



GEdge(Griffin-Edge) Platform

- 초저지연 지능형 클라우드 엣지 SW 플랫폼 -

# 엣지 기반 AI/ML 서비스 흐름관리 및 실행환경 최적화 기술

2022.12.20

GS-AI 프레임워크 코어개발자

서동윤 (dyseo@softonnet.com)

“GEdge Platform”은 클라우드 중심의 엣지 컴퓨팅 플랫폼을 제공하기 위한  
핵심 SW 기술 개발 커뮤니티 및 개발 결과물의 코드명입니다.

- Developer-Driven

**GEdge Platform Community 5<sup>th</sup> Conference** (GEdge Platform v3.0 Release) -

# Contents

---

- I 지능형 서비스 운용 프레임워크 개요
- II 지능형 서비스 워크플로 및 실행관리 기술
- III 추후 연구 내용

# GEdge 플랫폼 내 GS-Aiflow의 포지셔닝

## 초저지연 지능형 클라우드 엣지 플랫폼 (GEdge Platform)

### 클라우드 엣지 관리 플랫폼 (GM : GEdge Management Platform)

플랫폼 관리 도구 프레임워크 (GM-Tool)

Framework I/F

플랫폼 관리 기능 프레임워크 (GM-Center)

Platform I/F

### 지능형 서비스 운용 프레임워크 (GS-AI)

엣지 AI 서비스 환경  
(GS-Aiflow)

엣지 협업 학습 환경  
(GS-Optops)

### 서비스 협업 프레임워크 (GS-Link)

협업 게이트웨이  
(GS-Linkgw)

협업 정책 생성  
(GS-Linkhq)

Framework I/F

Framework I/F

### 초저지연 데이터 처리 프레임워크 (GS-Engine)

엣지 전용 스케줄러  
(GS-Scheduler)

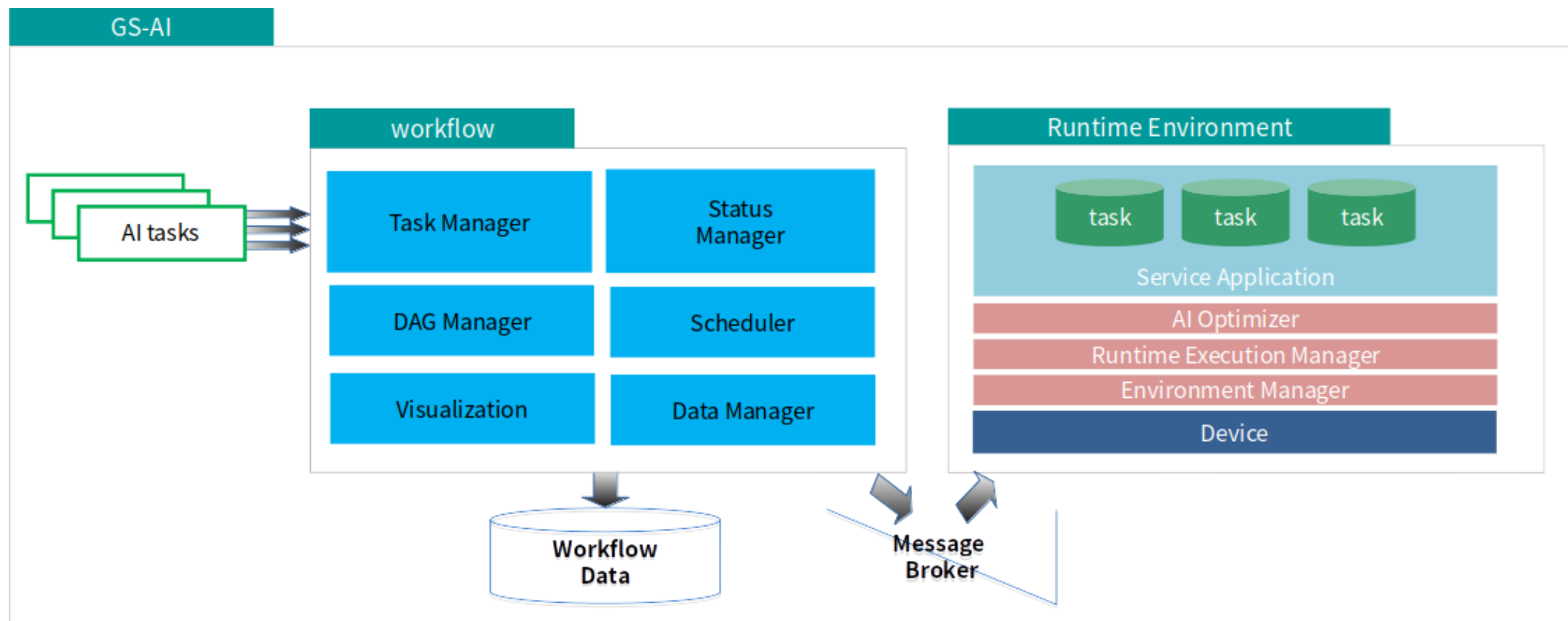
엣지 메시지 브로커  
(GS-Broker)

