

Herding **LLaMAs on** **Kubernetes**

**For Sovereignty, Operations
and Integration**

Max Körbächer

Founder & Cloud Native Advisor @ Liquid Reply

With a focus on Platform Engineering, Internal Developer Platforms & Cloud Native Engineering

- CNCF TAG Environmental Sustainability Initiator & Co-Chair
- CNCF Ambassador
- LF Europe Advisory Board
- Contributed 3y to the Kubernetes release team

 [maxkoerbaecher](#)

 [mkoerbi](#)

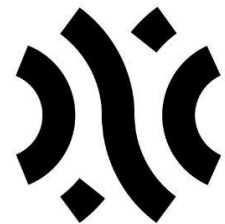


TAG ENVIRONMENTAL
SUSTAINABILITY

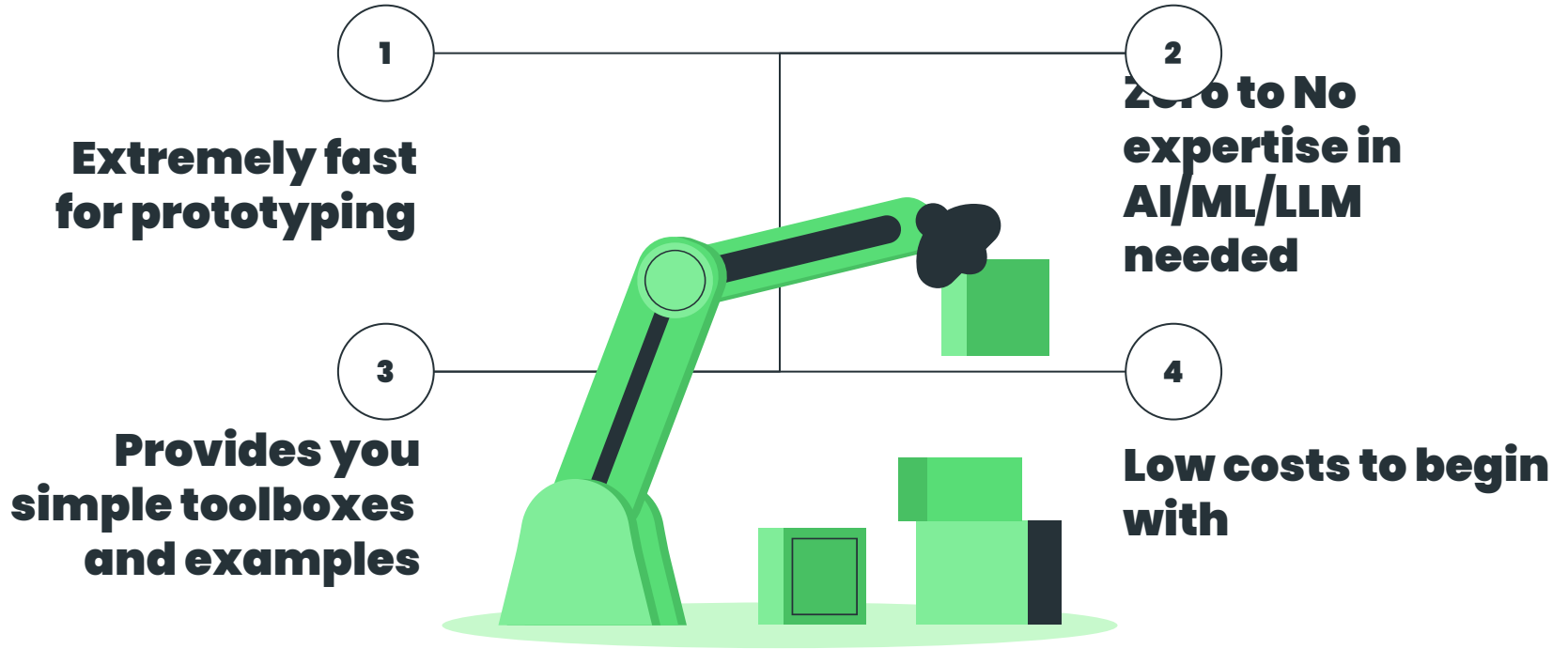
DIY or just buy?



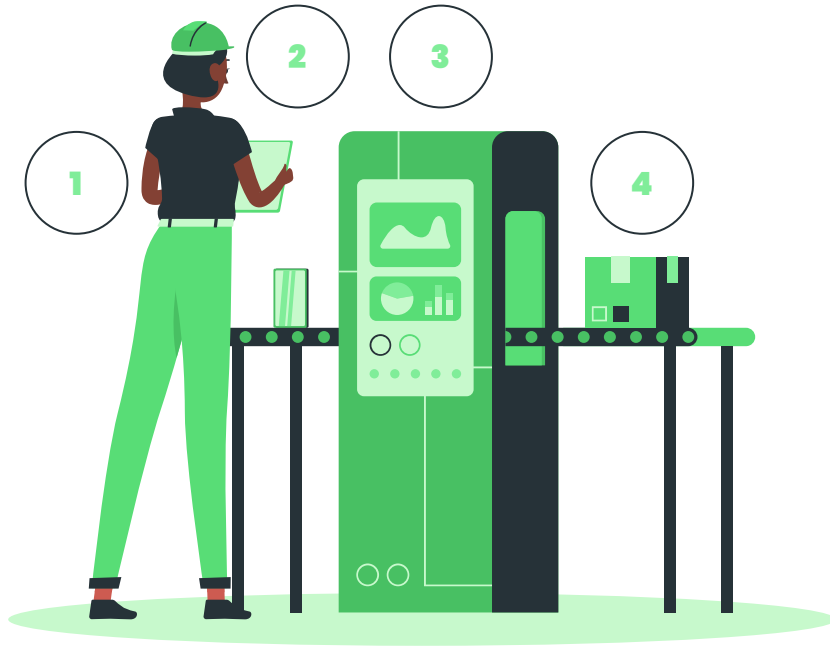

LLAMA 2



Buy-in ChatGPT & co



Self-hosting LLMs – the good
































































- 1 **Avoid lock-in, enable specialized use cases**
- 2 **Greater customization possibilities**
Enhanced Security, IP, Privacy, Compliance
- 3 **Flatter cost curve on the long-run**
- 4

