

Preliminary report #1 on Semester task WS23

NVIDIA NeMo model and runtime on minicomputers

Here is the description of the work I have done during my first month on the task

- 1) I have chosen one of the test sets of LibriSpeech ASR corpus, namely test-clean.tar.gz and test-other.tar.gz and written a python-code to unzip files and collect them into a set of samples to feed the model later
- 2) I have chosen one of the suggested STT En Citrinet 1024 model as a top one listed in the suggested models to use with NeMo-library
- 3) With nemo-toolkit I converted the model to the ONNX-model to later measure any changes in the model performance
- 4) Machines: one of the hardest part in my work was trying to create an environment to run my notebook for performance measuring on a variety of existing set of machines. Later in the report I will list all machines, actions taken, problems encountered and outcomes reached
- 5) Ready to use 3 Jupyter notebooks:
 - a. *ASR Librispeech Wav Audios Preparation* - unzipping a dataset in a *Librispeech* format and converting audio samples to a .wav-format suitable for NeMo-models
 - b. *ASR NeMo-ONNX Conversion* - converting an already downloaded .nemo-file to a .onnx-file with model to use in ONNX-runtime
 - c. *ASR Librispeech Wav Audios Traversal* - one cell with gathering audio samples from unzipped dataset of wav-s, one cell to feed these sample to an ONNX-runtime, the other cell to measure same samples in NeMo-runtime (test-other sample has been chosen)

Installing NeMo and ONNX-environments to use Citrinet model

NeMo library provides a mechanism to convert demo-model to onnx-model. I have used it successfully, although the inferencing would only produce word-piece-tokens and their probabilities. Therefore, I still needed to use some kind of preprocessing and decoding code before and after inferencing, using the classes from NeMo packages. Maintainers in NeMo-repository and forums suggest to use not the overall `nemo_toolkit[all]` package with model-training features, but the smaller one `nemo_toolkit[asr]` containing preprocessing and decoding classes.

Nevertheless, in order to find out the difference between original and converted models on machine I have tried to install all environments (ONNX, NeMo, NeMo-ASR).

For now I have got mixed results on trying to run initial/converted models on a range of devices (both minicomputers and regular PCs). I am listing all the test-cases I have tried and my conclusions in a chronological order:

Mainly I followed the guide from the [official NeMo guide](#), sometimes referring to issue-reports on specific device requirements in the [NeMo GitHub repository](#) and posts in [NVidia Developers Forum](#)

1. **Raspberry Pi 3 Model B** - No success

Multiple errors during installation, pip cannot find the proper version of wheels to install needed libraries multiple incompatibility issues

2. **Apple Silicon M1** - No success

Created a conda environment, installed all packages successfully, but the code produces errors. I didn't proceed to try, since these are not our target devices anyway

3. **Intel Core i7** - Success

Model path	Model size (Mb)	Peak size of memory blocks	Model load time (s)	Inference time (s)	Data set samples count	Dataset length (s)	RTF (s)	WER	MER	WIL	WIP	CER
models/stt_en_citrinet_1024.onnx	572.260998	586784744	0.35480785369873047	945.0373101234436	2939	19229.570124999973	0.0491	0.07080220850925625	0.0702812387400201	0.12485083142830666	0.8751491685716933	0.02867743567558055
models/stt_en_citrinet_1024.nemo	574.24896	25801311	4.89382791519165	5390.990655899048	2939	19229.570124999973	0.2803	0.07040100873087136	0.06988299102994444	0.12411761916397124	0.8758823808360288	0.028417131669831865

4. **Intel NUC 12 wshi3000** - Success (Not full table - I am going to run notebooks with more metrics again)

Model name	Samples count	Processing time (s)	WER
stt_en_citrinet_1024.onnx	2939	3391.3835377693176	0.0708786275146629
stt_en_citrinet_1024.nemo	2939	2918.6985421180725	0.07049653248762967

5. *Raspberry Pi 4 Model B - Partial Success*

ONNX model produce results, but NeMo model stops the notebook from running once

Model path	Model size (Mb)	Peak size of memory blocks	Model load time (s)	Inference time (s)	Dataset samples count	Data set length (s)	RTF (s)	WER	MER	WIL	WIP	CER
stt_en_citrinet_1024.onnx	572.260998	598059374	14.51574444770813	422.55275654792786	100	716.945000003	0.5894	0.07411630558722919	0.07336343115124154	0.12810456417258087	0.8718954358274191	0.023384049831512304

the inferencing starts - here is a quick snippet on a 100 samples, I am going to run all other samples later to gather more accurate *RTFs*, *WERs* etc.

6. *Nvidia Jetson Nano Kit - No success*

Unfortunately, the whole process of NeMo environment installation is very laborious, I didn't succeeded to use NeMo-toolkit on it. As mentioned in the [maintainer's post](#) they don't support Jetson Nano, but I have tried to install .[asr] part only and faced incompatibility issues as well

7. *Nvidia Jetson Xavier NX - No success*

The situation is similar to Jetson Nano. Possibly, the right way of implementing ASR on Jetson-devices is not python code, rather ready-to-use containers

Further work

In my further work I will try to further explore Nvidia NeMo model, decide how to not use nemo-library classes for preprocessing and decoding of inferencing results when using ONNX-model, try out suggested alternatives for running models on arm-based SBCs.