

# Accelerator-Aware Kubernetes Scheduler for DNN Tasks on Edge Computing Environment

ACM/IEEE 6<sup>th</sup> Symposium on Edge Computing (Poster), 2021.12



Jungae Park, Unho Choi, Kyungyong Lee (Kookmin University)  
and Seungwoo Kum, Jaewon Moon (KETI)

<http://reactor.kmubigdata.cloud>

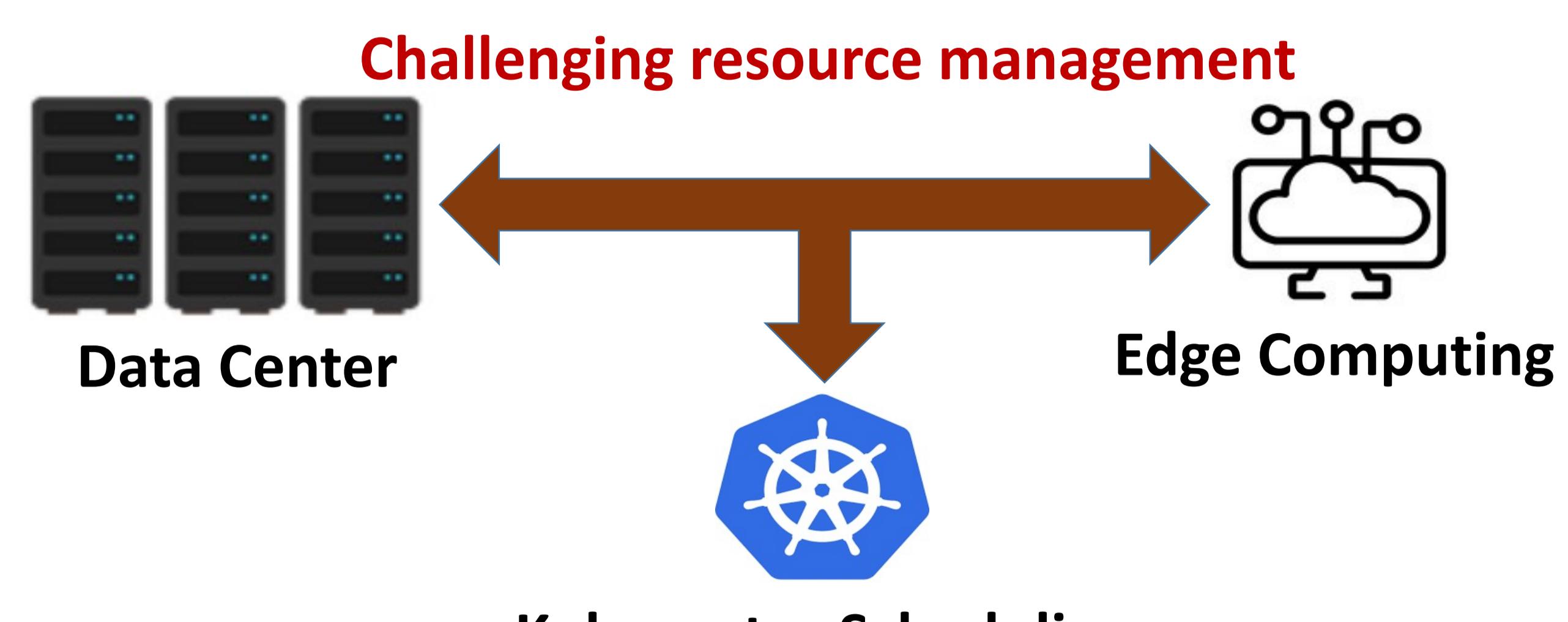
## Motivation and Challenge

### ❖ The challenge of edge computing

**Edge Computing** : data processing close to the source  
- Decentralized and globally located edge resources  
- Difficult to efficiently manage edge resources

### ❖ Opportunities and challenges of using kubernetes for edge resource management

**Kubernetes** : a well-established central resource management platform  
- Easy application deployment using containers  
- Limited resource information support from kubernetes