Self-hosting LLMs - challenges

Model Size	Large sizes, issues in loading them fast, poor deployment speed, worse agility
Memory Requirements	Extreme memory requirements, good memory management is crucial
Compute Resources	Requires a lot in production/many end users, high operational costs
Latency	High inference latency due to their complexity
Scalability	Scaling GPUs or replace workload fast isn't trivial

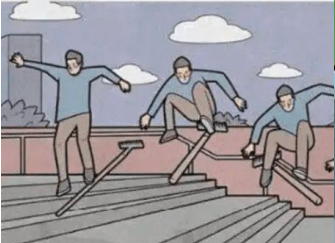
Serving LLMs on K8s

CNAI

Data Architecture 🔍	CI/CD - Delivery 🔍
                    	
General Orchestration 🔍	Workload Observability 🔍
 CNCF GRADUATED	  
Distributed Training 🔍	Governance, Policy & Security 🔍
       	   
Data Science 🔍	AutoML 🔍
     	  
Model/LLM Observability 🔍	Vector Databases 🔍
     	  
ML Serving 🔍	
    	



vm virtual machine



kubernetes

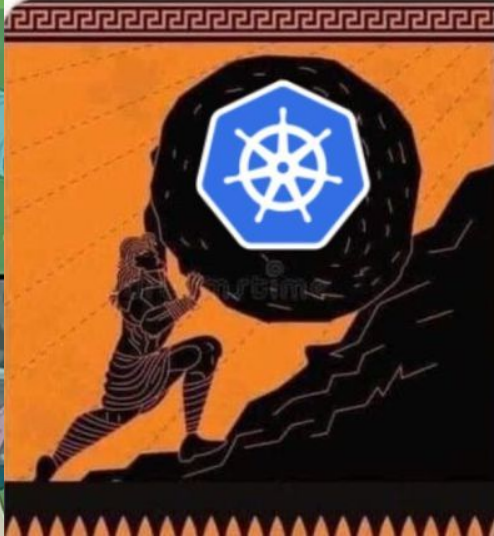
Expectation



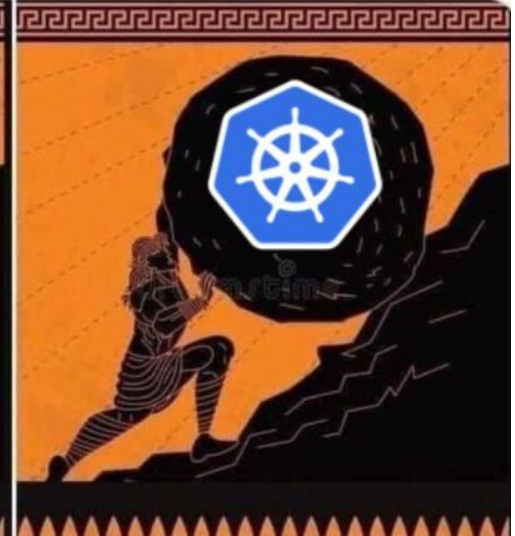
Reality



How it started:



How it's going:



K8s benefits running AI workload

- **Portability** - allows to (more or less) run anywhere
- **Scalability** - adjusts the workload based on its demand or adjust the workload to suit the resources
- **Resource utilization** - Optimize utilization of resources
- **Ecosystem** - Feature rich and active ecosystem for tools, best practices and community

true BUT boring

Kubernetes is a platform to build platforms

Infrastructure
Abstraction &
Automation

Developer Platform
(IDP)

Edge Workload
Handler (IoT)

AI/ML/LLM - the "new"
kid on the block



K8s – Supporting Capabilities for LLMs

Observability

Track performance & health

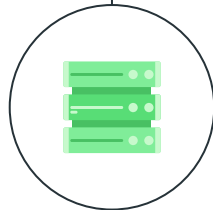


CI/CD – “GitOps”

Adopt best practices for automation

Data Management

Utilize persistent & dynamic volumes



Security

Build on a rich ecosystem for protection and separation/isolation

Observability

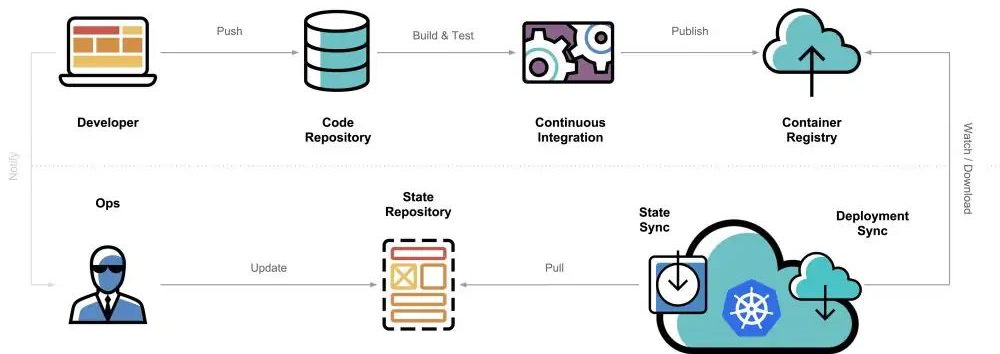
The K8s default provides a full transparency from infrastructure, through network, to the application layer.



CI/CD - "GitOps"

CI/CD and GitOps provided the foundation for MLOps and LLMOps.

