

CLOUD NATIVE INDORE: TECH TALK

# Build AI Cloud using Kubernetes

Anjul Sahu, CEO, CloudRaft

CLOUDRAFT



**Anjul Sahu**  
CEO, CloudRaft

# About Me

- Founder & CEO, **CloudRaft** - an AI & Cloud Native Consulting
- Organizer, Cloud Native Indore
- More than 16 years in Industry building large scale systems.  
Previously worked for Telco, Banks, Product & Startups
- Passionate about new technology



256 GH200 DGX Cluster  
1 Exaflop, 144 TB GPU  
2023

# In this Presentation

## Overview

- 01 What is AI Cloud?
- 02 Current Trends in AI Infrastructure
- 03 How Cloud Native helps in running AI
- 04 Architecture of AI Cloud
- 05 Cloud Native Projects for AI
- 06 Challenges
- 07 Q&A



An AI Cloud simplifies AI implementation for organizations by integrating it into daily operations. AI Clouds cover the AI lifecycle, from creating features and models to operating, monitoring, and sharing them throughout the organization. Platforms supporting the full AI lifecycle are known as AI platforms, and when available in scalable environments, they are termed AI Clouds.

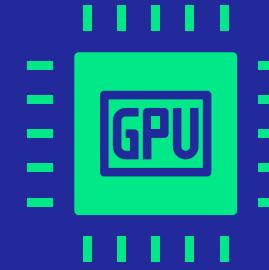
# Features of AI Cloud



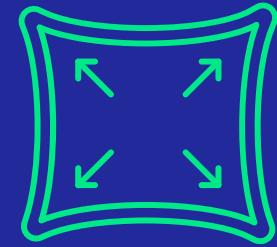
On-prem , Hybrid or Cloud



Support end-to-end  
lifecycle of AI



GPUs & High Performance



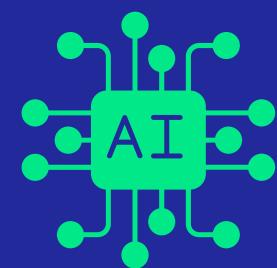
Scalable



Reliable



Billing or Chargeback



AI/ML Frameworks



Full stack: IaaS, PaaS, SaaS



Self-service

# Current Trends in AI Infrastructure

01

**Specialized Cloud**  
Eg: CoreWeave,  
Salad, RunPod,  
Nebius, Lambda labs  
etc

02

**2x Data in every 18 months**  
The demand for data to build better AI/ML models is increasing faster than Moore's Law, doubling every 18 months