### 클라우드 엣지 서비스 플랫폼 (GS : GEdge Service Platform)

# 지능형 서비스 운용 프레임워크 개요



- 지능형 서비스 운용 프레임워크 구조
  - 지능형 서비스의 실행 런타임 관리 기술과 지능형 서비스의 배포를 위한 워크플로 운영 기술로 구성

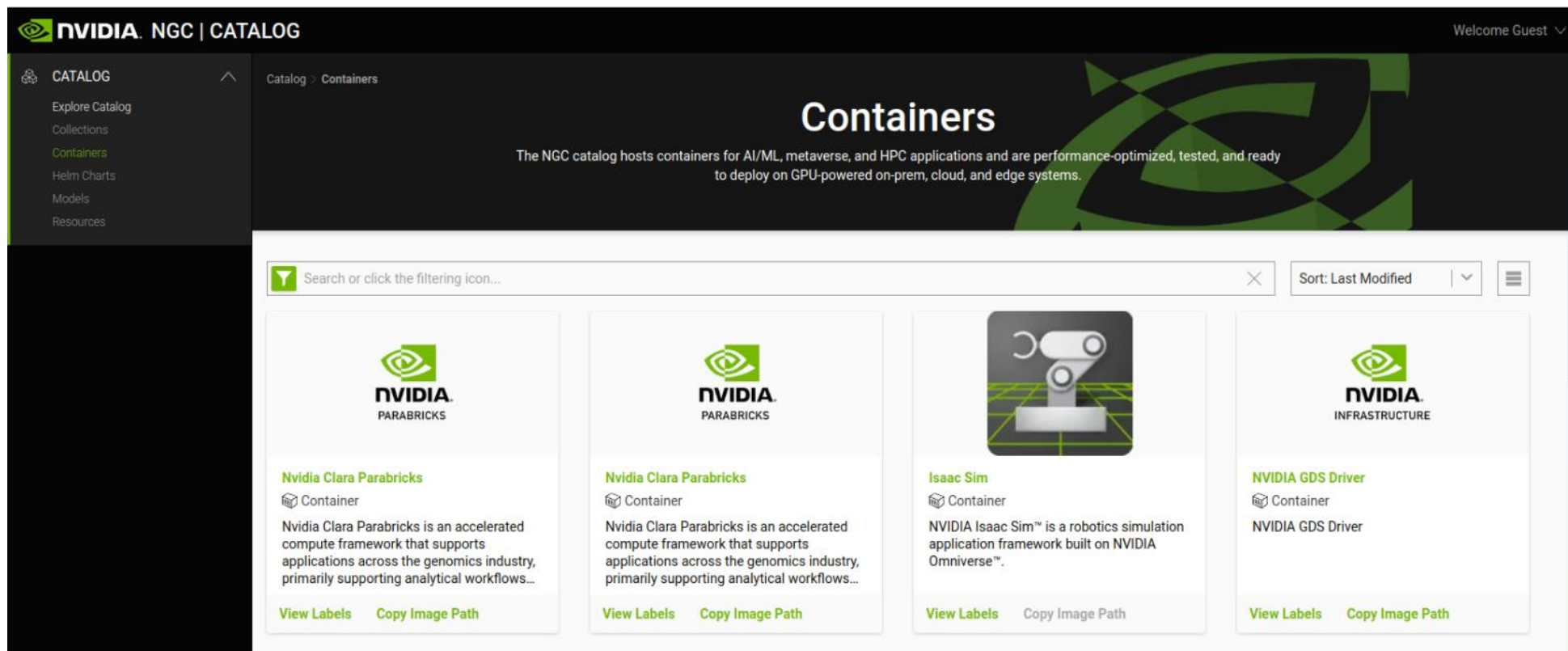




# 지능형 서비스 워크플로 및 실행 관리 기술



- GPU 워크스테이션, 서버 등 직접 구축
- NGC(NVIDIA GPU CLOUD) 컨테이너 사용
  - 클라우드 시스템에서 컨테이너 런타임으로 지능형 서비스를 배포하기 위해 제공되는 빌드 된 이미지 사용
  - Triton Inference Server 등의 이미지 추론 서버 구축도 컨테이너 형태로 제공



- 컨테이너 이미지의 크기의 거대화
- 릴리즈 이미지 별로 구성 라이브러리가 고정
  - OS, NVIDIA 라이브러리, 파이썬 언어, 딥러닝 프레임워크 버전의 고정
  - 여러 버전을 테스트 하기에 반복되는 이미지 풀링의 문제

