



Triton inference server on KISTI NEURON

☰ Contents	
🕒 Created	@July 19, 2022 9:46 AM
☰ Tags	Daily

1. Prepare a TensorRT model

- Tutorial by NVIDIA Jinho Lee
[https://github.com/leejinho610/TRT_Triton_HandsOn/blob/main/1. Model preparation \(starting Point\).ipynb](https://github.com/leejinho610/TRT_Triton_HandsOn/blob/main/1.%20Model%20preparation%20(starting%20Point).ipynb)
- **trtexec** command (ONNX → tensorRT)

```
!trtexec \  
  --onnx=model.onnx \  
  --explicitBatch \  
  --optShapes=actual_input_1:16x3x224x224 \  
  --maxShapes=actual_input_1:32x3x224x224 \  
  --minShapes=actual_input_1:1x3x224x224 \  
  --best \  
  --saveEngine=model.plan
```

- you may use TensorRT in the singularity image that I built on the KISTI server(/scratch/kedu05/leecy/tensorRT)
- Run the singularity image with “singularity run --nv tensorRT” and use the trtexec command inside the container

```

<model-repository-path>/
  <model-name>/
    [config.pbtxt]
    [<output-labels-file> ...]
    <version>/
      <model-definition-file>
    <version>/
      <model-definition-file>
    ...
  <model-name>/
    [config.pbtxt]
    [<output-labels-file> ...]
    <version>/
      <model-definition-file>
    <version>/
      <model-definition-file>
    ...
  ...

```

model directory should be in this format

2. Run tritonserver with the TensorRT model

- you may use the singularity image that I built on the KISTI server (/scratch/keu05/leecy/project1/triton.sif) or build your own singularity image with the corresponding docker image)
- Run the singularity image with “singularity run --nv triton.sif”

tritonserver --model-repository={/path/to/your/model} --allow-http=true --http-port=50712 --allow-grpc=true --grpc-port=50722 --allow-metrics=true --metric-port=50723 (SPECIFY YOUR PORT NUMBER)

- Port in singularity container will be automatically mapped to serverside port without specifying port numbers with --net flag

3. Client side request to the tritonserver

- You may run client side code with tritonserver_client both in the local environment as long as you connected your local port to the KISTI server by portforwarding. You can also run it on the KISTI server
- portforward both grpc port and http port (ex: ssh -L localhost:8000:gpu22:50712 kedu05@neuron.ksc.re.kr)
- I used docker environment in my local computer.
- docker pull [nvcr.io/nvidia/tritonserver:22.01-py3-sdk](https://hub.docker.com/r/nvcr.io/nvidia/tritonserver:22.01-py3-sdk)
- docker run -it --rm --net=host -v \$(pwd):/hello [nvcr.io/nvidia/tritonserver:22.01-py3-sdk](https://hub.docker.com/r/nvcr.io/nvidia/tritonserver:22.01-py3-sdk) (mounted local volume on the container to access to the inference dataset)
- run `client.py` (your code) in the docker container

```

1 import tritonclient.http as tritonhttpclient
2 import tritonclient.grpc as tritongrpcclient
3 import numpy as np
4 # from PIL import Image
5 # from torchvision import transforms
6 import json
7
8 ###CONFIGURATION#####
9 VERBOSE = False
10 input_name = 'modelInput'
11 input_shape = (1, 1, 128, 128, 128)
12 input_dtype = 'FP32'
13 output_name = 'modelOutput'
14 model_name = 'amyProg'
15 http_url = 'localhost:8000'
16 grpc_url = 'localhost:8001'
17 model_version = '1'
18 #####
19
20
21
22 image = np.load('/hello/target_nor.npy')
23 image = np.expand_dims(image, axis=0)
24 image = np.expand_dims(image, axis=0)
25 image = image.astype('float32')
26
27 # triton_client = tritonhttpclient.InferenceServerClient(url=http_url, verbose=VERBOSE)
28 triton_client = tritongrpcclient.InferenceServerClient(url=grpc_url, verbose=VERBOSE)
29 # model_metadata = triton_client.get_model_metadata(model_name=model_name, model_version=model_version) #You can remove this line
30 model_config = triton_client.get_model_config(model_name=model_name, model_version=model_version)
31
32 #input0 = tritonhttpclient.InferInput(input_name, input_shape, input_dtype)
33 input0 = tritongrpcclient.InferInput(input_name, input_shape, input_dtype)
34 #input0.set_data_from_numpy(image_numpy, binary_data=False)
35 input0.set_data_from_numpy(image)
36
37 #output = tritonhttpclient.InferRequestedOutput(output_name, binary_data=False)
38 output = tritongrpcclient.InferRequestedOutput(output_name)
39 response = triton_client.infer(model_name, model_version=model_version,
40                               inputs=[input0], outputs=[output])
41 logits = response.as_numpy(output_name)
42 logits = np.asarray(logits, dtype=np.float32)
43
44 print(logits)

```

client.py