03

**Cloud Native and Kubernetes is an accelerator for AI**

04

**GenAI: Bigger Models**  
model size is increasing that means more powerful infrastructure is required

05

**Data Sovereignty Requirements**  
Enterprise data loss risk, AI Safety and new Govt policies to keep data local

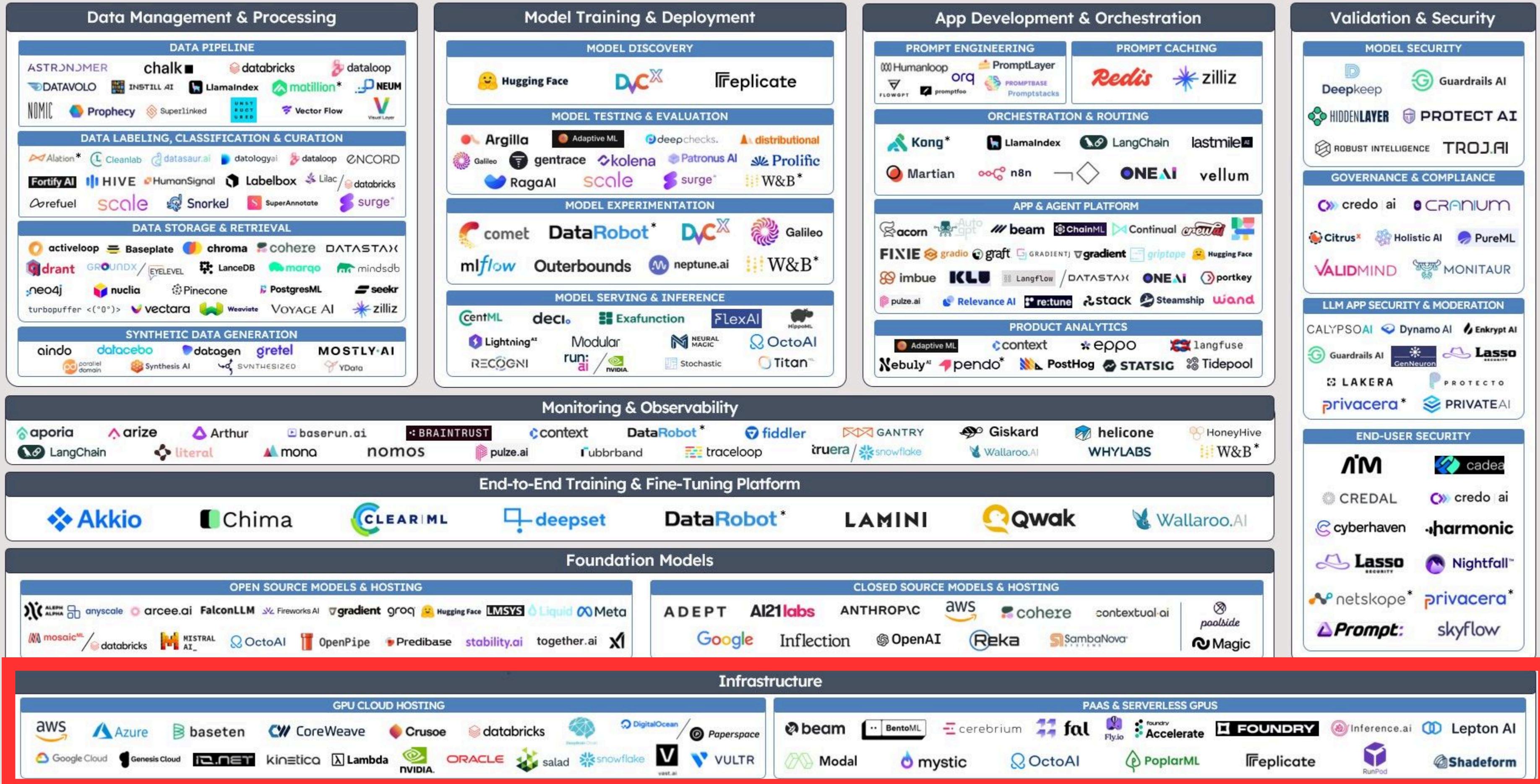
# AI Runs on ~~GPUs~~ Accelerators

- AI = matrix multiplications which is massively parallelizable
- GPUs are great at parallel programming
- CPU < 32 cores/threads, GPUs > 4000 cores/threads
- CPU is 10x slower at least
- Impractical to train or even run any reasonable AI model outside ASICs



