# 25+ Years where we come from...

**2023**
**AGRF**
Flash storage for genomics

**2023**
**Pawsey**
130PB Archive Cluster

**2023**
**PTT Thailand**
HPC Consulting

**2023**
**WSU**
Neuromorphic AI Cluster

**2022**
**Deakin University**
Clustered NVMe software defined storage for AI research

**2022**
**WEHI**
Clustered flash storage for genomics, cryo-em

**2022**
**XENON Cloud**
HPC-as-a-Service for HPC/AI, private & public cloud

**2021**
**Deakin University**
Design & deliver new HPC cluster for shared use across Deakin

**2021**
**Hong Kong Uni HPC**
HPC Cluster with XENON Cluster Stack

**2021**
**Murdoch Children's Research Institute**
Design and implement backup solution

**2021**
**XENON Cluster Stack**
Containerised HPC management solution

**2020**
**CSL Limited**
Services for HPC cluster implementation

**2020**
**Monash University**
HPC cluster & storage configured as shared cloud node for Australia Research Data Commons

**2020**
**Todd Energy**
HPC cluster for oil & gas exploration

**2019**
**Pawsey Supercomputing Centre**
New GPU cluster

**2018**
**Harrison.ai**
AI solution for IVF, radiology, healthcare

**2017**
**NCI**
Supercomputer for scientific research & technology innovation

**2017**
**WEHI**
Private cloud for next generation cancer, disease & medical research

**2016**
**Garvan Institute of Medical Research**
NVMe solution for medical research

**2015**
**Thales Defence**
High fidelity solution for Australian Army Tiger helicopter simulators

**2014**
**RCC**
FlashLite HPC cluster

**2012**
**Fujitsu**
Infiniband for Raijin Supercomputer

**2009**
**CSIRO**
Australia's first GPU Cluster

**2007**
**Victorian Partnership for Advance Computing**
HPC Cluster for Advance Computing

**2005**
**Thales Defence**
Image Generators for ASLAV simulators

**2000**
**Animal Logic**
Render Farm for "Matrix" movie trilogy

**1998**
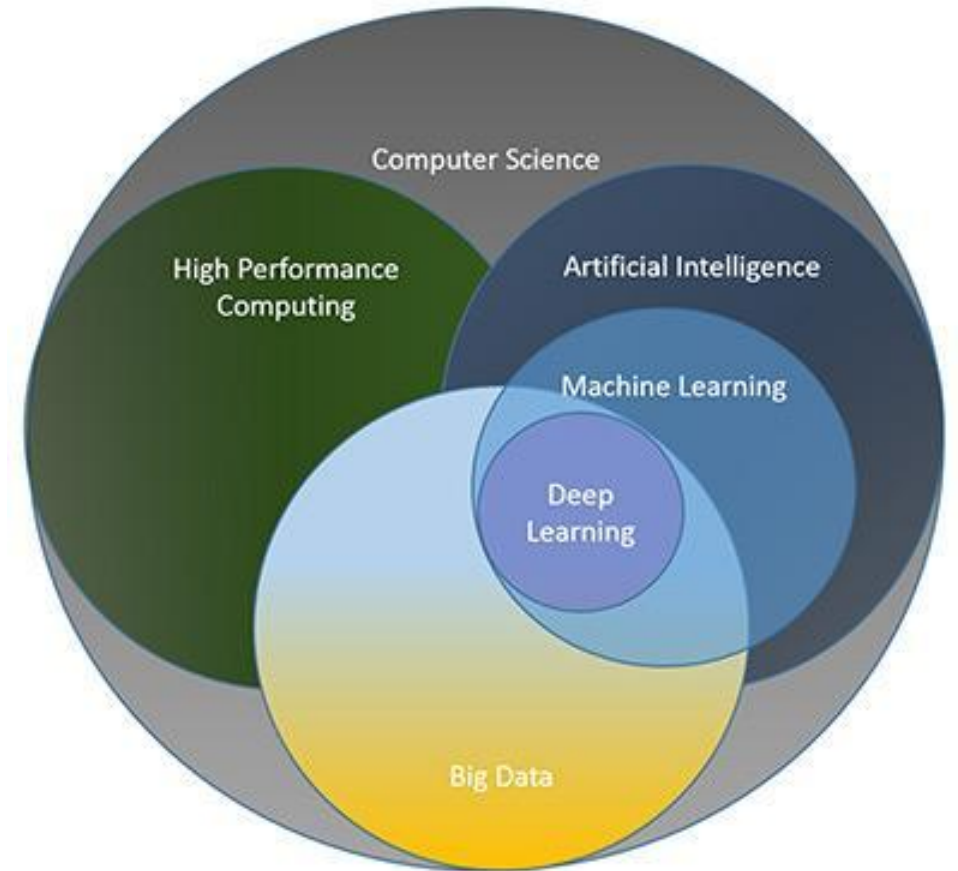**Telstra Digital Video Network**
Video Server equipment

**2020**
**2010**
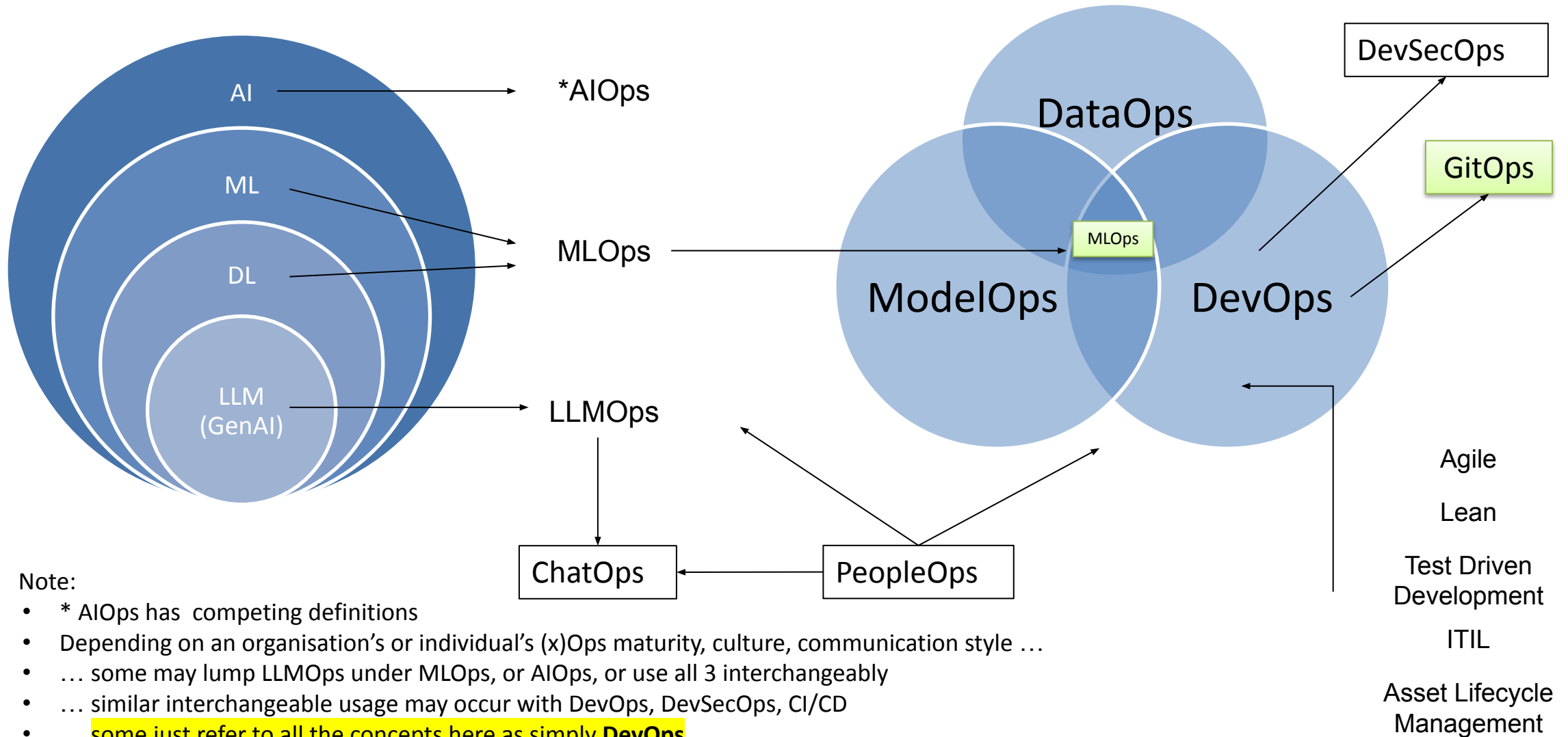**2000**

**Read the Case Studies**

# Introduction

- HPC and AI workload convergence

- Demystify and learn from container workflows, DevOps, MLOps, LLMOps, and Platform Engineering

- Can Kubernetes support:

  - Queuing and scheduling typical of HPC?