CI / CD with GitOps

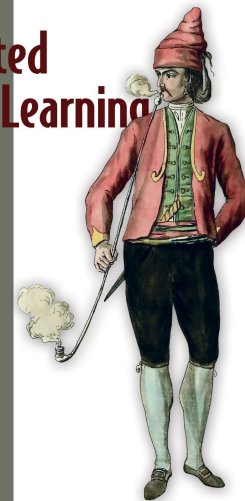


Icons: CC BY 3.0

Distributed Machine Learning Patterns

Yuan Tang

WANNING



Data Management



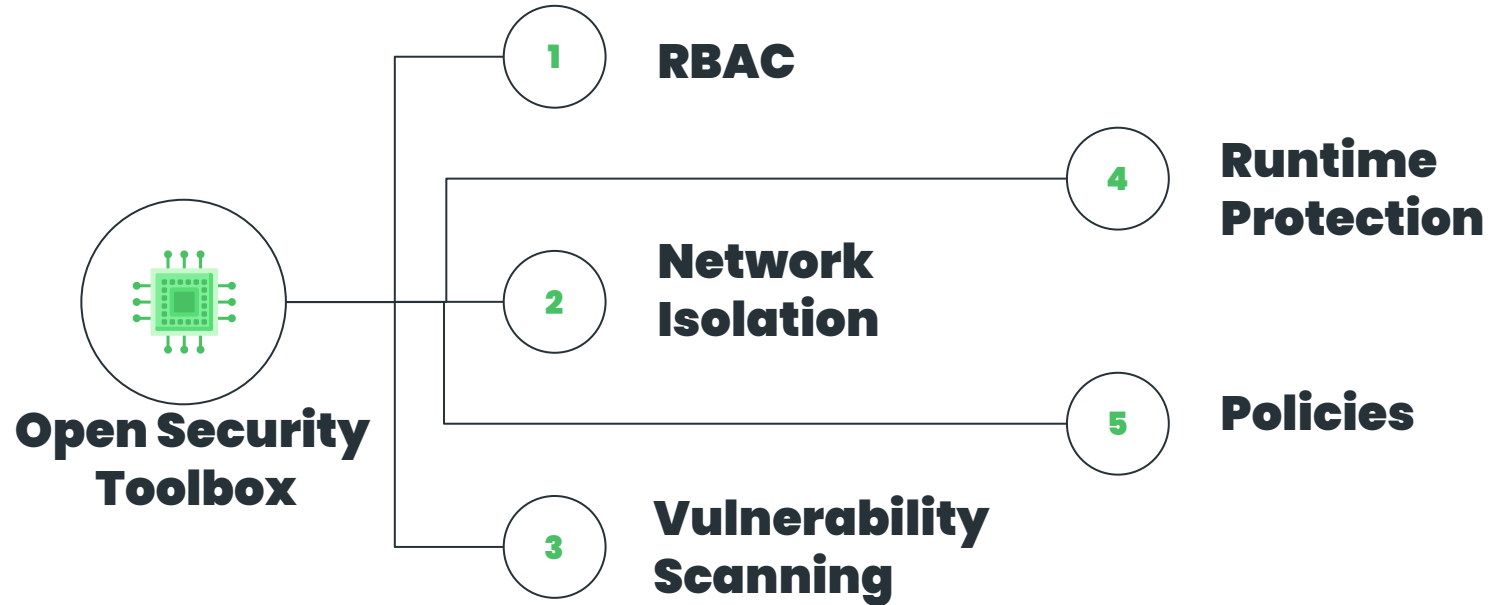
No matter where, no matter what, K8s provides a unified approach of integrating with the storage.

For the end user it is always just a claim, how it is done is out of the sight.



Re-claiming storage, backup of cluster and application state with the data and event-based triggers to adjust to certain metrics improve reliability drastically.

