

Model Serving

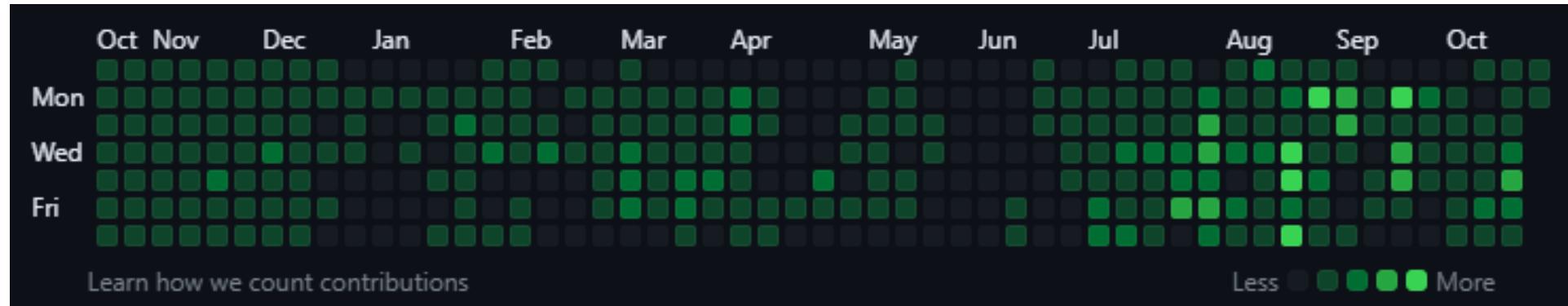
BOAZ MLOps 소모임

19기 분석 오효근

발표자 소개

▪ 오효근

- BOAZ 1971 분석 부문 수료
- AgileSoDA Machine Vision MLE 개작 중
- Algorithm, MLOps, 자동화에 관심이 많습니다!



GitHub



Blog

Contents

- 1. Introduction**
- 2. Preparing**
 - 1. Export**
 - 2. Triton Inference Server**
- 3. Model Serving**
 - 1. Docker Compose**
 - 2. Kubernetes**
 - 3. Ensemble**
 - 4. Gradio**

Introduction

■ Model Serving

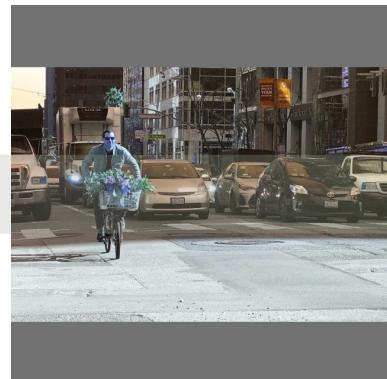
- Definition: 학습된 machine learning model을 실제 환경에서 사용할 수 있도록 하는 과정

▪ Tools / Platforms

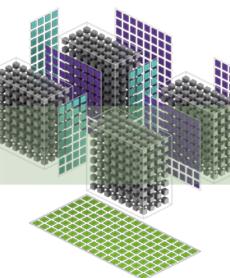
- TensorFlow Serving
- TorchServe
- Triton Inference Server
- MLflow



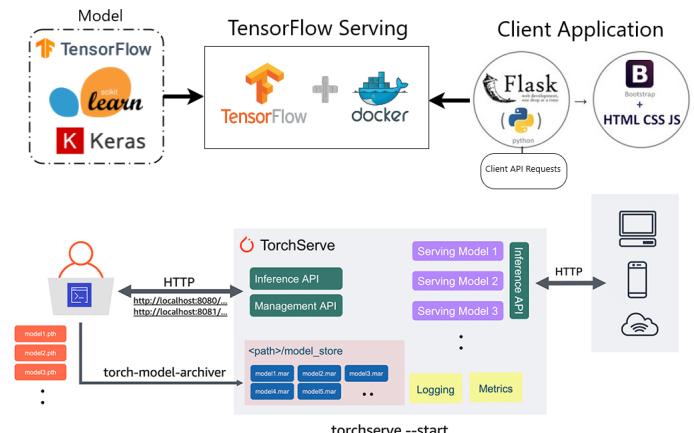
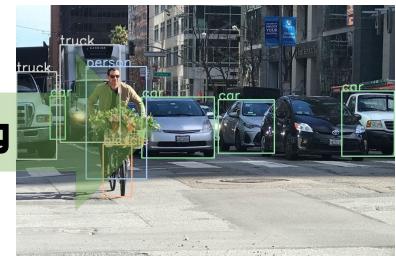
Pre-processing



