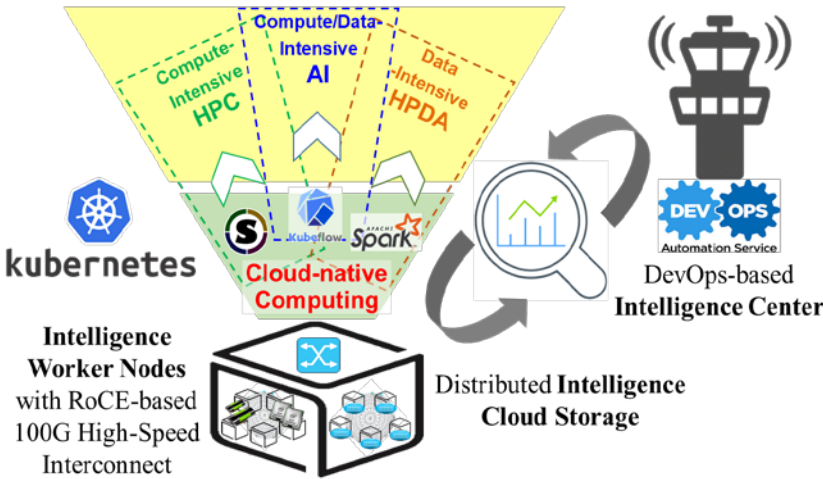


Cloud-native SmartX Intelligence Cluster for AI-inspired HPC/HPDA Workloads

NetCS Lab. (Networked Computing Systems Laboratory), GIST
Jungsu Han, GeumSeong Yoon and JongWon Kim
{jshan, gsyoon, jongwon}@nm.gist.ac.kr