Image Name	Timestamp	Size	Architectures	Layers	Size	Age	Size
22.11-tf2-py3	11/23/2022 5:24 AM	6.83 GB	2 Architectures	calico/node	v3.24.1	75392e3500e3	3 months ago 223MB
				nvcr.io/nvidia/tensorrt	22.06-py3	4b949560bb05	6 months ago 6.21GB
				anibali/pytorch	1.8.1-cuda11.1	8a146abe8a61	7 months ago 7.25GB
				nvcr.io/nvidia/pytorch	22.04-py3	6884f16521ea	8 months ago 14.1GB
				nvcr.io/nvidia/tensorflow	22.04-tf1-py3	b799c63410e6	8 months ago 14.4GB
				python	3.6.15-alpine	3a2e80fa4606	12 months ago 40.7MB

### PyTorch Release 22.10

The NVIDIA container image for PyTorch, release 22.10, is available on [NGC](#).

#### Contents of the PyTorch container

This container image contains the complete source of the version of PyTorch environment ( /opt/conda/lib/python3.8/site-packages/torch/ )

The container also includes the following:

- Ubuntu 20.04 including Python 3.8
- NVIDIA CUDA@ 11.8.0
- NVIDIA cuBLAS 11.11.3.6
- NVIDIA cuDNN 8.6.0.163
- NVIDIA NCCL 2.15.5 (optimized for NVIDIA NVLink@)
- NVIDIA RAPIDS™ 22.08.01 (For x86, only these libraries are included: cu)
- Apex
- rdma-core 36.0
- NVIDIA HPC-X 2.12.2tp1
- OpenMPI 4.1.5a1
- GDRCopy 2.3
- TensorBoard 2.10.0
- Nsight Compute 2022.3.0.22
- Nsight Systems 2022.4.2.1
- NVIDIA TensorRT™ 8.5.0.12
- Torch-TensorRT 1.1.0a0
- NVIDIA DALI@ 1.18.0
- MAGMA 2.6.2
- JupyterLab 2.3.2 including Jupyter-TensorBoard
- TransformerEngine 0.1.0

### PyTorch Release 20.11

The NVIDIA container image for PyTorch, release 20.11, is available on [NGC](#).

#### Contents of the PyTorch container

This container image contains the complete source of the version of PyTorch i environment ( /opt/conda/lib/python3.6/site-packages/torch/ ) i

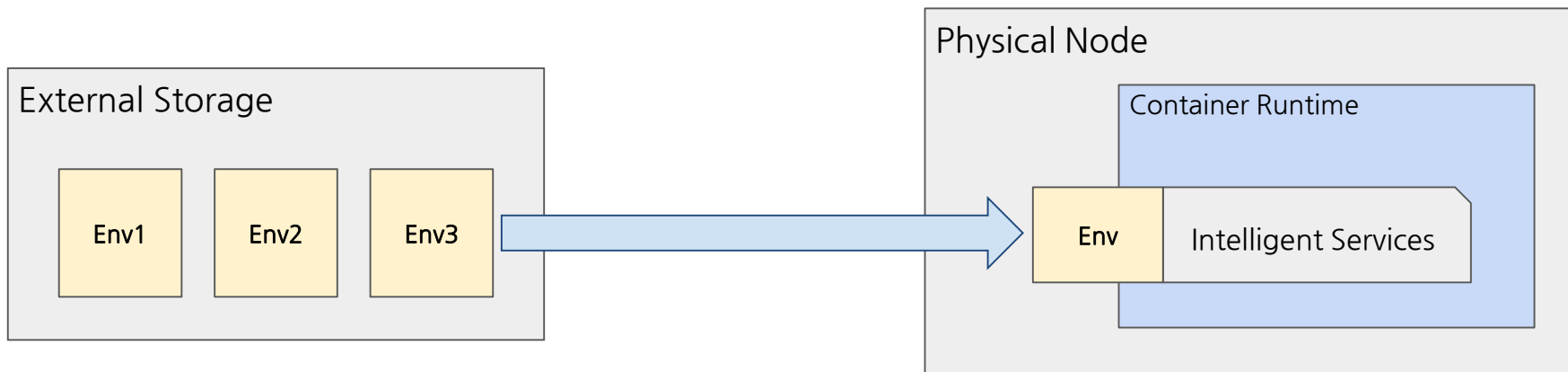
The container also includes the following:

- Ubuntu 18.04 including Python 3.6 environment
- NVIDIA CUDA 11.1.0 including cuBLAS 11.2.1
- NVIDIA cuDNN 8.0.4
- NVIDIA NCCL 2.8.2 (optimized for NVLink™)
- APEX
- MLNX\_OFED 5.1
- OpenMPI 4.0.5
- TensorBoard 1.15.0+nv20.11
- Nsight Compute 2020.2.0.18
- Nsight Systems 2020.3.4.32
- TensorRT 7.2.1
- DALI 0.27.0
- MAGMA 2.5.2
- DLProf 0.17.0
- PyProf r20.11
- Tensor Core optimized examples:



- 컨테이너 내부에 구현 될 런타임 환경에 대해 외부에서 이식가능하게 제공
  - 제공된 환경에 대해 다른 컨테이너 인스턴스에서 재사용이 용이해야 하며 재시작등의 관리의 편의성 고려
- 제공된 런타임 환경의 요소들을 미리 구축하고 필요한 요소만 선택적으로 제공
  - 컨테이너 런타임을 실행 시 사용자의 요구사항에 따라 원하는 환경 요소들을 구성할 수 있어야 한다.

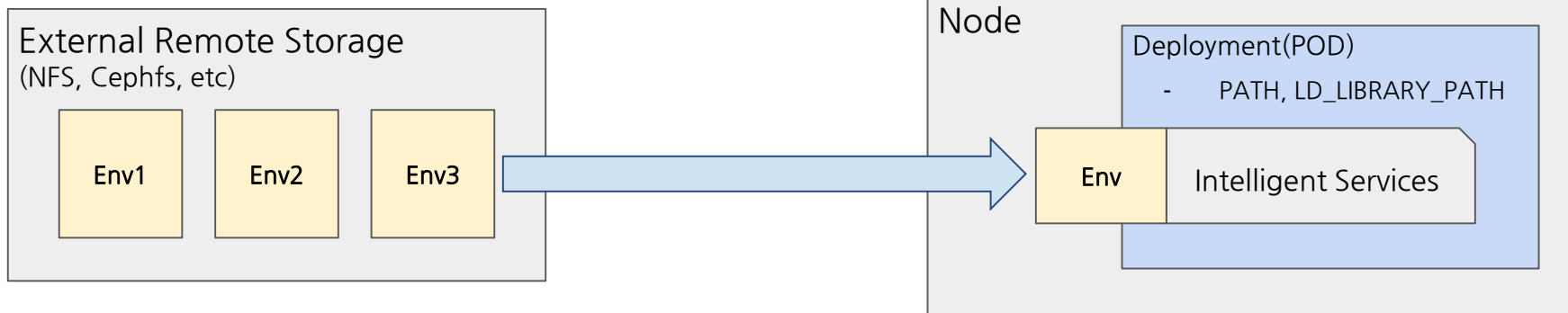
Implantable Storage + Additional processing



