

CHALLENGES OF INFERRING ON THE EDGE DEVICES :

Machine Learning is a broad domain in itself. An end to end process requires proper model building, evaluation and accurate deployment of those models. Time constraints have led us use the state of the art models and run them in our development environment but the deployment of these models is even harder than building a model pipeline. Here are some few challenges we faced while inference on the edge devices:

The various edge devices we used for inference are :

1. Edge TPU
2. Raspberry PI
3. Jetson Nano

1. **Edge TPU inference** : Inference on edge TPU requires our model to be fully quantized for weights and other parameters. Full quantization meaning converting to the unsigned integer 8 bits . For some of our processes we had already provided pretrained models by tensorflow model zoo. Object Detection, Pose Estimation, Face Detection, Image Classification and some others too. So it was quite easy using pre existing models trained on the dataset and using them for the inference. The results were quite astonishing. And deployment of the models on those devices were no pain within few lines of code using PYCORAL API and EDGERUNTIME Library provided by the Google.

However when we tried to convert these state of the art models to the smaller versions like the tensorflow lite and the edge TPU version we faced a lot of challenges while doing so :

1. Converting to tflite : The tensorflow documentation seems well defined and the examples show only few steps for the conversion of the model from .pb to .tflite but we faced a lot of issues while doing so :
 - a. The documentation did not well defined the conversion steps from .pb to tflite and hence it took quite a while converting from tf to tflite.

The model required some additional information or the metadata or The metagraph definition from the model. Which we were not sure whether it had been provided in the model or not. Since that model was not trained by us. And graph definition was also unknown. So we weren't quite sure what the input and output tensors were.

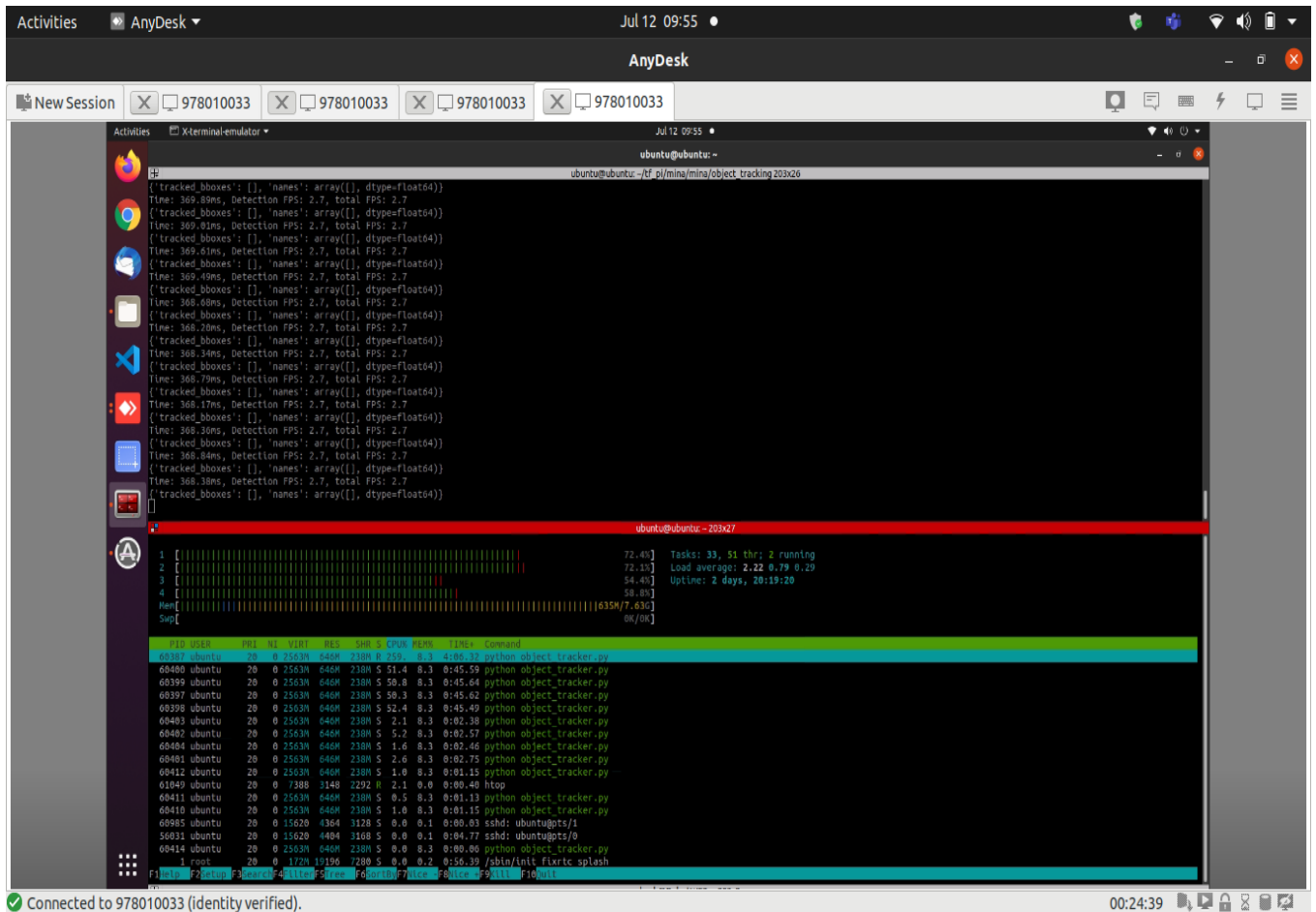
- b. The conversion path was also not well defined in official documentation we faced a lot of time solving that issue.
- c. The conversion from the tflite model to the quant process model required some dataset that were utilized during training .So it was not possible for us to train the quant model unless we went for the fine tuning and transfer learning.