Inferencing



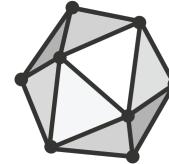
Post-processing



2.1. Export

■ ONNX (Open Neural Network Exchange)

- 다양한 플랫폼에서 실행 가능
- Deep learning에 대한 표준화 제공



ONNX

```
1 git clone https://github.com/ultralytics/yolov5
2 cd yolov5
3 conda create -n yolov5 python=3.8 -y
4 conda activate yolov5
5 pip install -r requirements.txt
6 python export.py --weights yolov5s.pt --include onnx
```

```
Export complete (4.7s)
Results saved to /home/zerohertz/Zerohertz/yolov5
Detect:    python detect.py --weights yolov5s.onnx
Validate:   python val.py --weights yolov5s.onnx
PyTorch Hub: model = torch.hub.load('ultralytics/yolov5', 'custom', 'yolov5s.onnx')
Visualize:  https://netron.app

ls
benchmarks.py  classify      data        export.py  LICENSE  README.md    requirements.txt  setup.cfg  tutorial.ipynb  val.py       yolov5s.pt
CITATION.cff  CONTRIBUTING.md detect.py  hubconf.py  models    README.zh-CN.md segment      train.py    utils           yolov5s.onnx
```

2.2. Triton Inference Server

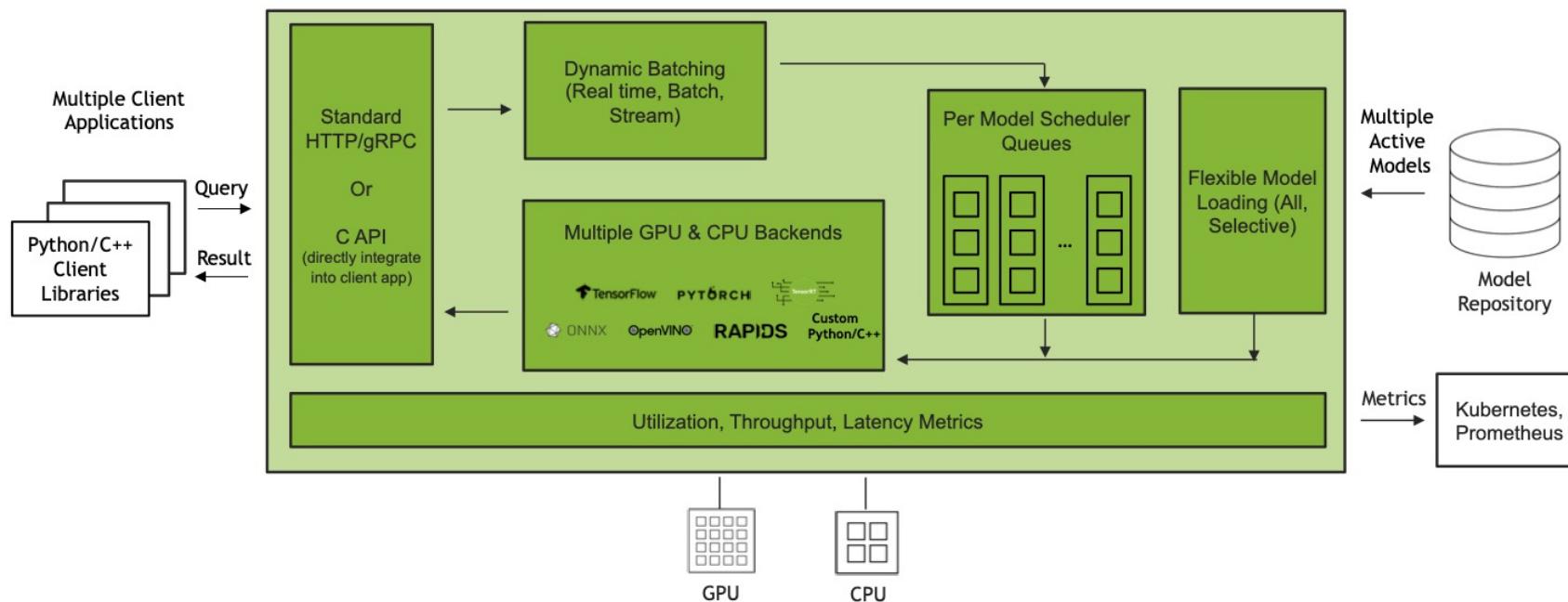
- **Multiple Client Applications:** 다양한 클라이언트 애플리케이션들이 Triton Inference Server에 큐리를 보낼 수 있습니다.
- **Python/C++ Client Libraries:** Triton은 Python과 C++을 위한 클라이언트 라이브러리를 제공하여 사용자들이 서버에 큐리를 쉽게 보낼 수 있게 합니다.
- **Standard HTTP/gRPC:** 클라이언트 애플리케이션은 HTTP 또는 gRPC를 통해 서버에 접속할 수 있습니다.
- **C API:** 클라이언트 애플리케이션은 C API를 사용하여 직접 Triton과 통합될 수 있습니다.
- **Dynamic Batching:** Triton은 실시간, 배치, 스트림과 같은 다양한 요청 유형을 동적으로 처리할 수 있는 능력을 가지고 있습니다.
- **Multiple GPU & CPU Backends:** 다양한 백엔드 프레임워크를 지원하여 GPU와 CPU에서 동작할 수 있습니다. TensorFlow, PyTorch, ONNX, OpenVINO, RAPIDS, Custom Python/C++ 백엔드 등을 지원합니다.
- **Per Model Scheduler Queues:** 각 모델마다 스케줄러 큐가 있어 추론 요청을 관리합니다.
- **Flexible Model Loading:** 모든 모델을 로드하거나 선택적으로 모델을 로드할 수 있는 유연성을 제공합니다.
- **Model Repository:** 모든 활성화된 모델들은 모델 저장소에서 관리됩니다.
- **Metrics:** 서버의 성능과 관련된 다양한 메트릭들을 제공합니다. (예: 활용도, 처리량, 지역 시간)
- **Kubernetes, Prometheus:** Triton은 Kubernetes와 Prometheus와 같은 플랫폼 및 도구와 통합하여 스케일링과 모니터링을 지원합니다.

