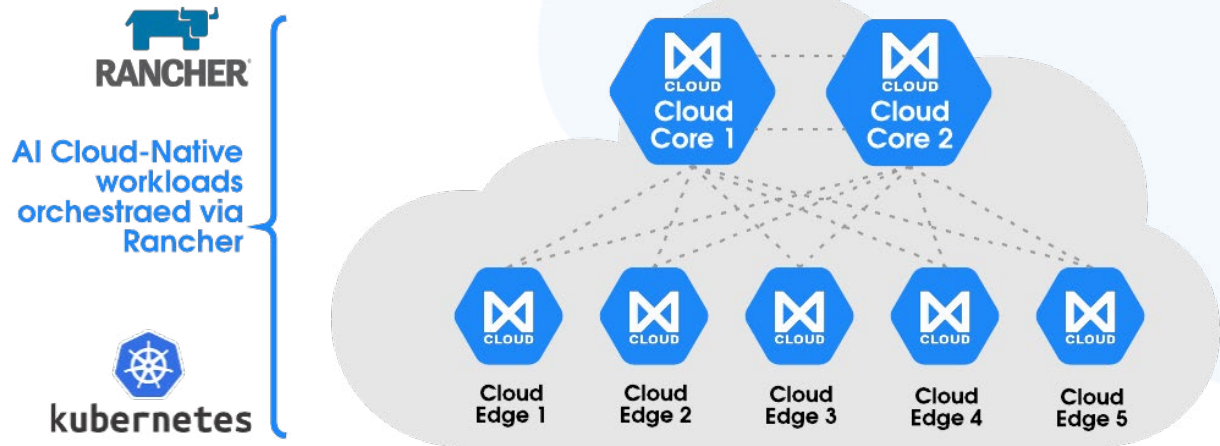*Fig 7. LMX AI Reference Archiecture*

*Fig 8. LMX AI Reference Architecture*

resilient core for AI training (including NVIDIA GPUs as well as NVMe), with all workloads orchestrated via k8s (using Rancher). Each of the 7 cloud regions operates a full LMX cloud environment with secure multi-tenancy. See fig. 8.

## CONCLUSION

Modern workloads are adapting AI tools and techniques to shift through ever growing datasets so as to extract insights in fields such as Life Science, Manufacturing, Academia and Energy, with a common goal of ultimately improving our way of life.

Our AI Reference Architecture is designed to suit any organisation, withing any industry that is looking to deploy an AI strategy -regardless of where they are on their AI journey and what budget restrictions they are under.

Our platform is built with the agility to respond and scale to any workload demands, ranging from compute intensive, complex simulations to the processing of large data sets, so you can keep your IT costs low, your resource utilisation high and the time to insights unbeatable.

*Future-proof your AI strategy with LMX.*

**info@define-technology.com**

**+44(0)203 034 5550**