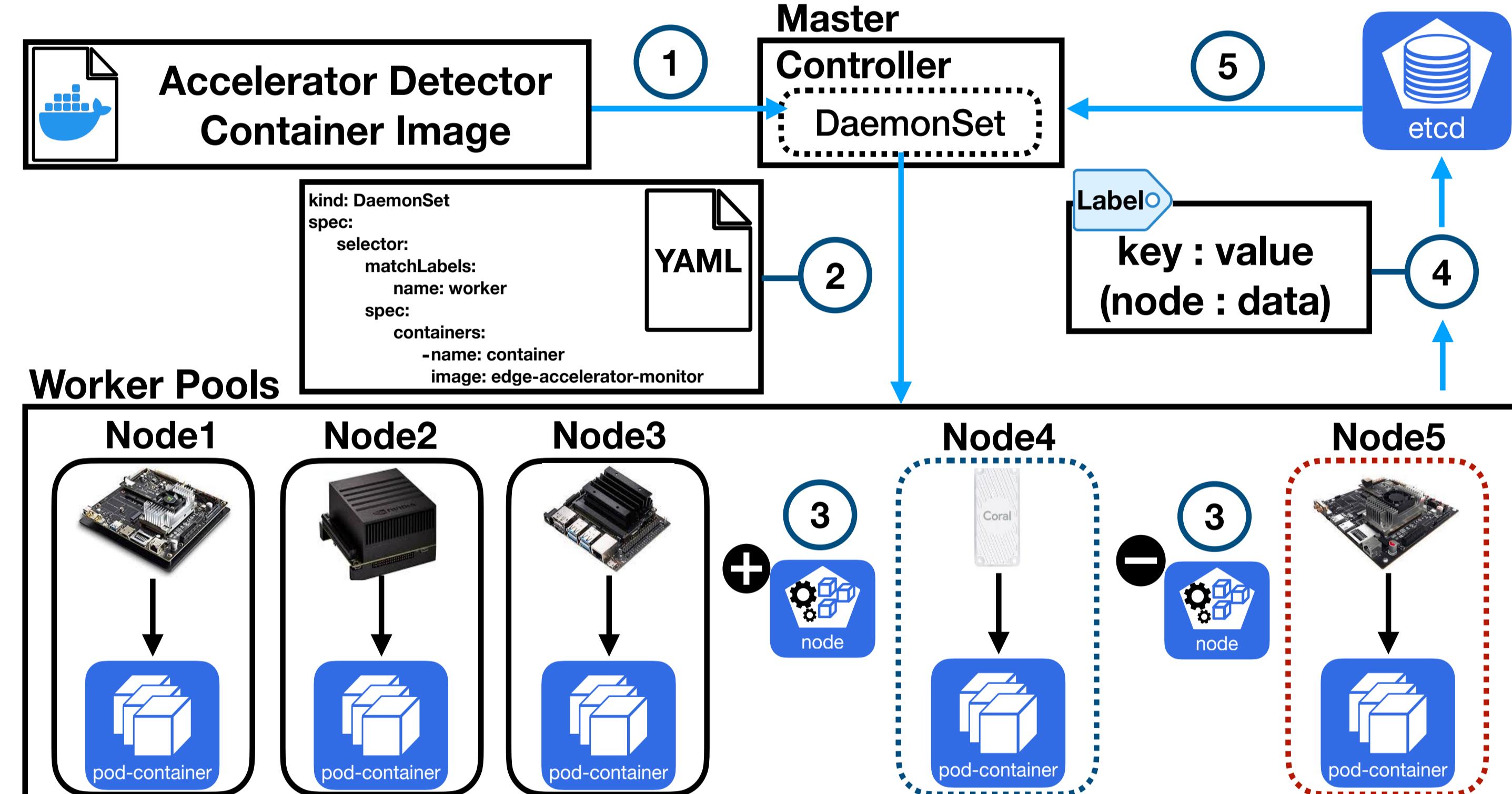
*Limited resource information support from kubernetes*