2.2. Triton Inference Server

- NVIDIA에서 제공하는 여러 AI 모델들을 동시에 서빙하도록 설계된 open source 추론 서버

NVIDIA TRITON INFERENCE SERVER ARCHITECTURE

Open-Source Software For Scalable, Simplified Inference Serving



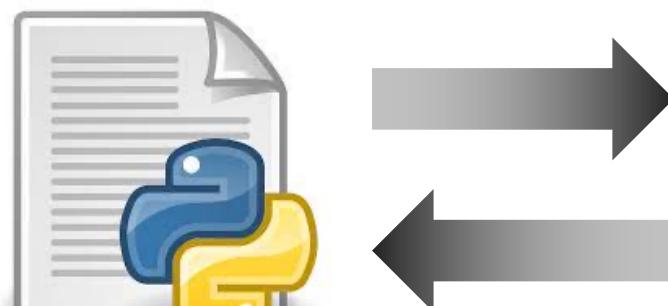
2.2. Triton Inference Server

■ Service Port

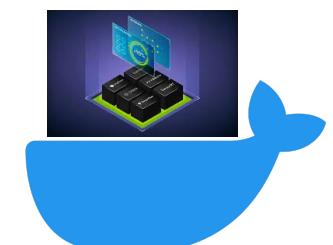
- 8000: RestAPI 를 위한 HTTPService 포트
- 8001: gRPC 를 위한 GRPCInferenceService 포트
- 8002: Metric 수집을 위한 포트

```
1 def inference(input_image: np.ndarray):
2     SERVER_URL = "0.0.0.0:8001"
3     MODEL_NAME = "YOLO"
4     cv2.imwrite("input_image.jpg", input_image)
5     input_image = input_image.astype("float32")
6     input_image = input_image.transpose((2, 0, 1))[np.newaxis, :]
7     input_image = np.ascontiguousarray(input_image)
8     print("input_image.shape:", input_image.shape)
9     with grpcclient.InferenceServerClient(SERVER_URL) as triton_client:
10         inputs = [
11             grpcclient.InferInput(
12                 "images", input_image.shape,
13                 np_to_triton.dtype(np.float32)
14             )
15             inputs[0].set_data_from_numpy(input_image)
16             outputs = [grpcclient.InferRequestedOutput("output0")]
17             response = triton_client.infer(
18                 model_name=MODEL_NAME, inputs=inputs, outputs=outputs
19             )
20             response.get_response()
21             output0 = response.as_numpy("output0")
22             print("output0.shape:", output0.shape)
23             return output0
```