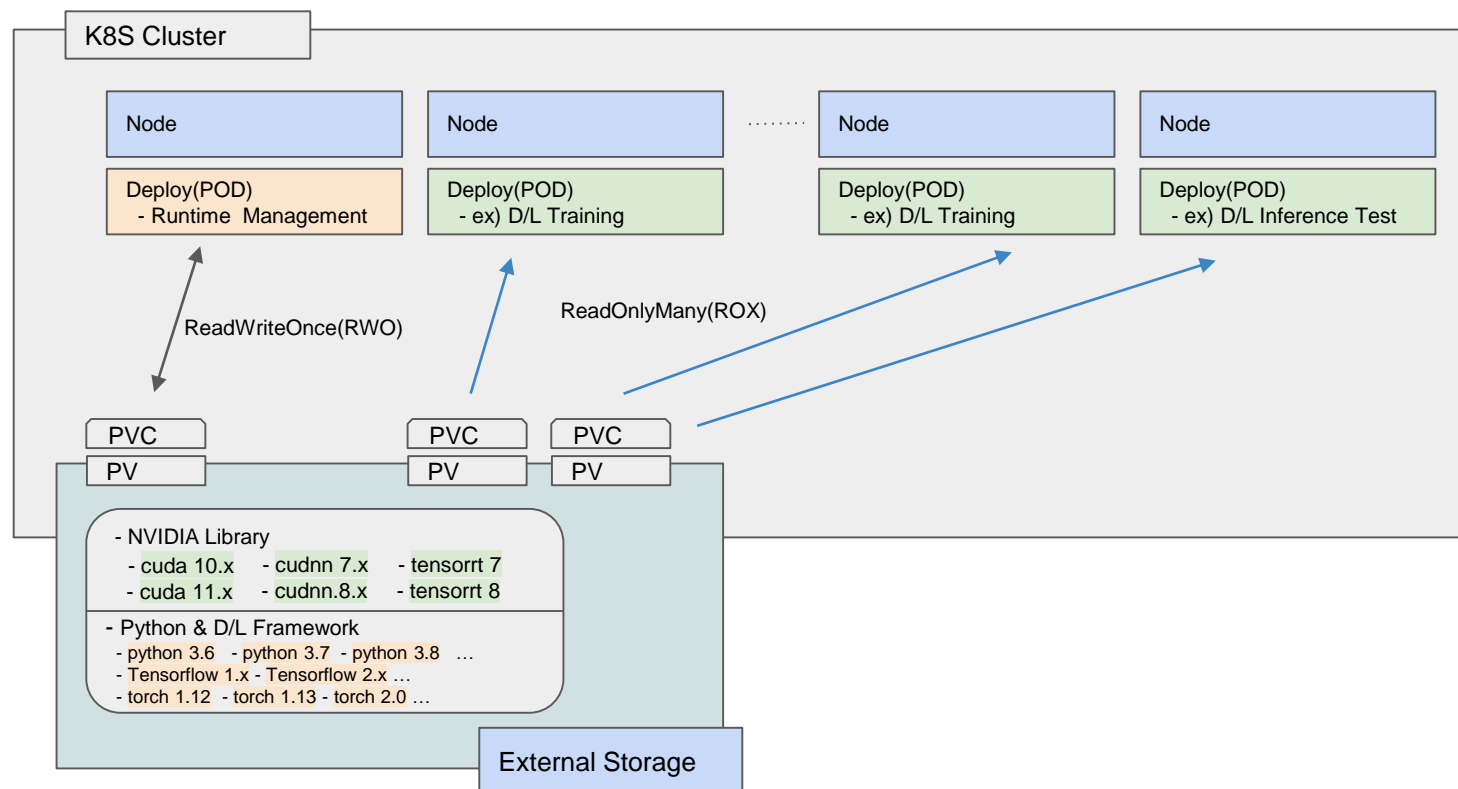
- 쿠버네티스 클러스터에서 접근 가능한 오브젝트 연구
  - 스토리지에 영속성 볼륨(Persistent Volume)으로 환경 구성요소들에 대해 사전 구성
  - 컨테이너 런타임 배포 시 영속성 볼륨 클레임(Persistent Volume Claim)을 이용해 필요한 볼륨을 이식
- 볼륨 연결 후 해당 환경을 사용하기 위한 추가 조치 필요
  - 리눅스 기반 시스템에서 라이브러리를 사용하기 위한 환경변수 등의 연결



File/Object Storage(NFS, Cephfs, etc)  
+ Binary PATH, LD\_LIBRARY\_PATH, etc



- 가상환경의 사전 구성
  - NVIDIA 라이브러리의 설치 방식 중 PATH 지정 방식 설치 활용
  - 파이썬 기반 딥러닝 프레임워크의 패키지 제공방식에 따라 버전 별로 제공
- 권한 분리로 인한 제공의 용이성 증대
  - 리눅스 기반 시스템에서 라이브러리를 사용하기 위한 환경변수 등의 연결



- 가상환경의 사전 구성
  - NVIDIA 라이브러리의 설치 방식 중 PATH 지정 방식 설치 활용
  - 파이썬 기반 딥러닝 프레임워크의 패키지 제공방식에 따라 버전 별로 제공
- 권한 분리로 인한 제공의 용이성 증대
  - 리눅스 기반 시스템에서 라이브러리를 사용하기 위한 환경변수 등의 연결

```
command: ["/bin/bash", "-c"]
args: ["source /root/path.sh; env; tail -f /dev/null"]
env:
- name: NEW_PATH
  value: '/root/volume/miniconda3/envs/tf2_py38/bin:/root/volume/cuda/cuda-11.2/bin'
- name: LD_LIBRARY_PATH
  value: '/root/volume/cuda/lib64:/root/volume/cudnn/lib64'
```

```
(base) smarteye@son:~/kube/nfs_test$ sudo kubectl exec -it pod-ubuntu-20.04 -- env | grep PATH
PATH=/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin
NEW_PATH=/root/volume/miniconda3/envs/tf2_py38/bin:/root/volume/cuda/cuda-11.2/bin
LD_LIBRARY_PATH=/root/volume/cuda/lib64:/root/volume/cudnn/lib64
```

```
volumeMounts:
- mountPath: "/root/path.sh"
  name: fileconfig
  subPath: path.sh
- mountPath: "/root/volume/cuda"
  name: nfs-volume-total
  subPath: cuda/cuda-11.2
- mountPath: "/root/volume/miniconda3"
  name: nfs-volume-total
  subPath: miniconda3
- mountPath: "/root/volume/cudnn"
  name: nfs-volume-total
  subPath: cudnn/cudnn-v8.2.1-11.x-linux-x64
volumes:
- name: nfs-volume-total
  persistentVolumeClaim:
    claimName: nfs-pvc-total
```