  - "Lift-and-shift" HPC jobs <u>onto</u> Kubernetes?

  - Ignoring microservices

- How do researchers access storage from Kubernetes?

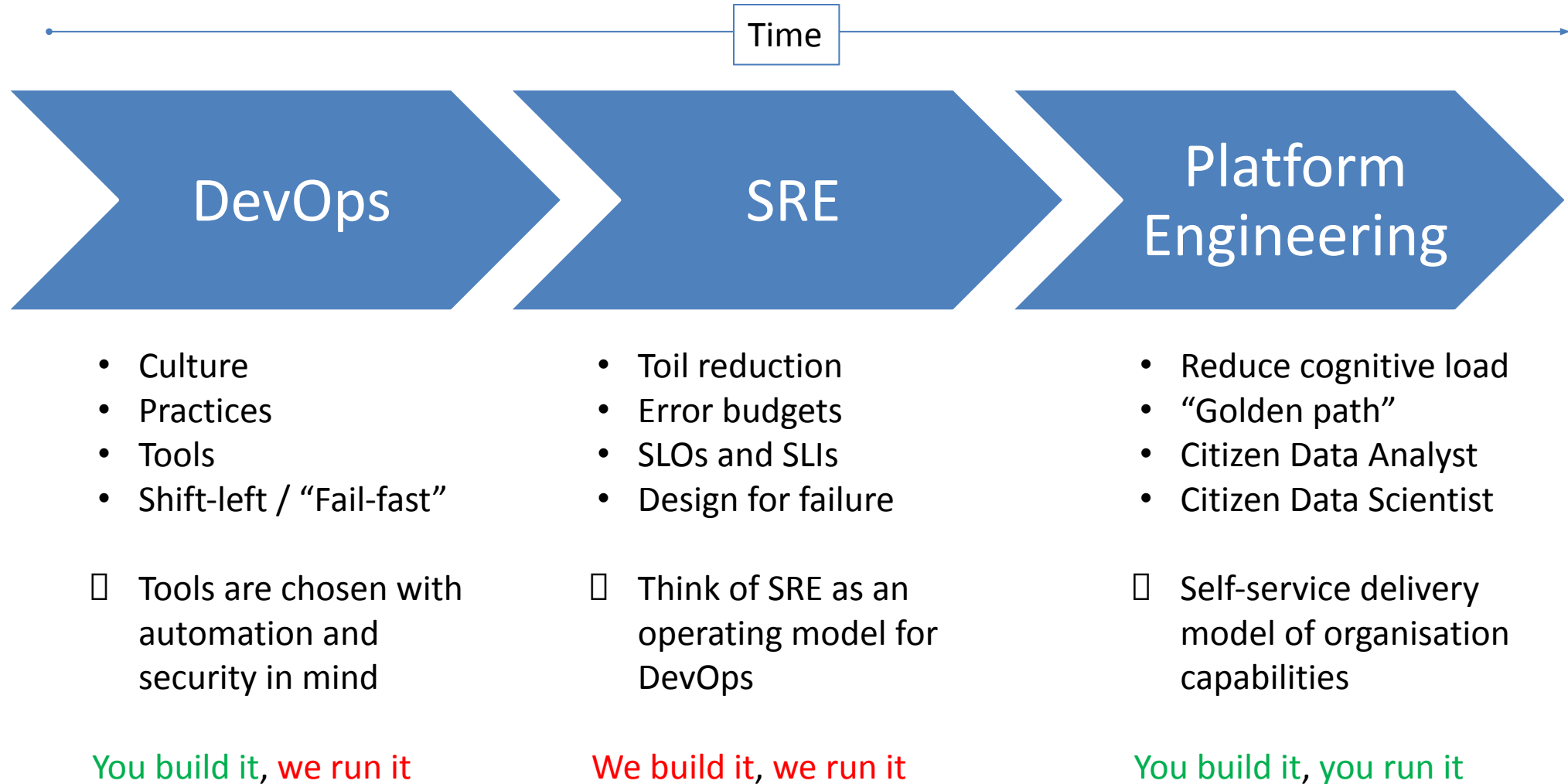- How does Kubernetes fit into a hybrid architecture encompassing cloud, on-prem, and IoT (lab devices)?