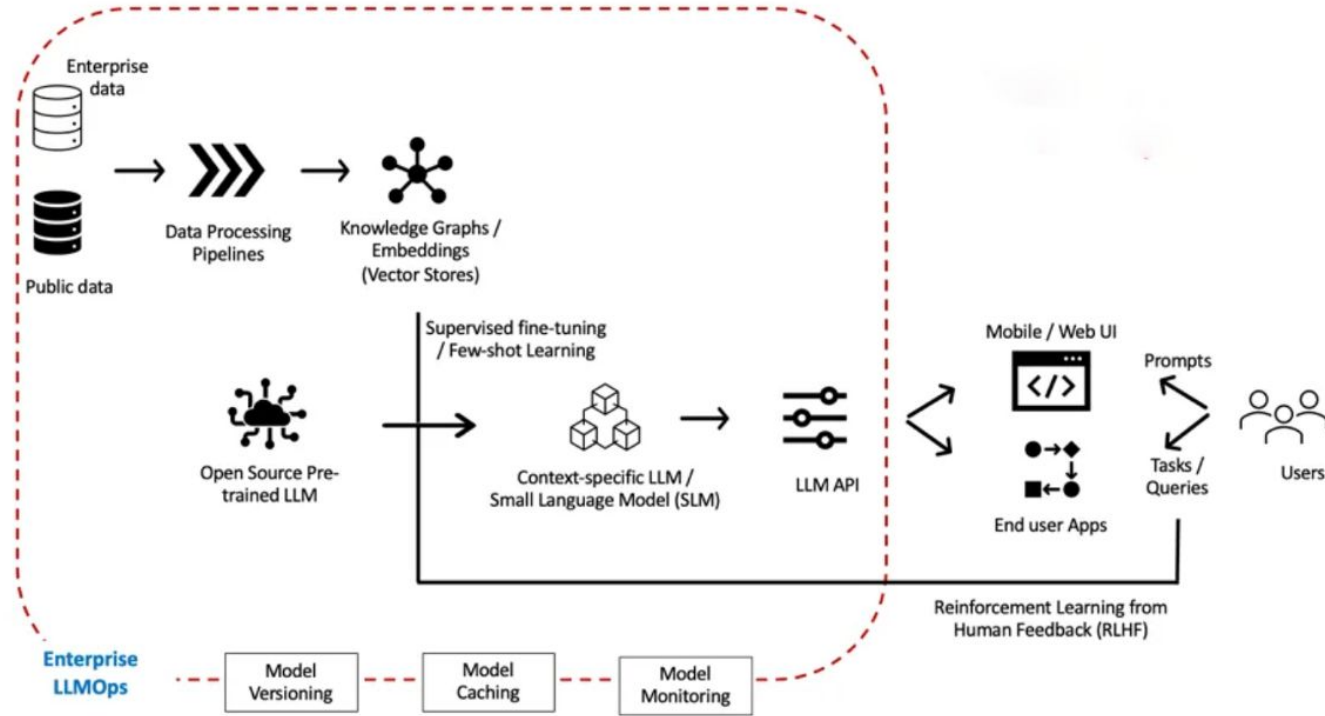
Security



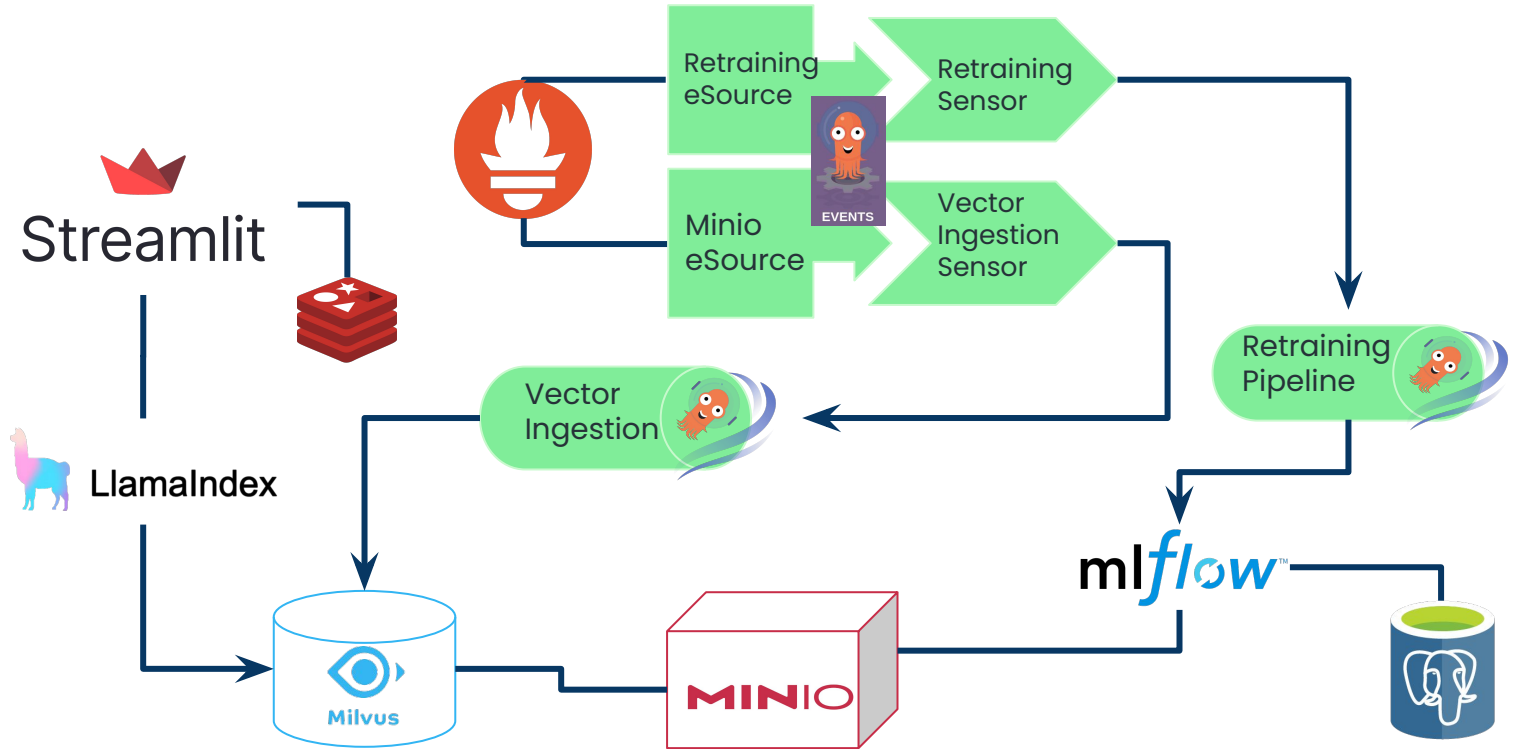


Let's DIY

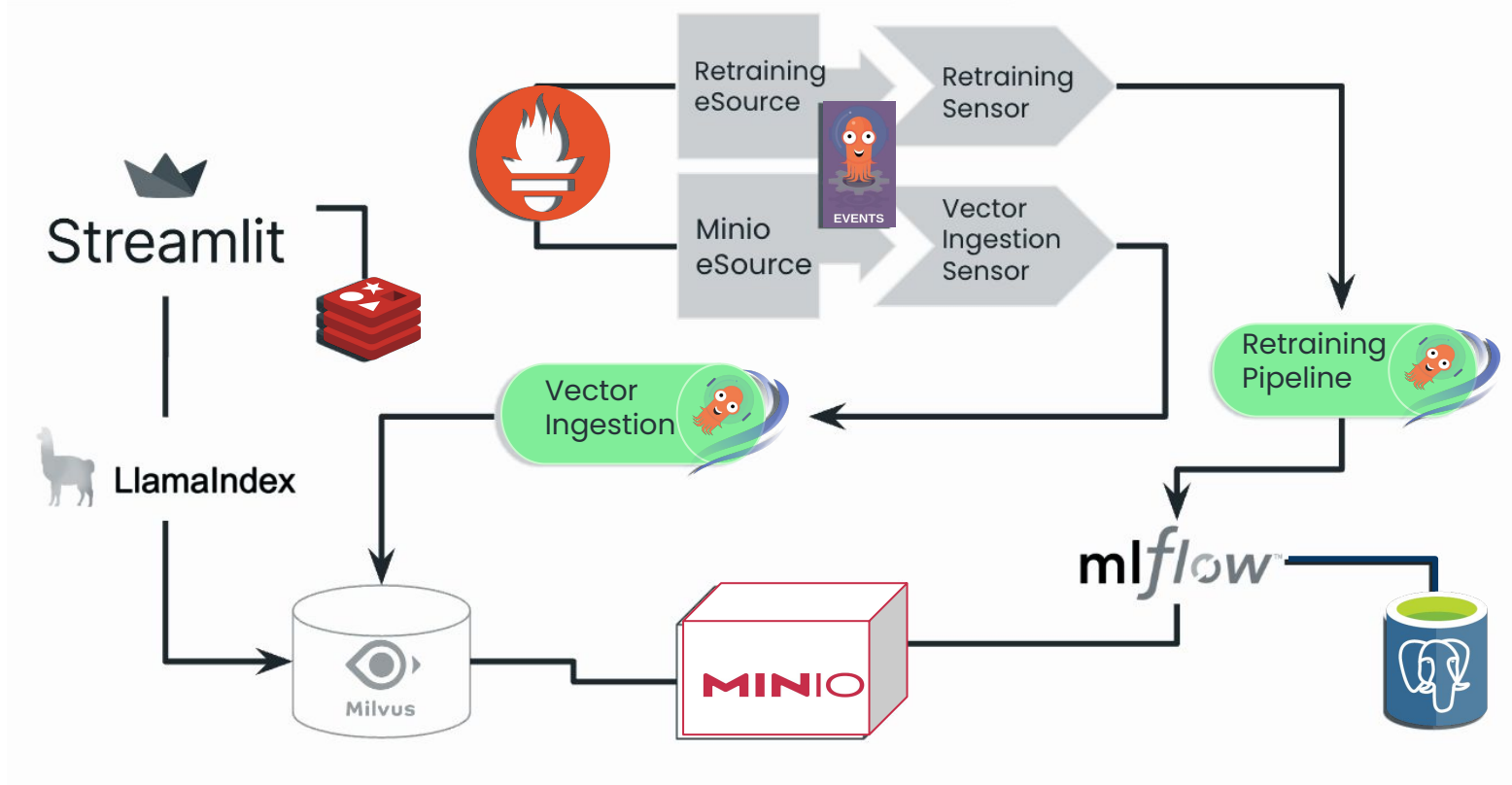
LLMOps Workflow



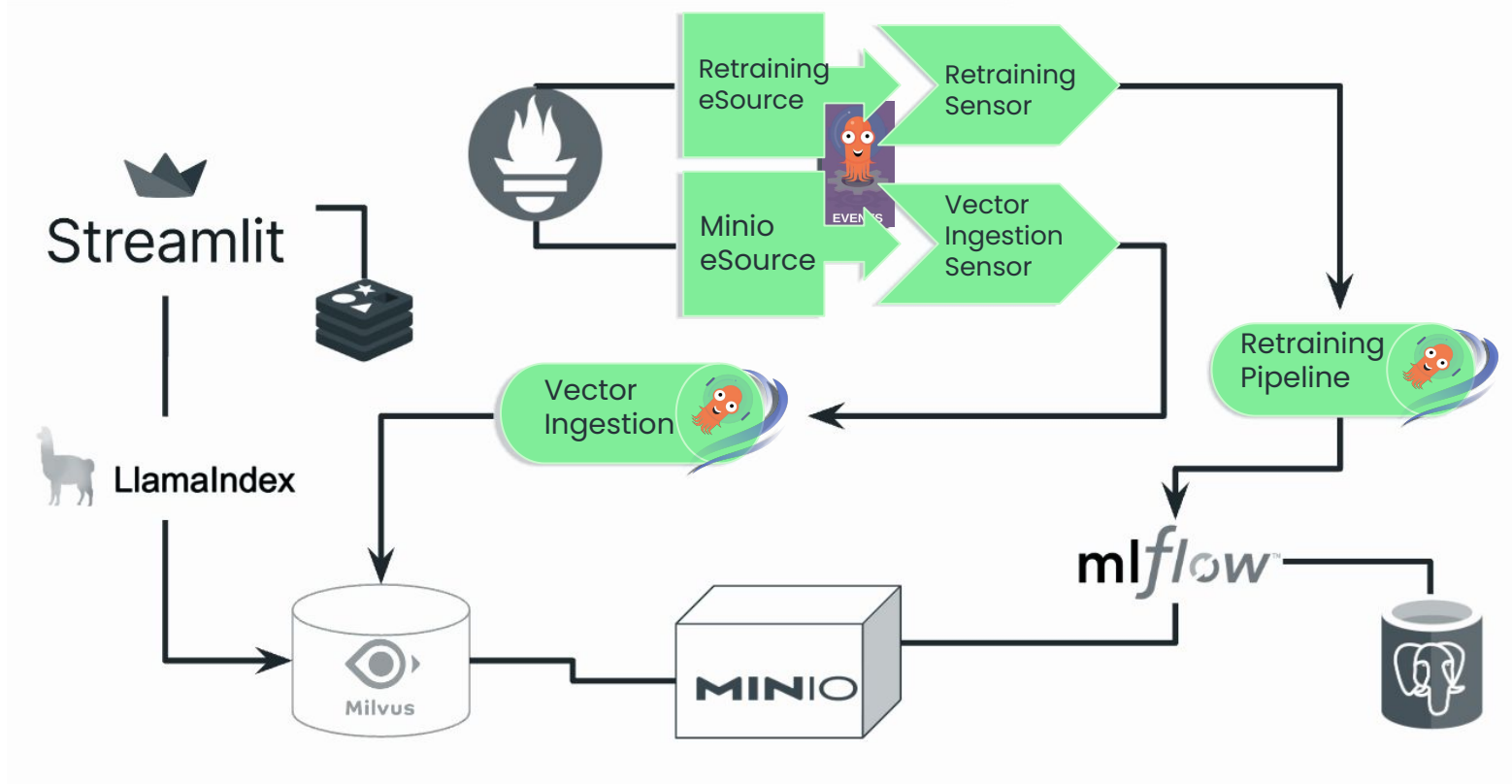
LLMOps example Architecture



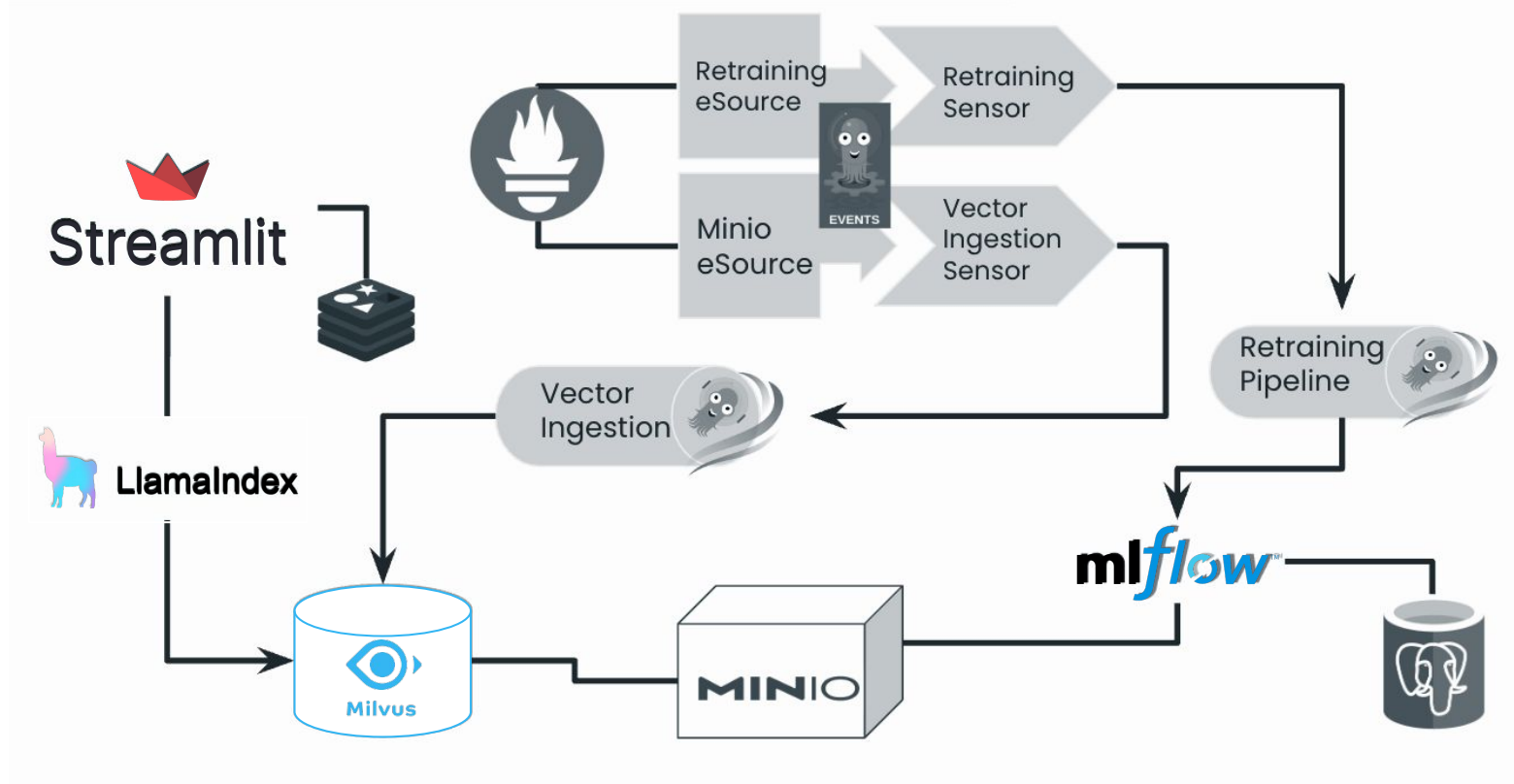
LLMOps example Architecture



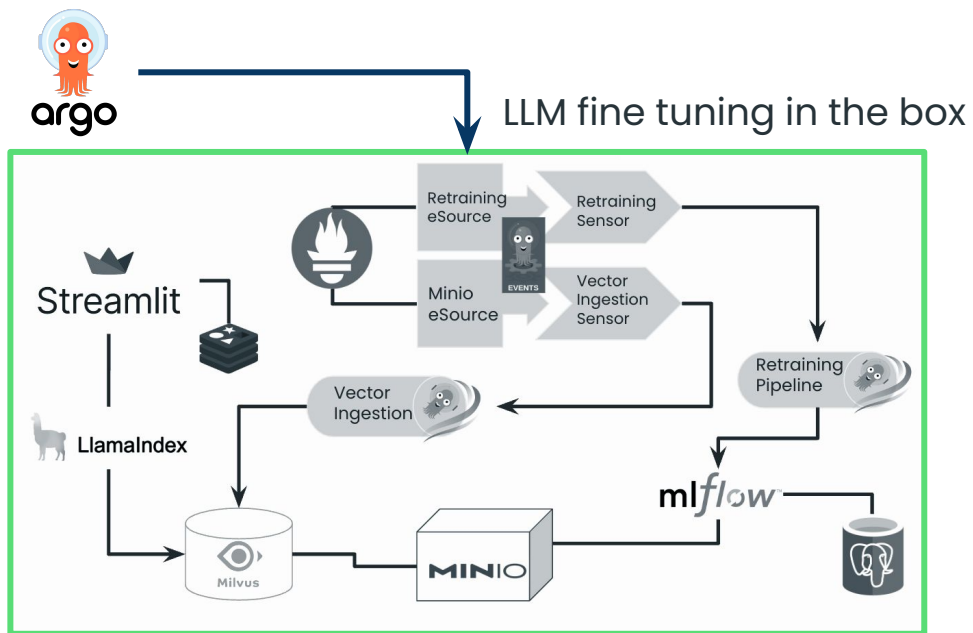
LLMOps example Architecture



LLMOps example Architecture