- 개발 및 서비스 배포환경 간 스위칭
  - 딥러닝 학습 환경, 어플리케이션 개발환경, 서비스 릴리즈 환경 등에 적용 가능
- 지능형 서비스 어플리케이션 릴리즈 비용의 경량화
  - Pyinstaller 기반 파이썬 어플리케이션 배포
  - 같은 어플리케이션이어도 빌드되는 환경에 따라 용량 차이 발생


NGC Based Container



Light OS Container + Virtual Env



- 클러스터 등록 및 클러스터 내 가동 오브젝트 모니터링
- 지능형 서비스의 워크플로 정의 및 서비스 배포 모니터링
  - 사전 정의 된 지능형 서비스로 테스트 중
  - 실행되고 있는 오브젝트의 로그 모니터링 기능



Monitoring - Project

Create - Project

Delete - Project

Monitoring - Cluster

Register - Cluster

### Register - Monitoring Cluster

#### Cluster Monitoring Register

Cluster Name

Cluster IP(Master Node)

Kubernetes API Port

Cluster API Token

등록

#### Cluster Monitoring Delete

cluster name

삭제

- 클러스터 오브젝트 모니터링
  - K8S API서버와 통신하여 클러스터 내 가동되는 오브젝트들의 스테이터스 모니터링

Deployment

Deployment Name		Collision	Available Replicas		Ready Replicas	Replicas
-----------------	--	-----------	--------------------	--	----------------	----------

Node IP(cluster\_test1)

Node IP	Type	Last Heartbeat	Last Transition	message	status
172.16.20.101	InternalIP	Thu, 15 Dec 2022 07:36:49 GMT	Fri, 28 Oct 2022 08:51:29 GMT	kubelet is posting ready status. AppArmor enabled	True
172.16.20.90	InternalIP	Thu, 15 Dec 2022 07:40:21 GMT	Fri, 28 Oct 2022 08:51:28 GMT	kubelet is posting ready status. AppArmor enabled	True
172.16.30.50	InternalIP	Thu, 15 Dec 2022 07:37:52 GMT	Tue, 13 Dec 2022 09:26:58 GMT	kubelet is posting ready status. AppArmor enabled	True
172.16.30.55	InternalIP	Thu, 15 Dec 2022 07:39:05 GMT	Fri, 28 Oct 2022 08:51:23 GMT	kubelet is posting ready status	True
172.16.30.56	InternalIP	Mon, 14 Nov 2022 04:51:01 GMT	Mon, 14 Nov 2022 04:51:47 GMT	Kubelet stopped posting node status.	Unknown

Pod Info (cluster\_test1,namespace=default)

default

Nodename	Podname	Host IP	Pod IP	Phase	StartTime
20-101	aieye-inference	172.16.20.101	10.244.4.73	Running	Thu, 15 Dec 2022 07:22:01 GMT
sapeon-test	aieye-mariadb	172.16.30.50	10.244.3.74	Running	Thu, 15 Dec 2022 07:20:56 GMT
20-101	aieye-media-server	172.16.20.101	10.244.4.72	Running	Thu, 15 Dec 2022 07:21:39 GMT
20-101	aieye-postprocessing	172.16.20.101	10.244.4.74	Running	Thu, 15 Dec 2022 07:22:28 GMT
sapeon-test	aieye-redisserver	172.16.30.50	10.244.3.73	Running	Thu, 15 Dec 2022 07:20:33 GMT
sapeon-test	aieye-web-dashboard	172.16.30.50	10.244.3.75	Running	Thu, 15 Dec 2022 07:23:05 GMT