- How to adopt?



Source (Feb 2025):
https://www.hpcadvisorycouncil.com/subgroups_hpc_ai.php

# xOps – Buzzword Explosion



Note:
- * AIOps has competing definitions
- Depending on an organisation's or individual's (x)Ops maturity, culture, communication style …
- … some may lump LLMOps under MLOps, or AIOps, or use all 3 interchangeably
- … similar interchangeable usage may occur with DevOps, DevSecOps, CI/CD
- … some just refer to all the concepts here as simply **DevOps**

# Evolution of Automation

**Time** →

**DevOps** → **SRE** → **Platform Engineering**

| DevOps | SRE | Platform Engineering |
|---|---|---|
| • Culture<br>• Practices<br>• Tools<br>• Shift-left / "Fail-fast" | • Toil reduction<br>• Error budgets<br>• SLOs and SLIs<br>• Design for failure | • Reduce cognitive load<br>• "Golden path"<br>• Citizen Data Analyst<br>• Citizen Data Scientist |
| ☐ Tools are chosen with automation and security in mind | ☐ Think of SRE as an operating model for DevOps | ☐ Self-service delivery model of organisation capabilities |
| You build it, we run it | We build it, we run it | You build it, you run it |

XENON
High Performance Computing

# DevOps vs SRE vs Platform Engineering

Gartner Hype Cycle for Software Engineering, November 2023



1. DevSecOps
2. Site Reliability Engineering (SRE)
3. Platform Engineering
   a. Internal Developer Portal (IDP)
4. GitOps

**What's the difference?**

**Choose 1 only?**

**Choose all 4?**

➡️ **They're complimentary, use all 4.**

Source (Feb 2025): https://www.gartner.com/

# Why Platform Engineering is Important

Note: The size of each box shows the <u>respective size and importance</u> for a functional MLOps platform



Code is a <u>tiny</u> part of an AI/ML system

Source (Feb 2025): <u>Hidden Technical Debt in Machine Learning Systems</u>

# Principles and Standards

## 1. The Twelve-Factor App

**I. Codebase**
One codebase tracked in revision control, many deploys

**II. Dependencies**
Explicitly declare and isolate dependencies

**III. Config**
Store config in the environment

**IV. Backing services**
Treat backing services as attached resources

**V. Build, release, run**
Strictly separate build and run stages

**VI. Processes**
Execute the app as one or more stateless processes

**VII. Port binding**
Export services via port binding

**VIII. Concurrency**
Scale out via the process model

**IX. Disposability**
Maximize robustness with fast startup and graceful shutdown

**X. Dev/prod parity**
Keep development, staging, and production as similar as possible

**XI. Logs**
Treat logs as event streams

**XII. Admin processes**
Run admin/management tasks as one-off processes

## 2. Semantic Versioning 2.0.0 (SEMVER)

MAJOR.MINOR.PATCH, e.g. v3.2.8

1. MAJOR version when you make incompatible API changes
2. MINOR version when you add functionality in a backward compatible manner
3. PATCH version when you make backward compatible bug fixes

Additional labels for pre-release and build metadata are available as extensions to the MAJOR.MINOR.PATCH format.

## 3. Open Container Initiative (OCI)

- **Image** Specification: Defines the structure and format of container images (e.g. Dockerfile).

- **Runtime** Specification: Specifies how a container should be executed and managed.

- **Distribution** Specification: Outlines standards for distributing container images.