### ❖ Automatic edge accelerator hardware detector for kubernetes

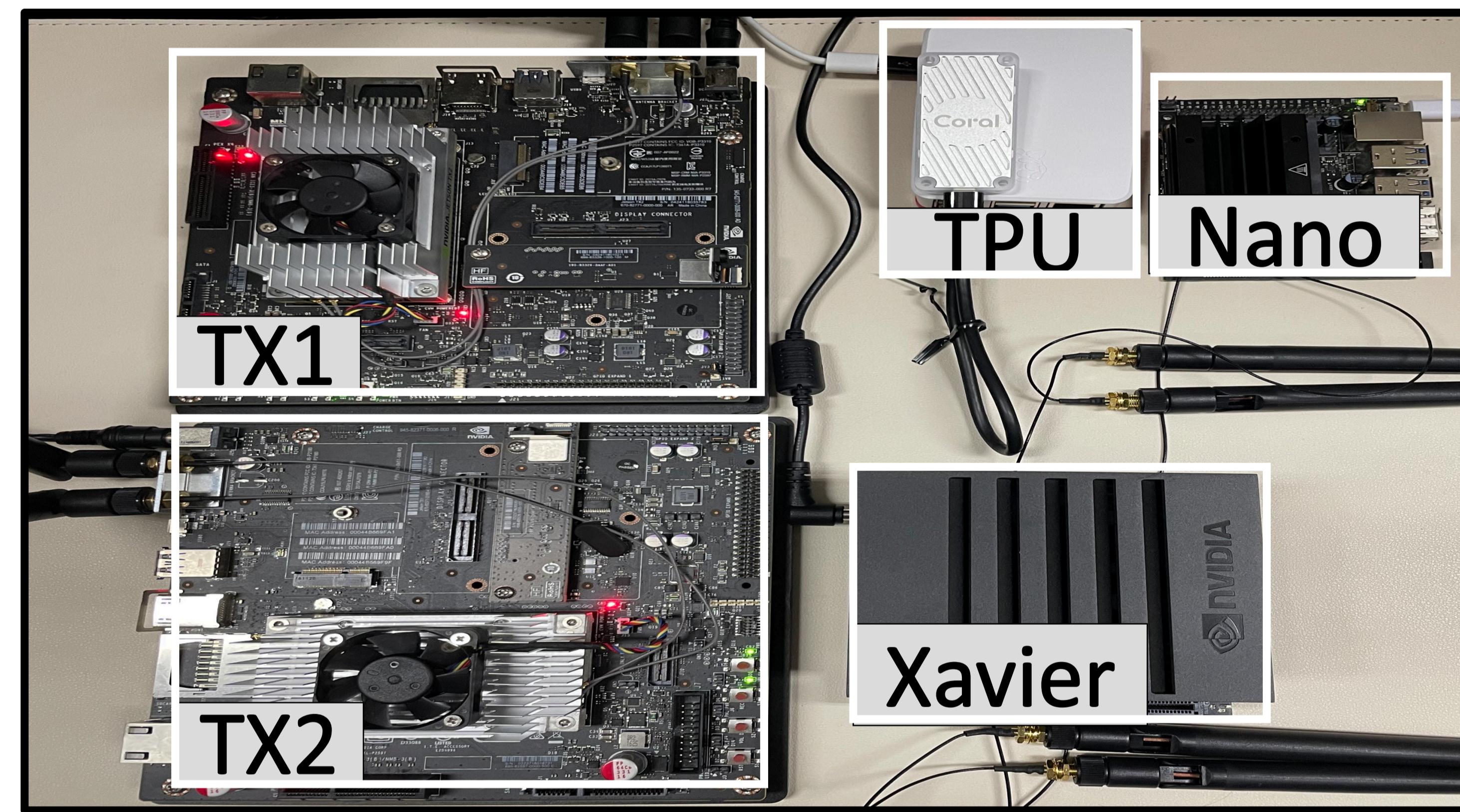
- Providing rich hardware information for kubernetes scheduler

## Implementation

### Kubernetes cluster



### ❖ Prototype implementation



### ❖ Modules for automatic accelerator detector

- Hardware information extractor container image
- DaemonSet for automatic detector container deployment
- ServiceAccount for automatic labeling

### ❖ Extracting accelerator information

- GPU driver (/etc/nv\_tegra\_release file)
- GPU model (/proc, /sys file)
- GPU resource (tegrastats, nvidia-smi command)

### ❖ Hardware information automatic labeling

```
root@node1:~# kubectl get node --show-labels
NAME    STATUS   ROLES    AGE     VERSION   LABELS
node1   Ready    master   9d     v1.18.14  kubernetes.io/arch=amd64,kubernetes.io/hostname=node1,
node2   Ready    worker   43h    v1.18.14  node-role.kubernetes.io/master=master
node3   Ready    worker   9d     v1.18.14  kubernetes.io/arch=arm64 gpu.driver=5.1,
node4   Ready    worker   9d     v1.18.14  gpu.model=Jetson-TX1,gpu.resource=tegrastats,
node5   Ready    worker   9d     v1.18.14  kubernetes.io/arch=arm64 gpu.driver=5.1,
node6   Ready    worker   9d     v1.18.14  gpu.model=Jetson-Nano,gpu.resource=tegrastats,
node7   Ready    worker   9d     v1.18.14  kubernetes.io/arch=arm64,gpu.model=Google-Coral-TPU,
node8   Ready    worker   9d     v1.18.14  kubernetes.io/arch=arm64,gpu.driver=3.1,
node9   Ready    worker   9d     v1.18.14  gpu.model=Jetson-AGX,gpu.resource=tegrastats,
```