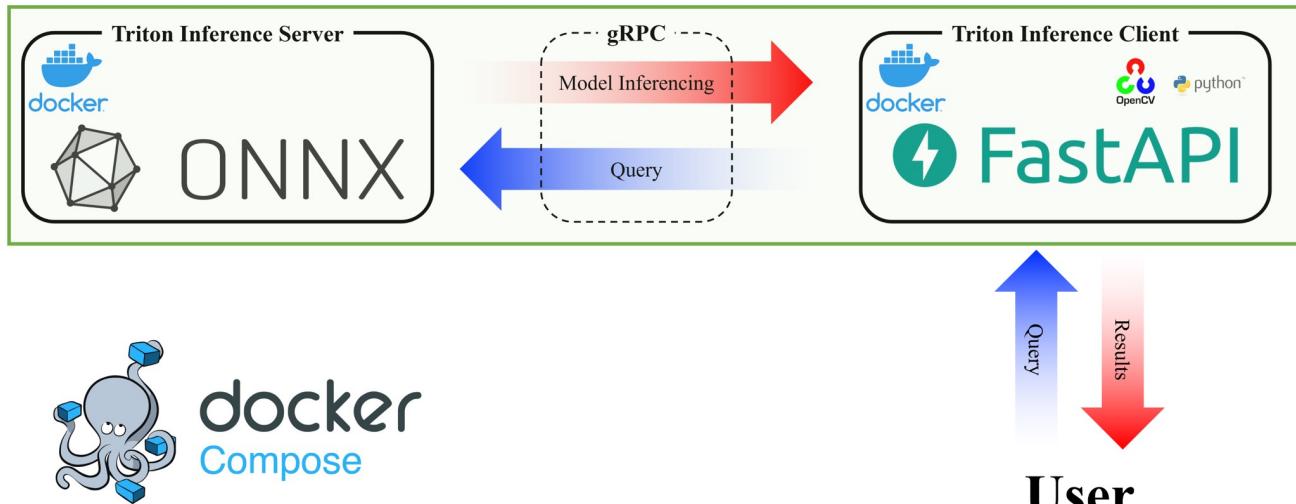
```
1 # I1023 1 grpc_server.cc:2445] Started GRPCInferenceService at 0.0.0.0:8001
2 # I1023 1 http_server.cc:3555] Started HTTPService at 0.0.0.0:8000
3 # I1023 1 http_server.cc:185] Started Metrics Service at 0.0.0.0:8002
```



 Client.py

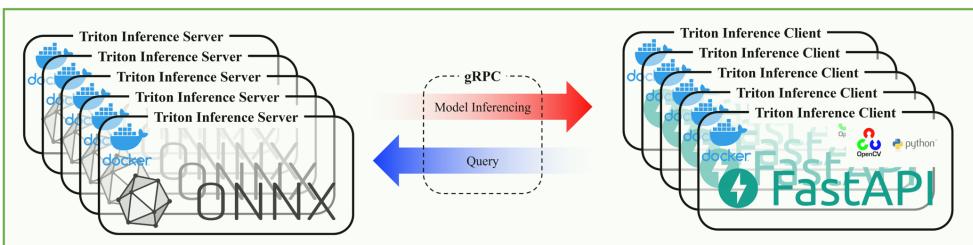
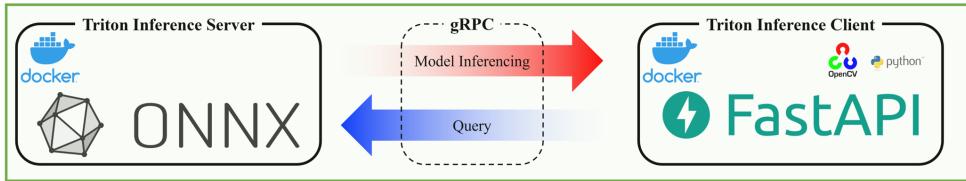
 docker
Triton Inference Server
Container