➡️ Think in terms of OCI compliant images for "lift and shift" enablement

**XENON**
High Performance Computing

# AI Platform

Weights & Biases

NVIDIA NIM OPERATOR

Kubeflow

PyTorch

Apache Spark

**Platform Customers**

- Accessible by Customer Users (some)
- Accessible by Customer Admins (all)
- Administered by Customer Admins (all)
- Customised by Customer Admin (all)

## Platform

jupyter

YUNIKORN

KEYCLOAK

- Administered by Platform Team (all)
- Customised by Platform Team (only)

CERT MANAGER

NVIDIA GPU OPERATOR

### K8S Operators

- CNCF, Non-CNCF, Apache projects
- Administered by Platform Team (all, only)
- Customised by Platform Team (all, only)

## Kubernetes

containerd

- Administered by Platform Team (all, only)
- Customised by Platform Team (all, only)

### Hardware

- Administered by Platform Team (some/all/only)

# PaaS (I)

PaaS solutions can be used as a shortcut to building your platform

**Rafay**



PaaS Reference Architecture for GPU Clouds

**Run:AI**



Source (Feb 2025): <u>Rafay Platform - GPU PaaS Reference Architecture for Nvidia Cloud Partners & Enterprises</u>

Source (Feb 2025): <u>https://www.run.ai/</u>

# PaaS (II)

PaaS solutions can be used as a shortcut to building your platform

**KubeFlow**



**NVIDIA AI Enterprise**



Source (Feb 2025): https://www.kubeflow.org/docs/started/architecture/

Source (Feb 2025): https://docs.nvidia.com/ai-enterprise/overview/latest/platform-overview.html

# Next-gen Data Centre Management



Soft Multi-tenancy

Hard Multi-tenancy

# App Catalogue

# Batch Schedulers



A light-weight universal resource scheduler for container orchestrator systems.

- App-aware scheduling
- Hierarchical queues
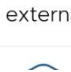- Gang scheduling
- Job ordering & queuing
- Resource fairness
- Resource reservation
- Preemption
- Max application enforcement



Kueue



Slinky = Slurm on Kuberntes

- Offering flexibility and ease of use for both HPC and cloud-native users
- Run and manage Slurm clusters on Kubernetes
- Manages the scaling of Slurm nodes within Kubernetes
- Job allocation/accounting/dependencies
- Fair-share, and priority scheduling
- "Lift-and-shift" potential for unified infrastructure

# Batch Workloads



YuniKorn

+

Kubeflow
Training Operator

- MPI
- Pytorch
- Tensorflow
- Ray Clusters / Jobs / Services
- Spark
- Flink
- PaddlePaddle
- XGBoost
- JAX

# Cluster Mesh – Example Scenarios



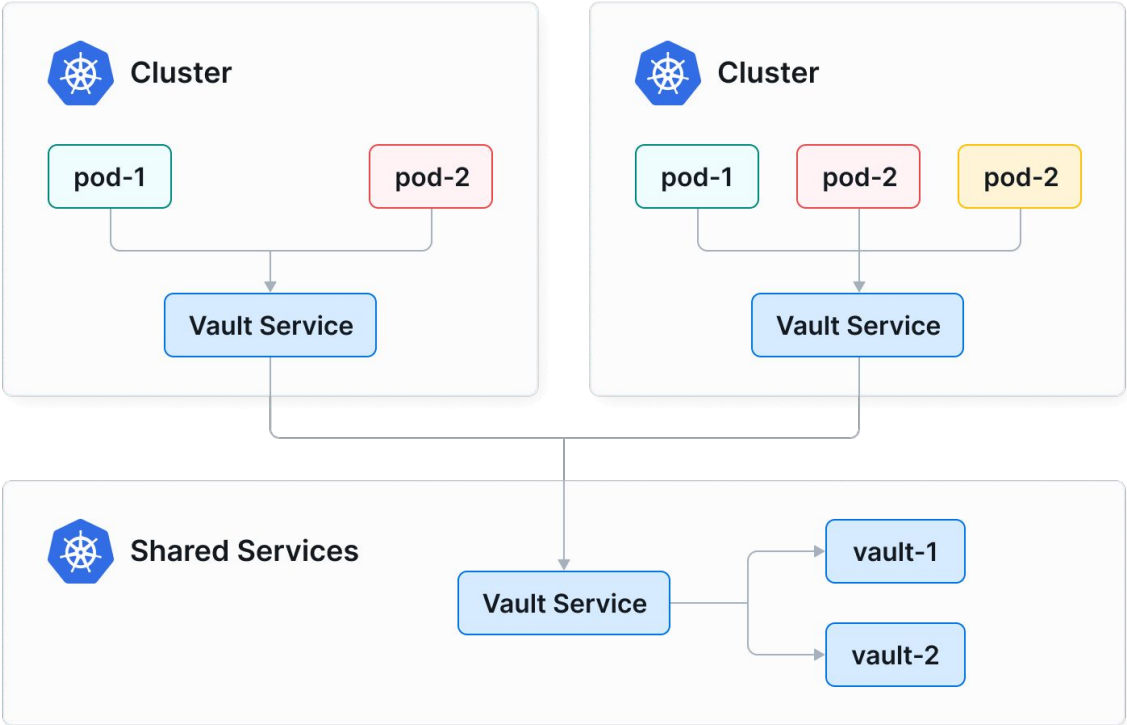…"an open source, cloud native solution for providing, securing, and observing network connectivity between workloads, fueled by the revolutionary Kernel technology eBPF."