Kubernetes Empowerment



Kubernetes



Observability



GPU

Operator

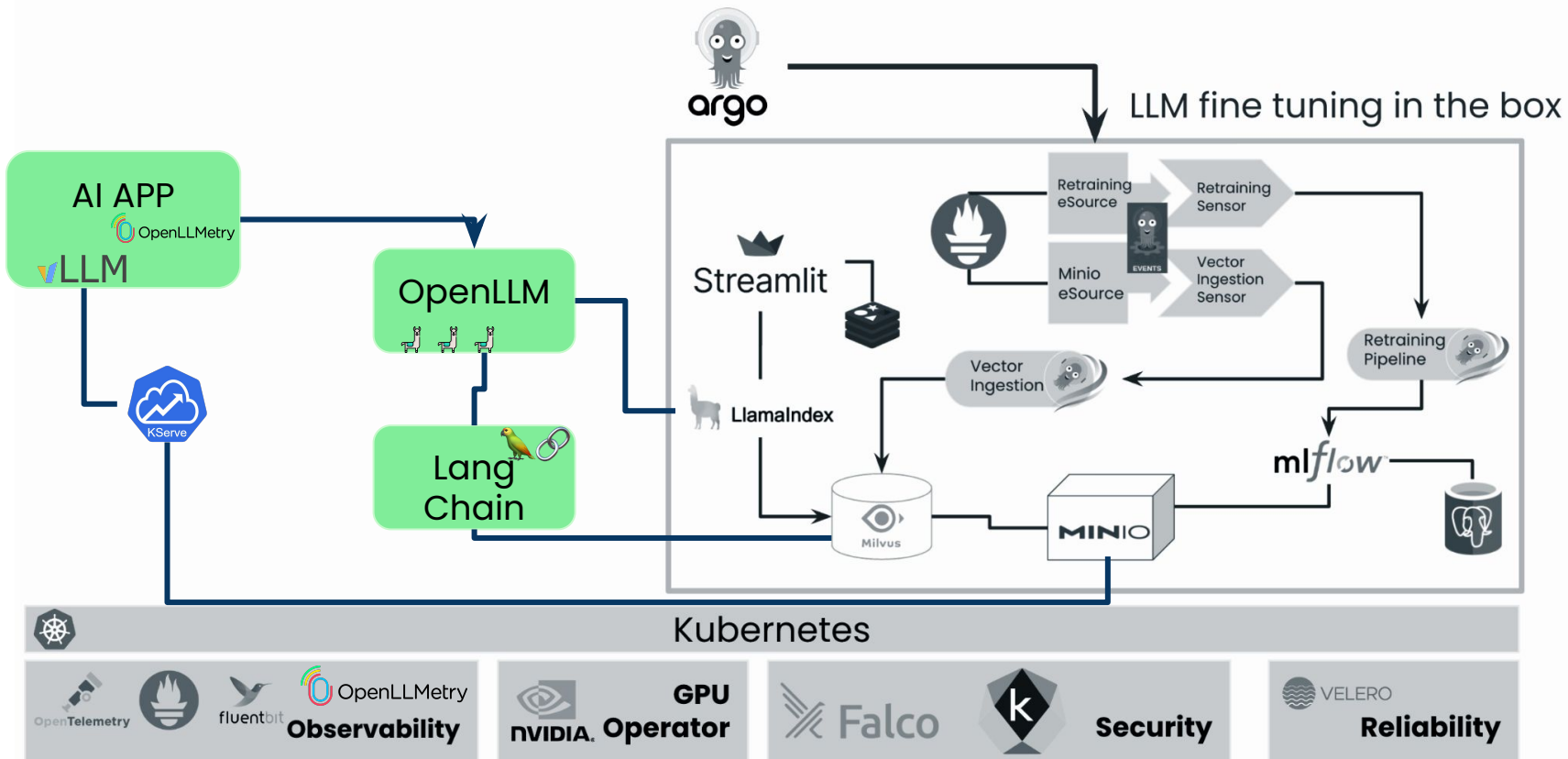


Security

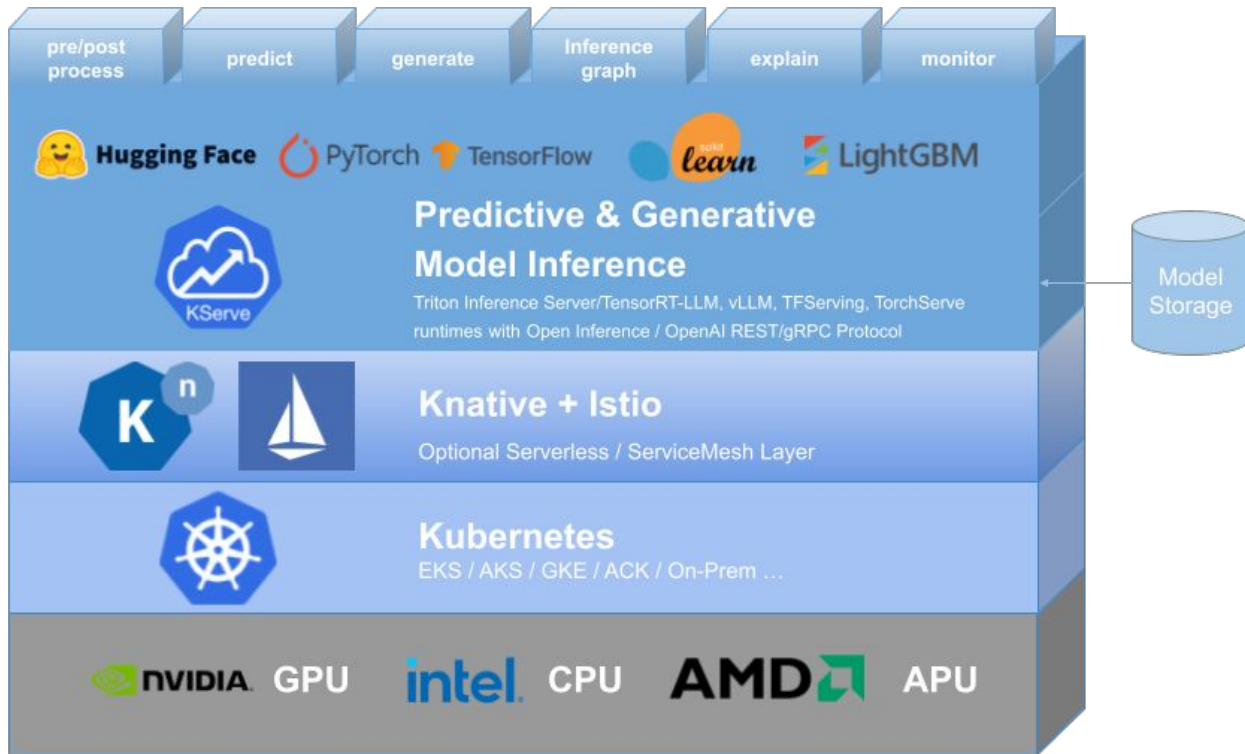


Reliability

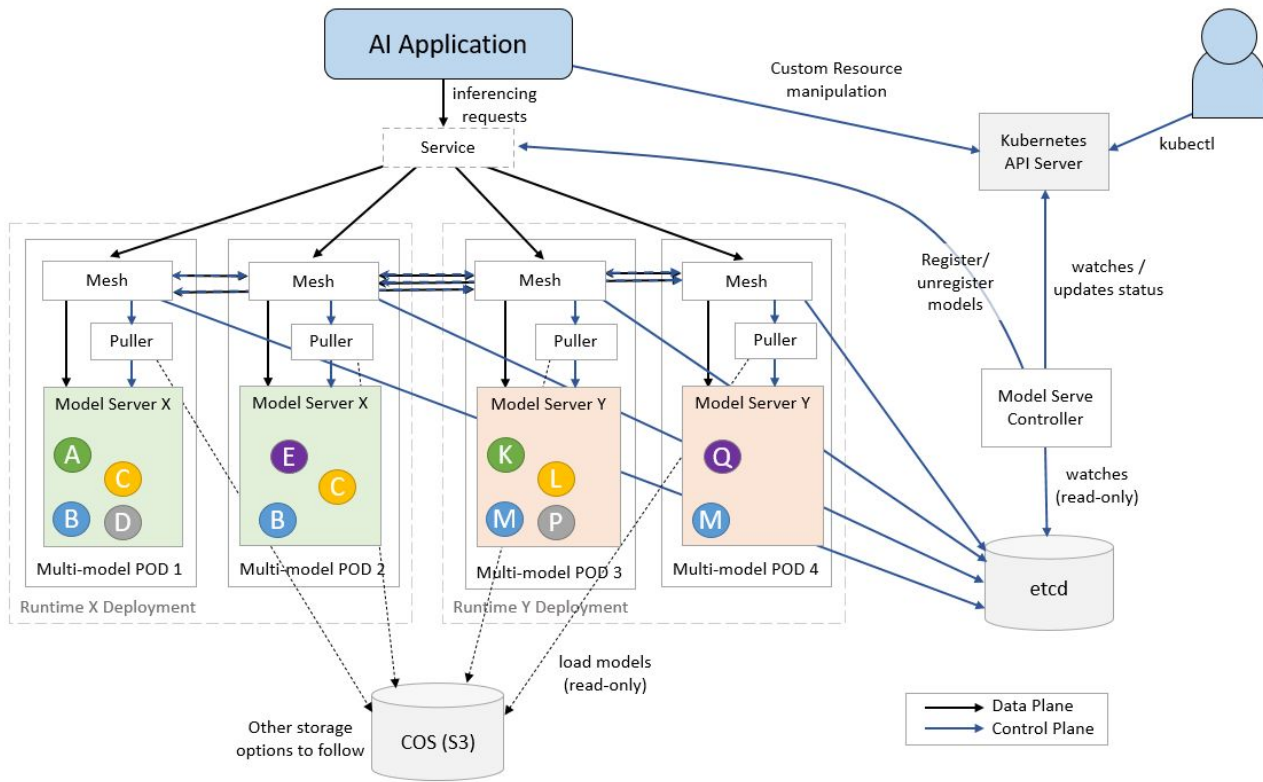
Using LLMs on K8s



10k feet view kserve



10k feet view kserve

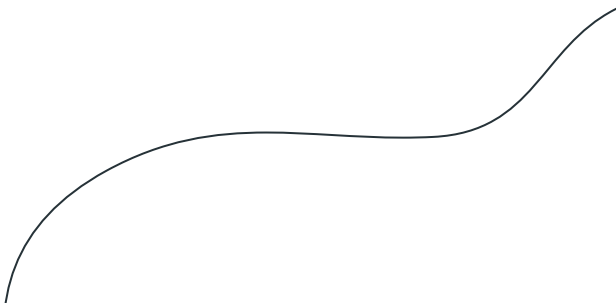




An

Outlook

**We (will) have
various
specialized LLMs to
serve**