2.Raspberry PI : Inference on raspberry PI was not easy either because Installing tensorflow on Raspberry PI took 3 days . And that also we weren't sure whether It will work or not . Tensorflow on raspberryPI required to be built from the source using bazelisk .

- 1.Building bazel was a very hard task at first.
- 2.The binary file execution error was consistent.
- 3.Once build later the file vanished
- 4. It took 3 days to finally install Tensorflow on our PI system. Later we directly build the tensorflow from the wheel files found on the github.

We were still not sure whether it would work on our system or whether inference could be performed with such limited space and memory limitations.It wasn't so easy. However we managed to perform our experiment on PI reducing our Input size to 320 from 416
And changing batch size from 4 to 8

Some of the results of the process are screenshot and kept here :



We obtained an inference time of about 2.7 FPS .The above picture explains the use of the Memory which is also all consumed by our process.

The screenshot shows a remote desktop connection to a system running Ubuntu. The terminal window displays the following information:

```
Tasks: 33, 54 lib; 4 running
Load average: 2.35 1.08 0.52
Uptime: 2 days, 20:25:39
627M / 63M
```

Below this, a table lists system resources and processes:

pid	user	ppid	uid	gid	st	ti	pr	ni	in	cs	ss	nm	cmd
61388	ubuntu	20	0	2563M	639M	239M	R	47.8	8.2	0:46.55	python	object_tracker.py	
61385	ubuntu	20	0	2563M	639M	239M	R	46.5	8.2	0:46.41	python	object_tracker.py	
61387	ubuntu	20	0	2563M	639M	239M	R	45.2	8.2	0:46.45	python	object_tracker.py	
61386	ubuntu	20	0	2563M	639M	239M	R	44.6	8.2	0:46.40	python	object_tracker.py	
61391	ubuntu	20	0	2563M	639M	239M	S	3.9	8.2	0:02.89	python	object_tracker.py	
61392	ubuntu	20	0	2563M	639M	239M	S	3.2	8.2	0:02.54	python	object_tracker.py	
61389	ubuntu	20	0	2563M	639M	239M	S	2.0	9.2	0:02.55	python	object_tracker.py	
61049	ubuntu	20	0	7896	3684	2292	R	1.9	0.0	0:07.54	htop		
61399	ubuntu	20	0	2563M	639M	239M	S	1.3	8.2	0:01.14	python	object_tracker.py	
61390	ubuntu	20	0	2563M	639M	239M	S	0.6	8.2	0:02.23	python	object_tracker.py	
61397	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:01.16	python	object_tracker.py	
61398	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:01.14	python	object_tracker.py	
60985	ubuntu	20	0	15620	4364	3128	S	0.0	0.1	0:00.40	sshd	ubuntu@pts/0	
61401	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:00.07	python	object_tracker.py	
56031	ubuntu	20	0	15620	4404	3168	S	0.0	0.1	0:04.87	sshd	ubuntu@pts/0	
1718	root	20	0	80940	3664	7708	S	0.0	0.0	0:40.97	/usr/sbin/irqbalance	--foreground	
1536	root	RT	0	273M	16084	6300	S	0.0	0.2	0:24.62	/sbin/multipathd	-d -s	
1540	root	RT	0	273M	16084	6300	S	0.0	0.2	0:11.82	/sbin/multipathd	-d -s	
1	root	20	0	372M	6110	7200	S	0.0	0.2	0:56.40	/sbin/init	flirtie splash	
61381	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:00.07	python	object_tracker.py	
61382	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:00.07	python	object_tracker.py	