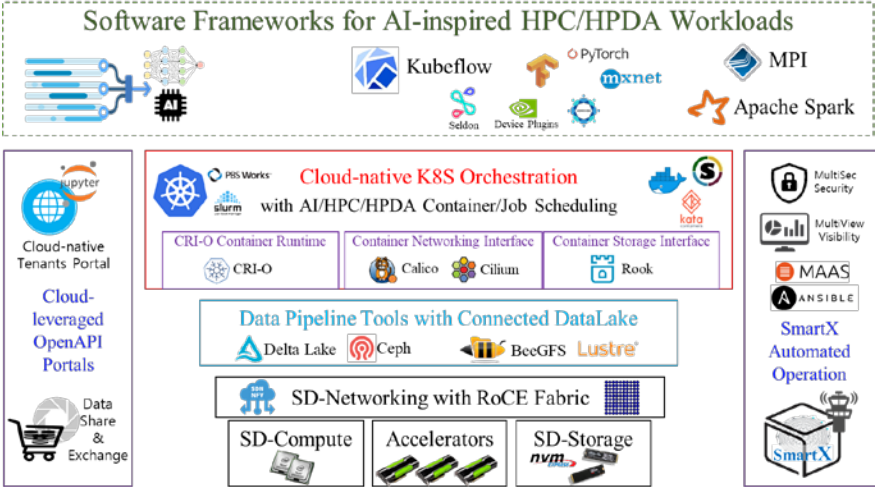


Motivation



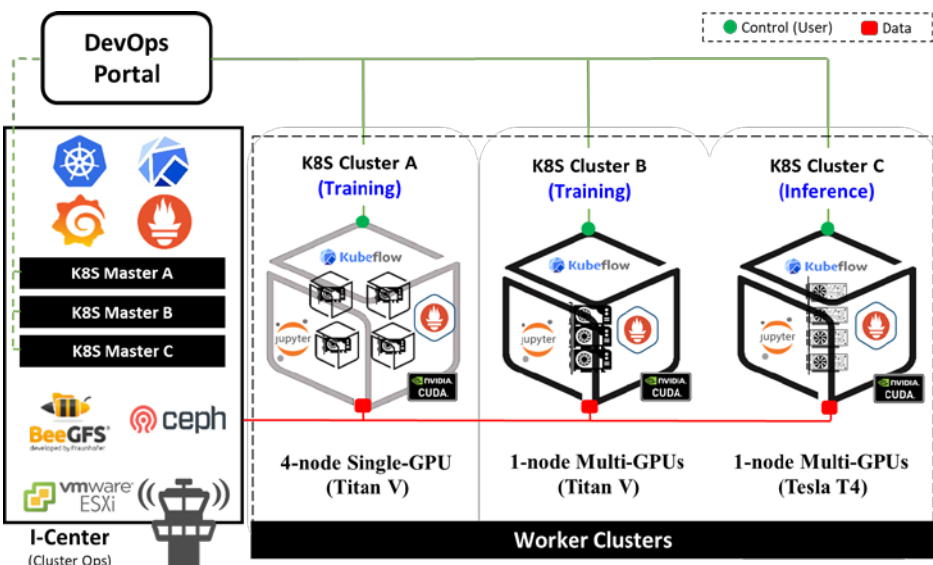
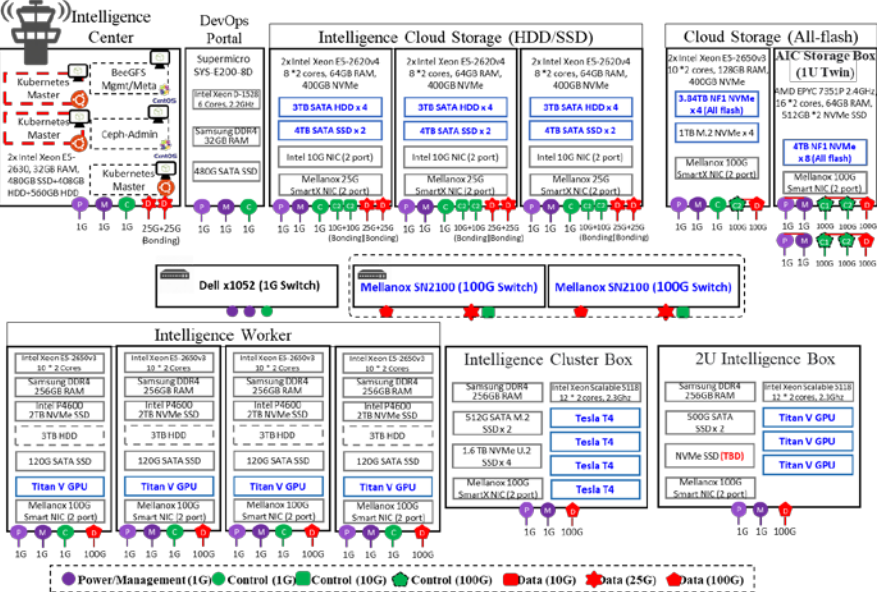
- Diversified & Affordable AI-inspired HPC/HPDA Workloads
- A scalable and flexible cluster satisfying Data 3Vs (Velocity, Volume, Variety) for intelligence services
- Emerging cloud-native paradigm with containerization and container orchestration for agile, flexible, and resource efficient computing

Concept and Design



- A shared resource pool of compute/storage/networking with hyper-converged boxes and networking fabric
- Kubernetes-based cloud-native computing with open source software collections
- Improving workload performance by coordinating hardware/software configuration with practical experiments

Hardware and Software



Issues and Solutions

Issue	Options	Solution	Reason
High-speed Physical Interconnect	Ethernet, RoCE	RoCE	▪ RoCE shows at most 30 times lower networking delay
Container Runtime for HPC/HPDA/AI	Docker, Singularity, Shifter	Docker	▪ Docker is a popular runtime due to its versatility and usability
Container Networking Interface (CNI)	Calico, Weave	Calico	▪ Calico (with BGP-based overlay networking) shows higher networking performance
File System & Container Storage Interface (CSI)	BeeGFS, NFS, CephFS	BeeGFS, CephFS	▪ BeeGFS for a caching layer in intelligence worker nodes ▪ CephFS for storing/serving massive amount of data