That's where K8s plays strong

Flexible provisioning & integrations

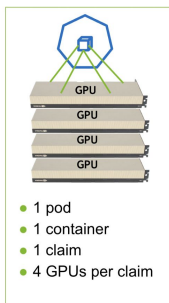
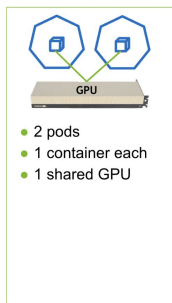
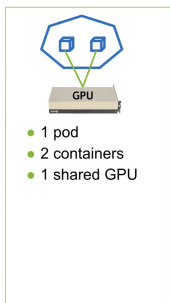
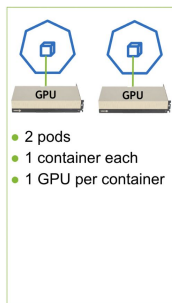
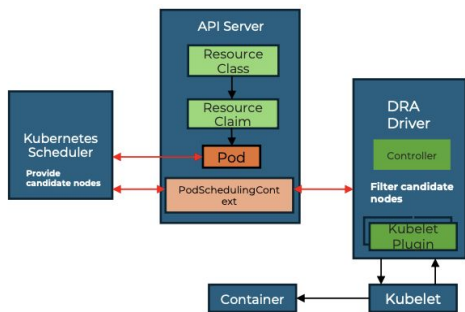
K8s support a very dynamic scaling and shifting of workload based on metrics.

Adjust resources allocation - DRA

Dynamic Resource Allocation to improve GPU and memory utilization. -> better for the business

Dynamic Resource Allocation

Dynamic Resource Allocation (DRA) allows users to create MIG slices



Slice 0	Slice 1	Slice 2	Slice 3	Slice 4	Slice 5	Slice 6
7g.40gb						
4g.20gb						
3g.20gb		✗	3g.20gb			
2g.10gb		2g.10gb		2g.10gb		
1g.10gb	✗	1g.10gb	✗	1g.10gb	✗	1g.10gb
1g.5gb	1g.5gb	1g.5gb	1g.5gb	1g.5gb	1g.5gb	1g.5gb
2g.10gb		1g.5gb	1g.5gb	3g.20gb		
4g.20gb				2g.10gb		1g.10gb
		1g.5gb				1g.5gb

Example A100-40GB MIG layouts

Wasm & Inferencing

Problems

- Efficiency
- Language Dependencies
- Platform Dependencies



Power

- Minimum Resource Footprint
- Dev flexibility
- Secure
- Close to data

Combining DRA, kserve & WASM

Runs "anywhere"
any scale



Extreme resource
efficient

Very cost efficient

Highly dynamical
and adjustable



**Thank
you!**