# GEN AI APP INFRASTRUCTURE STACK

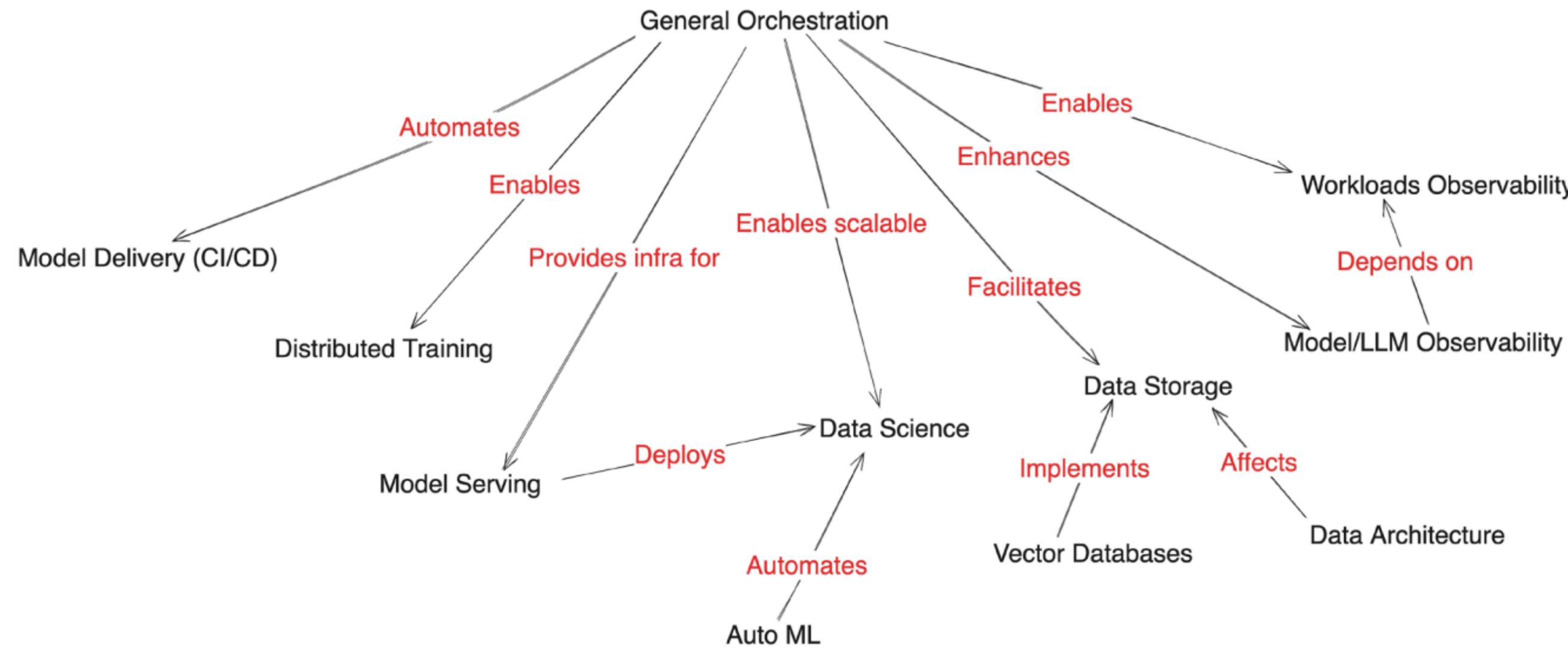
SAPPHIRE  
VENTURES



\* Denotes current or exited Sapphire portfolio company

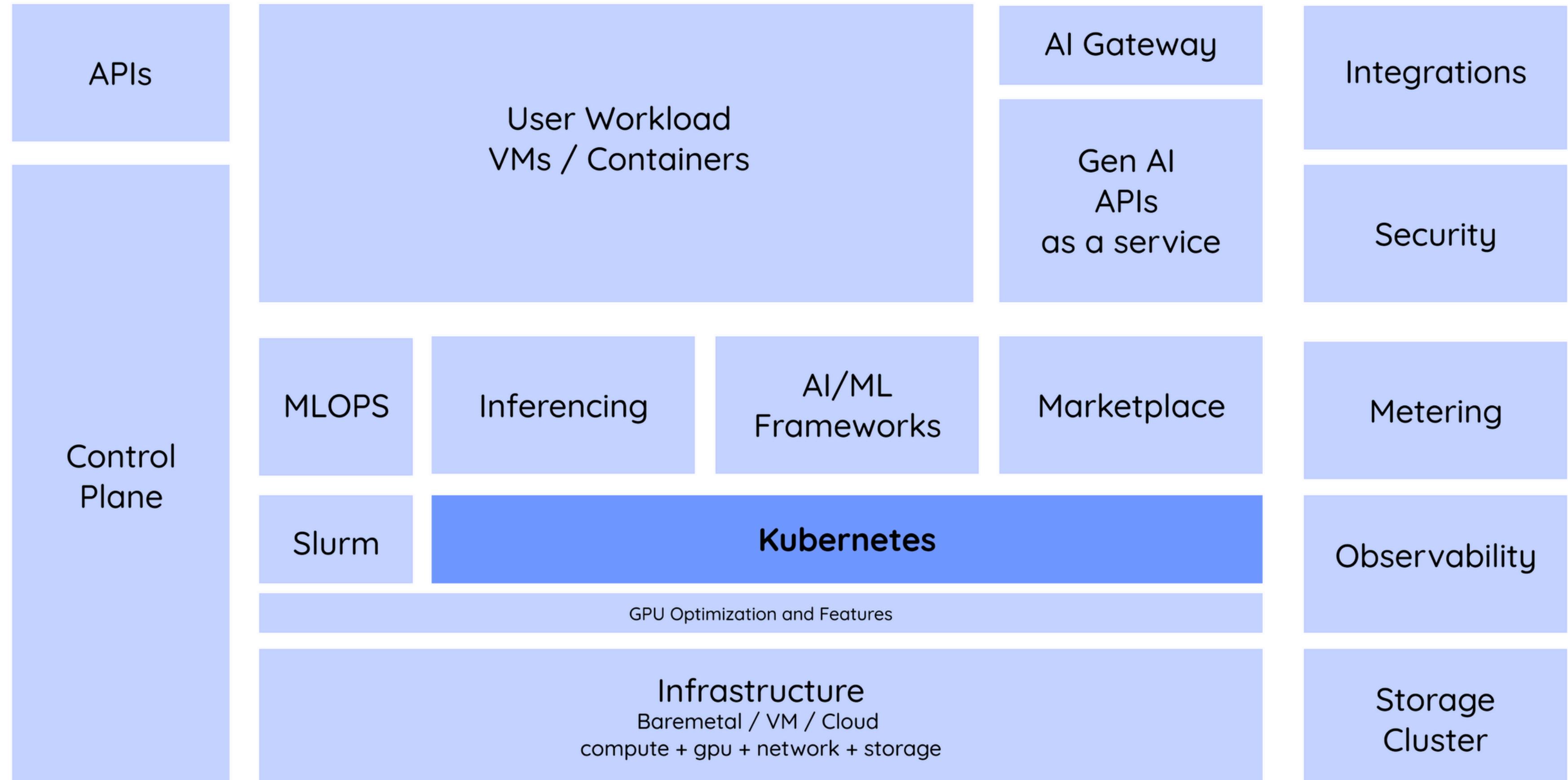
Such companies are a representative sample of portfolio companies in which Sapphire has invested and which the author believes fit the criteria of the market, which do not reflect all investments made by Sapphire. Sapphire has internally assigned categories to all companies described above in the market map. Although Sapphire believes that such categorizations are reasonable, they are subjective in nature and may be categorized differently by other market participants.

# How Cloud Native helps in running AI Workload



"Research teams can now take advantage of the frameworks we've built on top of Kubernetes, which make it easy to launch experiments, scale them by 10x or 50x, and take little effort to manage."