3.1. Docker Compose



```
1 version: "3.5"
2
3 services:
4   tritoninferenceserver:
5     container_name: server
6     image: nvcr.io/nvidia/tritonserver:23.06-py3
7     environment:
8       - NVIDIA_VISIBLE_DEVICES=0
9     volumes:
10      - ${PWD}/src/server:/models
11     command: tritonserver --model-
12   repository=/models
13   fastapi:
14     container_name: client
15     build:
16       context: .
17       dockerfile: ./src/client/Dockerfile
18     ports:
19       - "80:80"
20
21 networks:
22   default:
23     name: mlops-network
```

3.2. Kubernetes

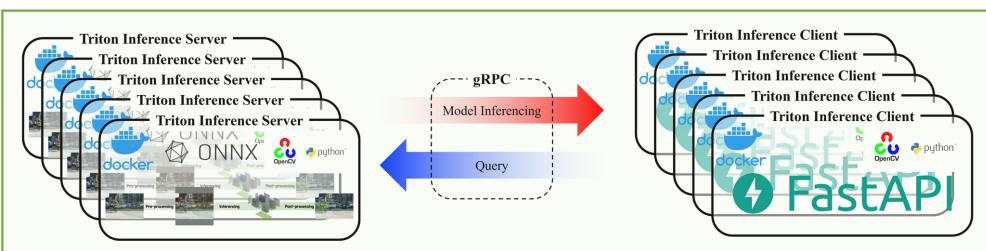
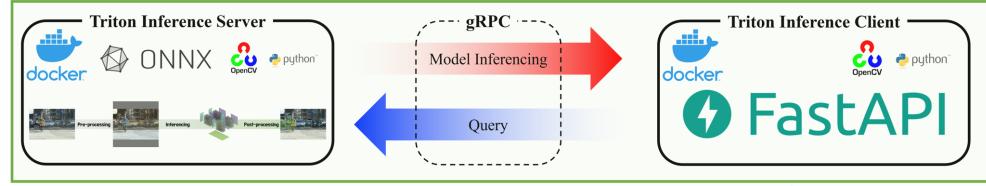


```
1 apiVersion: apps/v1
2 kind: Deployment
3 metadata:
4   name: triton-inference-server
5 spec:
6   replicas: 5
7   selector:
8     matchLabels:
9       app: triton-inference-server
10    template:
11      metadata:
12        labels:
13          app: triton-inference-server
14      spec:
15        containers:
16          - name: server
17            image: yolo-server:latest
18            imagePullPolicy: IfNotPresent
19            env:
20              - name: NVIDIA_VISIBLE_DEVICES
21                value: "0"
22            command:
23              ["tritonserver", "--model-repository=/models", "--log-
24 verbose=2"]
25            resources:
26              requests:
27                memory: "1024Mi"
28                cpu: "1000m"
29              limits:
30                memory: "2048Mi"
31                cpu: "2000m"
32  ---
33 apiVersion: v1
34 kind: Service
35 metadata:
36   name: triton-inference-server-svc
37   spec:
38     selector:
39       app: triton-inference-server
40     ports:
41       - protocol: TCP
42         port: 8001
43         targetPort: 8001
```