## High Availability  and Fault Tolerance



## Shared Services Across Clusters

# Cluster Mesh – Uniform Network Policy Enforcement

# IoT Edge Gateways and Devices

## KubeEdge



Provides fundamental infrastructure support for network, app deployment and metadata synchronization between cloud and edge.

 Good for edge gateways

## K3S



Kubernetes distribution built for IoT & Edge computing

 Good for lightweight vClusters
 Good for devices

Note: Cilium works on both distributions for cluster mesh scenarios

# Storage

**CSI Drivers**

- NFS
- SMB
- Local file (on-node)
- Cloud hyperscalers

**CSI Drivers for native integration with Vendor storage**



DELLTechnologies

HAMMERSPACE

WEKA

ddn

VAST

# Simple Adoption Idea – 1. POC

## Objective

- Prove the technology, i.e. batch jobs as containers on Kubernetes

**Step 1**

**NVIDIA.**
GPU OPERATOR

**Kubeflow**
Training Operator

✔ Driver management
✔ GPU in containers

✔ MPI Jobs
✔ Pytorch Jobs

**Step 2 (Optional)**

**NVIDIA.**
NETWORK OPERATOR

✔ GPU Direct Storage
✔ RDMA

**Tips**:

- Use a single node Kubernetes "cluster" (no VM) to simplify and focus on the objective
- Focus on a <u>single</u> Kubeflow Training Operator workload
- Test the Nvidia Network Operator once you have the job running

# Simple Adoption Idea – 2. MVP

**Objective**
- Prove scheduling workflows, understand additional job types



Training Operator

GPU OPERATOR

NETWORK OPERATOR

✔ Batch Scheduling
✔ Gang scheduling
✔ Fair use
✔ Queues

✔ MPI Jobs
✔ Pytorch Jobs
✔ Tensorflow Jobs
✔ Spark Jobs

✔ Driver management
✔ GPU in containers

✔ GPU Direct Storage
✔ RDMA

# Production considerations

- HA/DR
- Tighter user access controls
- Secure UIs
- Certificates for securing traffic
- Multi-node
- Platform services

# Other Tips

1. Leaky Tap Strategy
2. Blueprints and Deployment stamps
    1. Blueprints tend to refer to software stacks (Rafay and Nvidia use this term)
    2. Deployment stamps refer to Infrastructure as Code (IaC) – Ansible and Pulumi
    3. Sometimes the lines are blurred between the two in cloud-native approaches
    4. Both are essentially templates that rely on configuration to stand-up resources (use config as code!)
3. Scaffold repos
    1. 1x functional repo per use case
    2. Standards embedded in repo design
    3. Configuration locations pre-defined and required to trigger DevOps processes
    4. Clone from to execute and learn, and modify to start new projects in a new repo
    5. Good to teach new users
    6. Fast way to learn for those new to Platform Engineering (even if you're a senior)
4. vClusters
    1. Give you access to blue-green deployment patterns for whole clusters
5. GitOps
    1. Tools provide access to blue-green deployment patterns for cluster resources
6. Reference architectures
    1. Azure MLOps v2, AWS, Nvidia, Xenon AI Sandpit

# Connect with Xenon

# Thank You

**XENON Systems Pty Ltd**
10 Westall Road, Springvale, Victoria 3171, Australia

**www.xenon.com.au**

**P**   +61 3 9549 1111
**F**   +61 3 9549 1199
**E**   info@xenon.com.au

*A  member of the XENON Technology Group*
*www.xtg.com.au*