— CHRISTOPHER BERNER, HEAD OF INFRASTRUCTURE FOR OPENAI



# AI Cloud Reference Architecture

# Cloud Native Projects for AI

## Distributed Training

## Model / LLM Observability

## ML Serving

## Vector Databases

## Data Architectures

## Governance and Policy

## General Orchestration

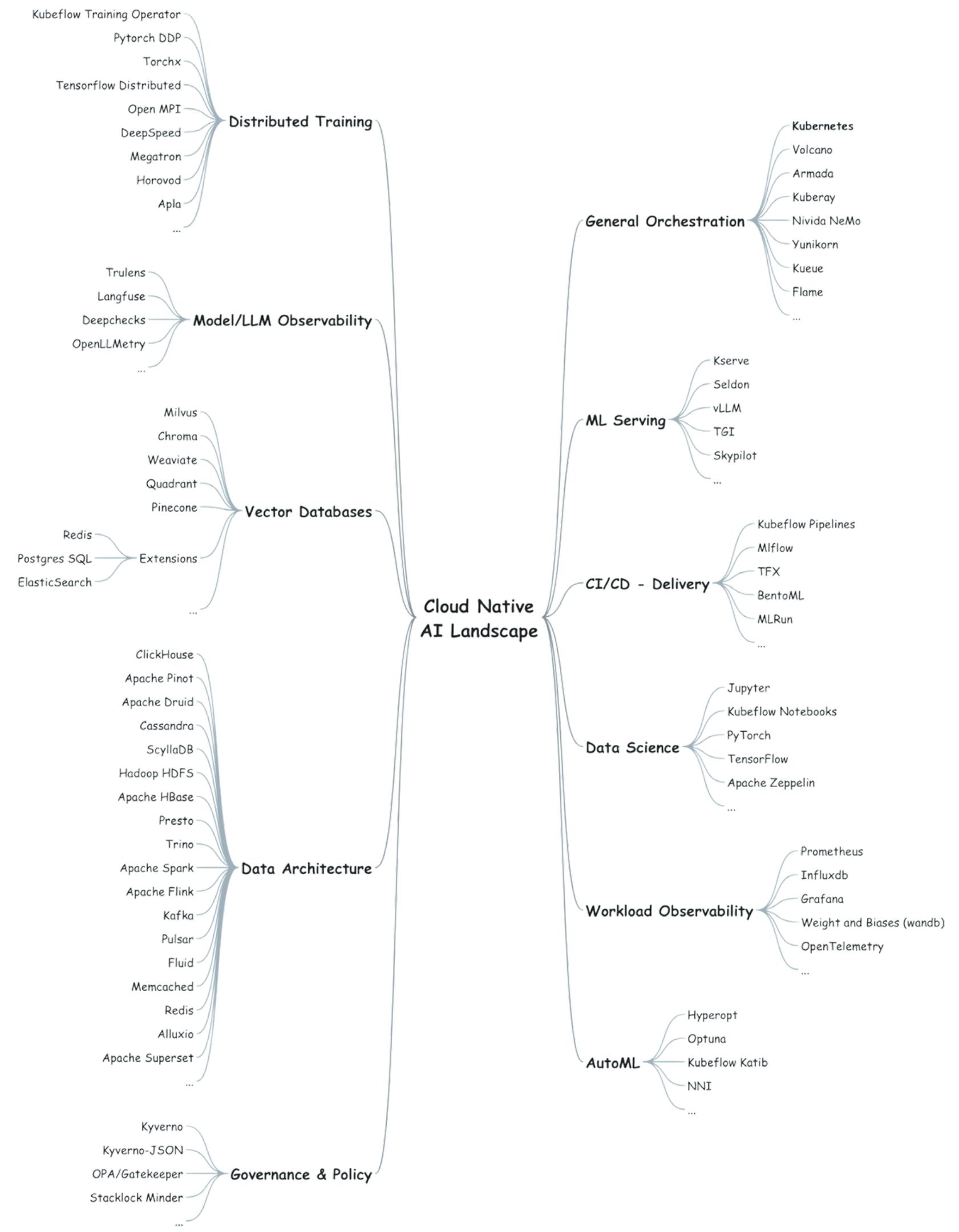
## CI/CD Delivery

## Workload Observability

## AutoML

## Security

*Ecosystem is evolving fast...*



# Challenges in Building AI Cloud

- Building an AI Cloud is a large investment
- GPU supply chain issues
- Skill issues
- High reliability required for long running distributed training jobs
- Unknown security threats and AI Risk in the fast evolving ecosystem
- Sustainability - Each H100 energy consumption is more than avg household
- Some of the hardware limitations becomes bottlenecks such as storage or the network

# Why we need AI Cloud?

- Data Privacy
- AI is making humans more productive
- AGI is possible
- Cost is still less as compared to hyperscalers
- It is a game changer for many enterprises

This talk is based on **our recent work**.  
And it was not possible without the ground breaking  
innovations done by  
**Kubernetes, NVIDIA and CNCF foundation**

See our insights on AI  
[clouddraft.io/blog](http://clouddraft.io/blog)



**kubernetes**



# Q & A

*"Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last, unless we learn how to avoid the risks."*

*-Stephen Hawking, Theoretical Physicist*

CLOUDRAFT