3.3. Ensemble

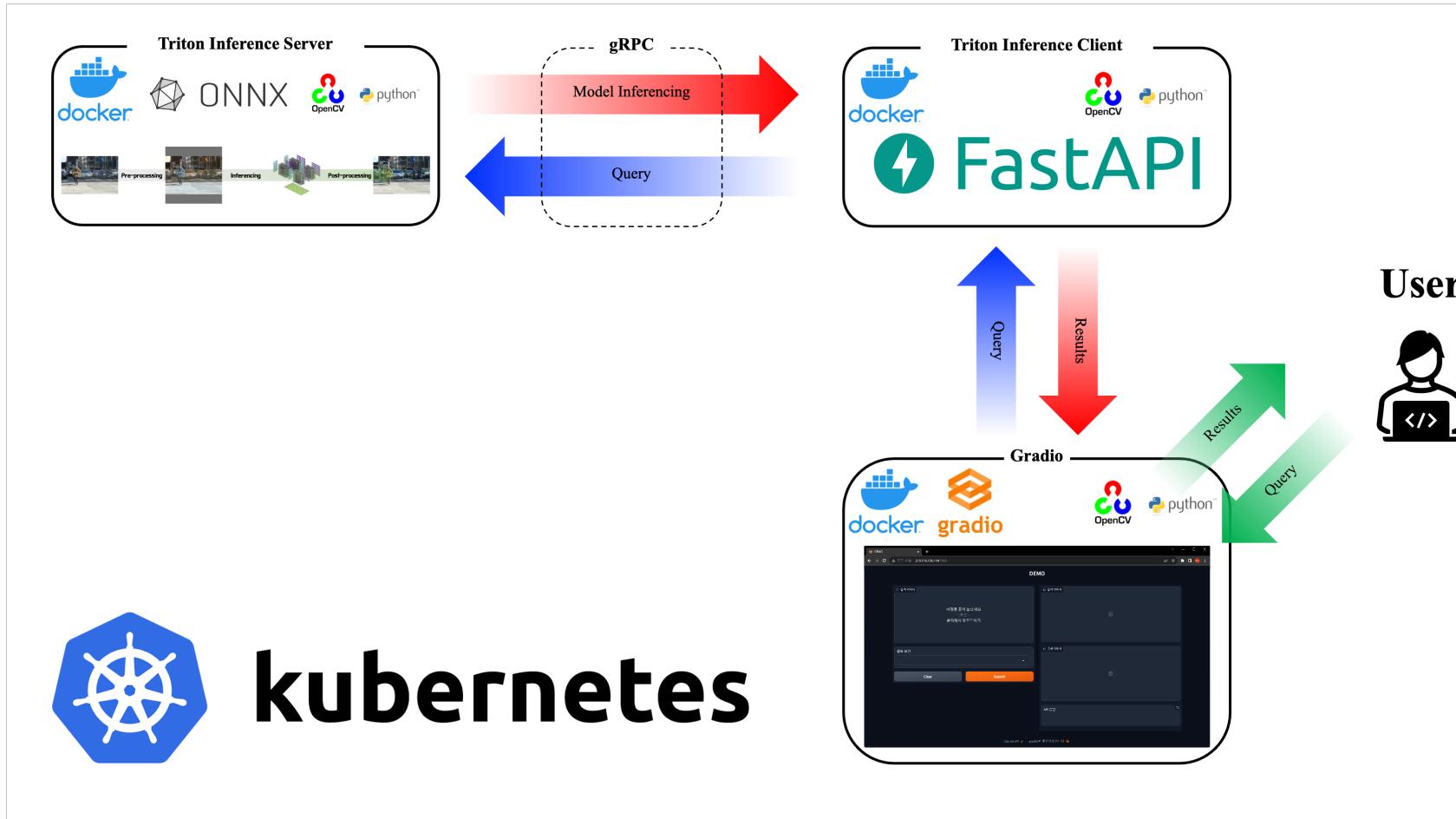
- 여러 모델들을 연결하여 하나의 단일 모델처럼 작동하도록 설정

```
1 name: "YOLO"
2 platform: "ensemble"
3
4 input [
5   {
6     name: "INPUT_IMAGE"
7     data_type: TYPE_UINT8
8     dims: [-1,-1,-1]
9   }
10 ]
11
12 output [
13   {
14     name: "RESULTS"
15     data_type: TYPE_FP32
16     dims: [-1,-1,-1]
17   },
18   {
19     name: "VISUALIZE"
20     data_type: TYPE_UINT8
21     dims: [-1,-1,-1]
22   }
23 ]
24
25 ensemble_scheduling {
26   step [
27     {
28       model_name: "preprocess"
29       model_version: -1
30       input_map {
31         ...
32       }
33       output_map {
34         ...
35       }
36     },
37     {
38       model_name: "model"
39       model_version: -1
40       input_map {
41         ...
42       }
43       output_map {
44         ...
45       }
46     },
47     {
48       model_name: "postprocess"
49       model_version: -1
50       input_map {
51         ...
52       }
53       output_map {
54         ...
55       }
56     }
57   ]
58 }
```