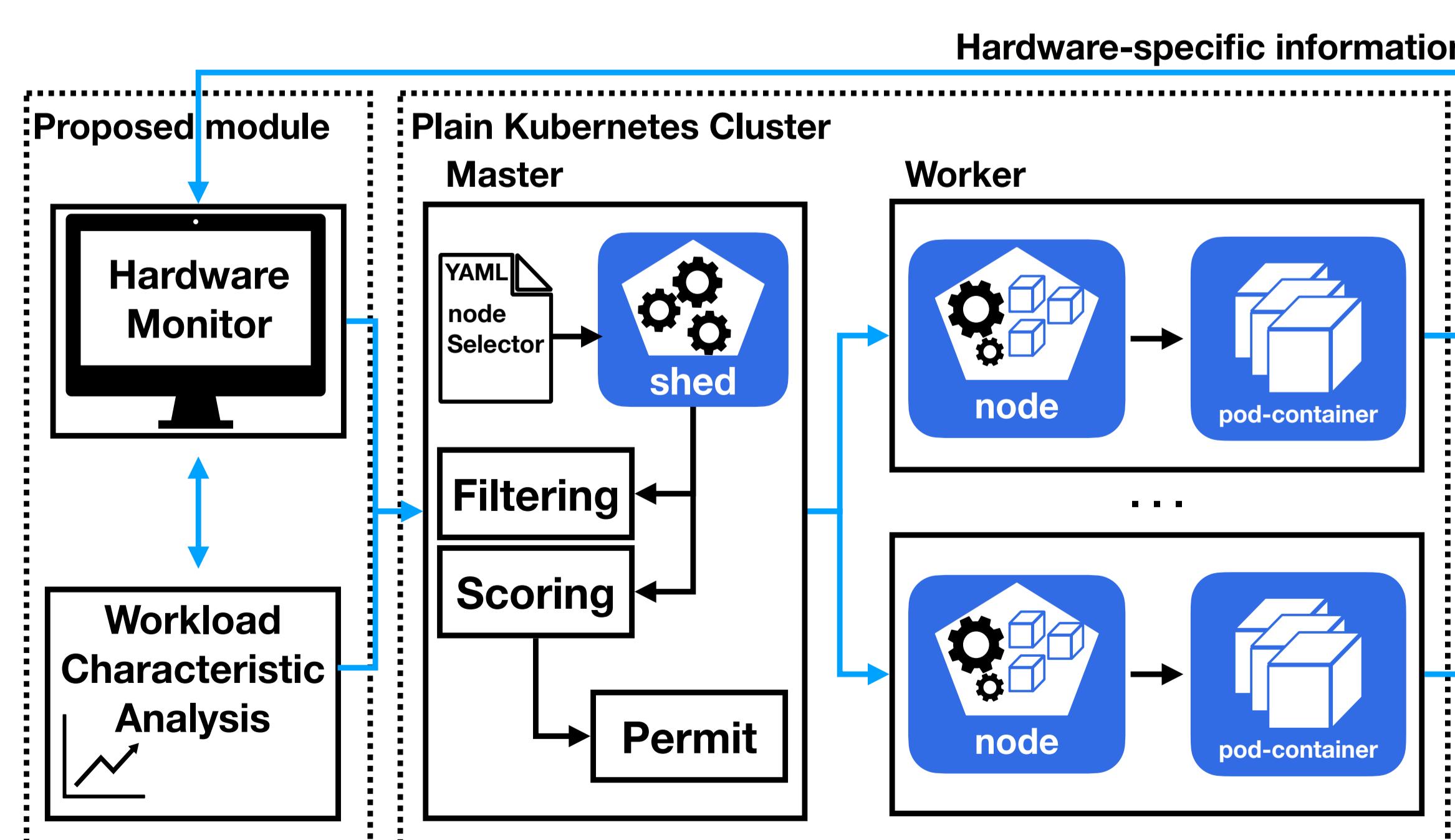
## Discussion and Future Work

### ❖ Hardware monitor module

- Gathering edge accelerator device hardware information

### ❖ Workload character analysis module

- Modeling workload characteristics (ex: DNN inference)



### ❖ Plain kubernetes scheduler

- Limited scheduling capability
- Workload-ignorant scheduling

### ❖ Workload-aware scheduler for kubernetes

- Using rich hardware and workload information

*Improvement of existing kubernetes scheduling, covering various workloads*