61383	ubuntu	20	0	2563M	639M	239M	S	0.0	8.2	0:00.07	python	object_tracker.py	
7280	root	20	0	1120M	28396	15432	S	0.0	0.4	0:45.32	/usr/lib/snapd/snapd		
7317	root	20	0	1120M	28396	15432	S	0.0	0.4	0:03.60	/usr/lib/snapd/snapd		
7418	root	20	0	1120M	28396	15432	S	0.0	0.4	0:04.89	/usr/lib/snapd/snapd		
30500	root	20	0	1120M	28396	15432	S	0.0	0.4	0:03.62	/usr/lib/snapd/snapd		
1722	root	20	0	231M	5588	5764	S	0.0	0.1	0:08.66	/usr/lib/accountsservice/accounts-daemon		
7385	root	20	0	1120M	28396	15432	S	0.0	0.4	0:05.79	/usr/lib/snapd/snapd		
7312	root	20	0	1120M	28396	15432	S	0.0	0.4	0:03.93	/usr/lib/snapd/snapd		
1714	root	20	0	231M	5588	5764	S	0.0	0.1	0:08.99	/usr/lib/accountsservice/accounts-daemon		
1537	root	RT	0	273M	16084	6300	S	0.0	0.2	0:02.87	/sbin/multipathd	-d -s	
984	root	19	0	67656	12220	11896	S	0.0	0.2	0:02.20	/lib/systemd/systemd-journald		
929	root	20	0	20248	4916	3404	S	0.0	0.1	0:01.43	/lib/systemd/systemd-udev		
1538	root	RT	0	273M	16084	6300	S	0.0	0.2	0:00.00	/sbin/multipathd	-d -s	
1539	root	RT	0	273M	16084	6300	S	0.0	0.2	0:00.45	/sbin/multipathd	-d -s	
1541	root	RT	0	273M	16084	6300	S	0.0	0.2	0:00.00	/sbin/multipathd	-d -s	
1542	root	RT	0	273M	16084	6300	S	0.0	0.2	0:00.00	/sbin/multipathd	-d -s	
1578	systemd-t	20	0	80152	4508	5624	S	0.0	0.1	0:00.02	/lib/systemd/systemd-timesyncd		
1571	systemd-t	20	0	80152	4508	5624	S	0.0	0.1	0:01.15	/lib/systemd/systemd-timesyncd		
1028	systemd-n	20	0	26352	7080	5988	S	0.0	0.1	0:01.71	/lib/systemd/systemd-networkd		
1630	systemd-r	20	0	24340	13524	8192	S	0.0	0.2	0:01.84	/lib/systemd/systemd-resolved		
1751	root	20	0	231M	5588	5764	S	0.0	0.1	0:00.07	/usr/lib/accountsservice/accounts-daemon		
1715	messagebus	20	0	8484	4492	3400	S	0.0	0.1	0:06.54	/usr/bin/dbus-daemon	--system --address=systemd: --nofork --nopidfile --system-activation --syslog-only	

3.JETSON NANO : The most challenging task was installing the dependencies and facing the challenge of installing the correct version . We faced a lot of version issues while installing tensorflow and inferencing on Jetson Nano. One of the main issue was OOM(out of memory).Our Nano couldn't accept tensor of large sizes which we fixed through resizing the input size and increasing the batch dimensions. The following outputs were obtained while inferencing on Jetson Nano .

We also tried running sort model in jetson Nano but some issues with that as well. Recently we tried inferencing from our mobilenet ssd Detections to our deep sort algorithm but that requires the frame captured from the edge device to be compatible to the Frames captured by the nano device.

We experienced the same graph definition issues in jetson nano from conversing from the heavy model to the light models.

We are still exploring the alternate algorithms that could be used for the object Tracking. Optical flow, mean shift, Open Cv library also includes many of such algorithms but those have various issues regarding:

1. Occlusion
2. Id switching

Till date the best algorithm in accordance to our subject is the SOTA Deep Sort algorithm . However we are exploring more on the context.