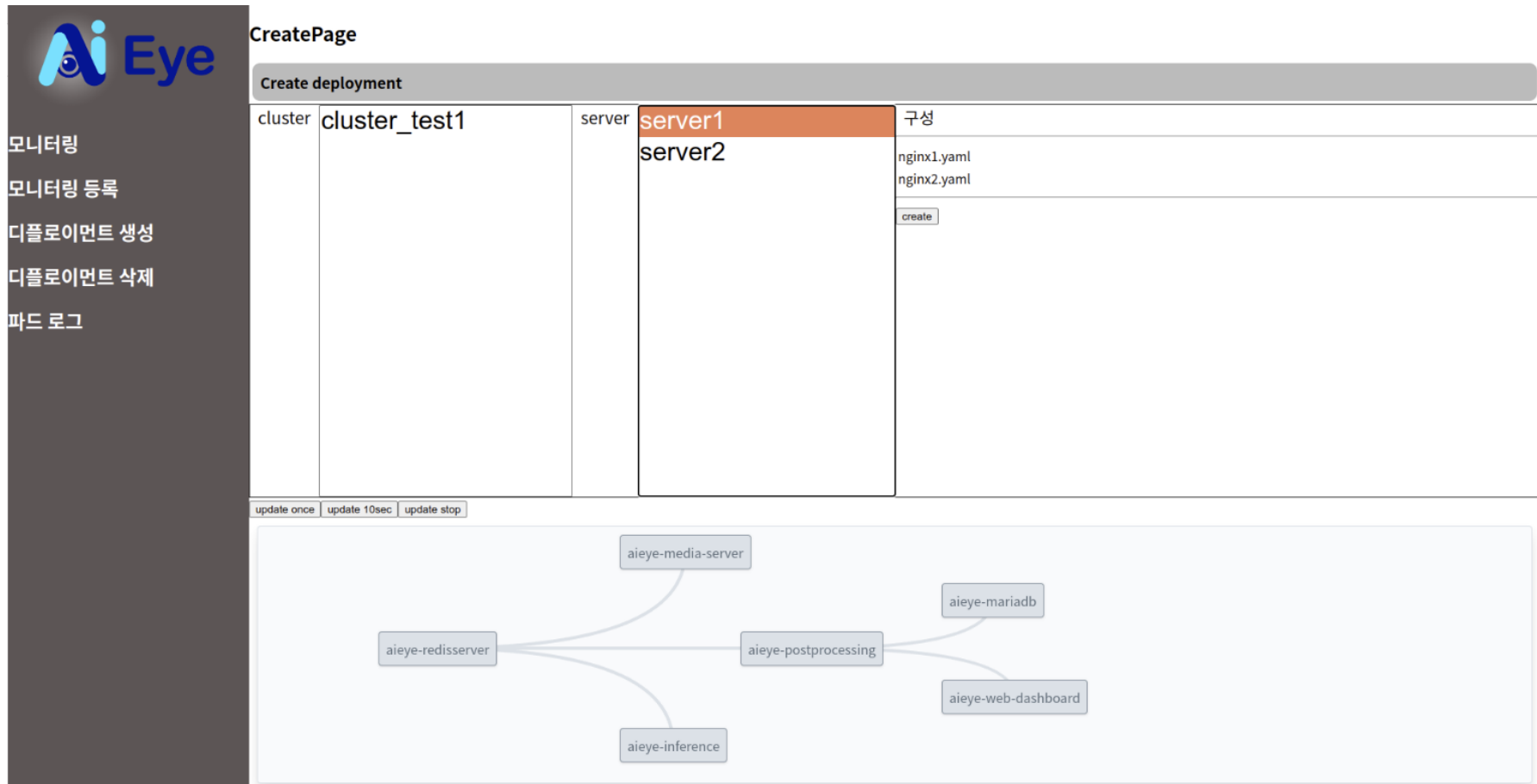
PV Info(cluster\_test1)

PV Name	PV Namespace	PV AccessMode	PV Capacity	PV Volumemode	PV Localpath	PV Phase	PV Createiontime	PV Claim Name	PV Claim Namespace
---------	--------------	---------------	-------------	---------------	--------------	----------	------------------	---------------	--------------------

SC Info(cluster\_test1)

SC Name	SC Namespace	AllowVolumeExpansion	AllowedTopologies	Provisioner	Reclaim Policy	Volume Binding Mode	Creationtime
aieye-local-sc				kubernetes.io/no-provisioner	Delete	WaitForFirstConsumer	Tue, 01 Nov 2022 09:03:08 GMT

- 사전 정의 한 지능형 서비스의 배포
  - 지능형 서비스의 워크플로 구성도 확인 가능하도록 정의



The screenshot displays the 'CreatePage' interface for 'Ai Eye'. On the left is a dark sidebar with the 'Ai Eye' logo and a menu containing: '모니터링' (Monitoring), '모니터링 등록' (Monitoring Registration), '디플로이먼트 생성' (Deployment Creation), '디플로이먼트 삭제' (Deployment Deletion), and '파드 로그' (Pod Logs). The main area is titled 'CreatePage' and features a 'Create deployment' button. Below this is a table for configuring the deployment:


cluster	server	구성
cluster_test1	server1	nginx1.yaml
	server2	nginx2.yaml

Below the table, there are buttons for 'update once', 'update 10sec', and 'update stop'. At the bottom, a workflow diagram shows the following components and their connections:

- 'aieye-redisserver' connects to 'aieye-media-server', 'aieye-inference', and 'aieye-postprocessing'.
- 'aieye-media-server' connects to 'aieye-postprocessing'.
- 'aieye-postprocessing' connects to 'aieye-mariadb' and 'aieye-web-dashboard'.
- 'aieye-inference' connects to 'aieye-postprocessing'.