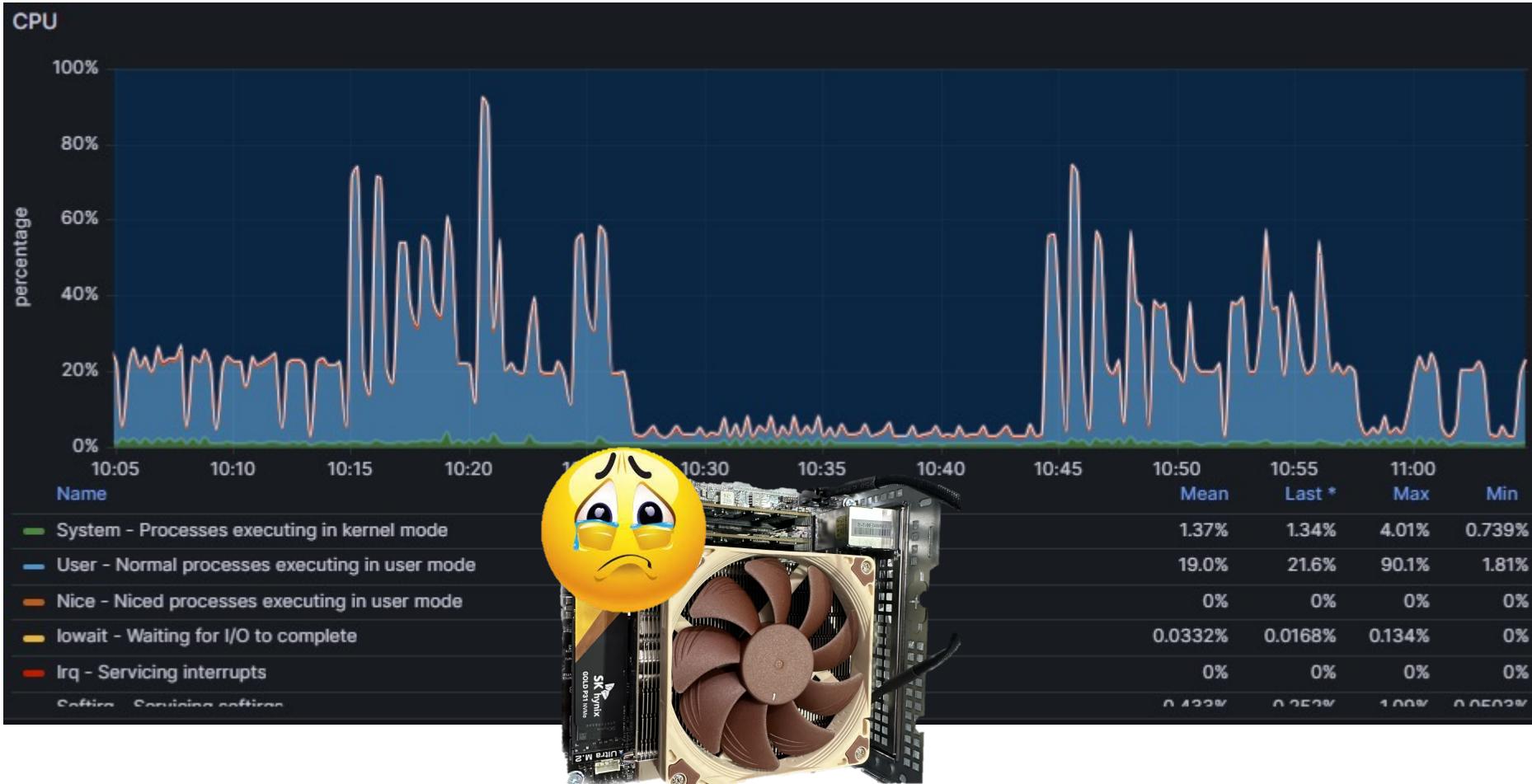


3.4. Gradio

- Machine learning model을 위한 user interface 생성 라이브러리



Special Thanks



Q & A