- 사전 정의 한 지능형 서비스의 모니터링
  - 지능형 서비스를 구성하는 파드의 상태 실시간 모니터링 및 로그 표출



모니터링

모니터링 등록

디플로이먼트 생성

디플로이먼트 삭제

파드 로그

**logviewer**

cluster  
cluster\_test1

namespace  
default

pod  
aieye-imp

This container image and its contents are governed by the NVIDIA Deep Learning Container License. By pulling and using the container, you accept the terms and conditions of this license: <https://developer.nvidia.com/ngc/nvidia-deep-learning-container-license>

NOTE: The SHMEM allocation limit is set to the default of 64MB. This may be insufficient for TensorFlow. NVIDIA recommends the use of the following flags:  
docker run --gpus all --ipc-host --ulimit memlock=-1 --ulimit stack=67108864 ...

sed: -e expression #1, char 22: unknown option to 's'

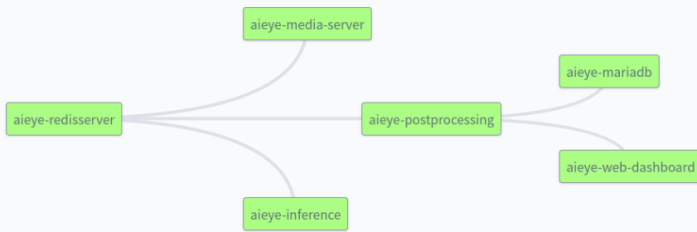
Checking servers

```
[2022-12-15 17:08:26,420]-[INFO]-[117:140278679820096]-[multilingual] > Language has set English
[2022-12-15 17:08:26,420]-[INFO]-[117:140278679820096]-[main] >
[2022-12-15 17:08:26,420]-[INFO]-[117:140278679820096]-[main] > Starting Analyze Server v1.1.4050
[2022-12-15 17:08:26,421]-[INFO]-[117:140278679820096]-[main] >
[2022-12-15 17:08:26,423]-[INFO]-[117:140278679820096]-[main] > License Checking.....
[2022-12-15 17:08:26,424]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License File Path : bin/license_analyze_for_etri.dat
[2022-12-15 17:08:26,424]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License Mode : common.license.edition.2
[2022-12-15 17:08:26,424]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License macaddr : 00:00:00:00:00:00
[2022-12-15 17:08:26,424]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License description : use for Doosan(Industrial)
[2022-12-15 17:08:26,424]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License expire date : 2022-12-31
[2022-12-15 17:08:26,425]-[INFO]-[117:140278679820096]-[LicenseKeyEngine] > License state : Normal
[2022-12-15 17:08:26,425]-[INFO]-[117:140278679820096]-[main] > Starting in daemon mode.
[2022-12-15 17:08:26,427]-[DEBUG]-[117:140278679820096]-[database] > 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:26,427]-[ERROR]-[117:140278679820096]-[main] > Failed to connect to database: 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:31,433]-[DEBUG]-[117:140278679820096]-[database] > 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:31,434]-[ERROR]-[117:140278679820096]-[main] > Failed to connect to database: 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:36,440]-[DEBUG]-[117:140278679820096]-[database] > 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:36,441]-[ERROR]-[117:140278679820096]-[main] > Failed to connect to database: 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:41,447]-[DEBUG]-[117:140278679820096]-[database] > 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
[2022-12-15 17:08:41,448]-[ERROR]-[117:140278679820096]-[main] > Failed to connect to database: 2003 (HY000): Can't connect to MySQL server on 'localhost:9333' (99)
```

**Cluster Info**

Cluster Name	MasterNode IP	Port	Token Exist	Search Node&Pod Info
cluster_test1	172.16.20.90	6443	True	<input type="text" value="Search"/>

server1 ▼  
update once update 10sec update stop





## 추후 연구 내용



- ML/DL
  - 범용 및 플랫폼 동작환경에 맞춘 네트워크 모델 최적화 기술 고도화
    - Deployment 프레임워크: TensorRT
- 런타임 환경 관리 기술
  - 프레임워크 별 런타임 환경 구성/배포 자동화 기술 개발
- 워크플로 모니터링 도구
  - 런타임 환경 관리 기술과 연계한 서비스 내 환경 배포 기능 개발
  - 사전 정의(pre-defined)된 지능형 서비스 고도화
    - 다양한 pre-defined 서비스 추가
    - 지능형 서비스 성능 개선

# 감사합니다.

<http://gedge-platform.github.io>



GS-Aiflow 코어 개발자

서동윤 (dyseo@softonnet.com)

## Welcome to GEdge Platform

An Open Cloud Edge SW Platform to enable Intelligent Edge Service

### GEdge Platform will lead Cloud-Edge Collaboration