



# Lightning-Fast Modulation Classification with Hardware-Efficient Neural Networks

10/27/2021

The A(MC) Team:

Jakob Krzyston

PhD Student @ GT, Research Engineer, GTRI  
jakobk@gatech.edu

Dr. Rajib Bhattacharjee

Principal Engineer, DeepSig Inc  
raj@deepsig.ai

Dr. Andrew Stark

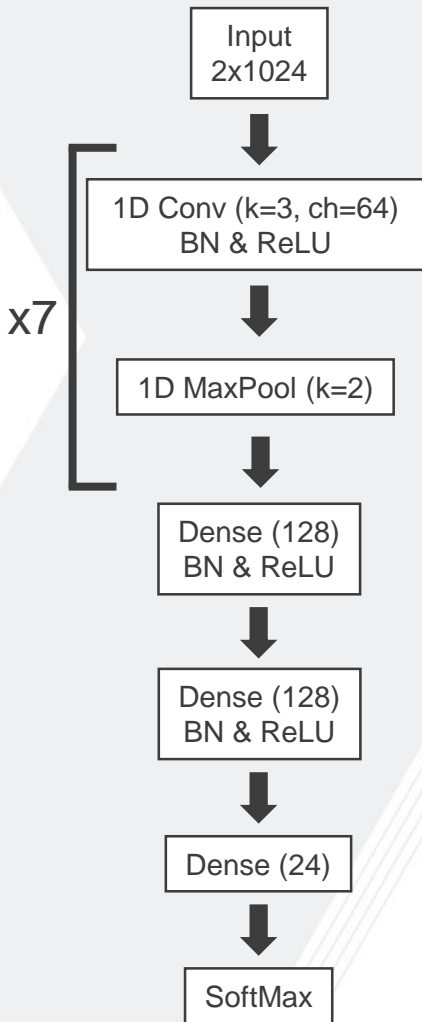
Senior Research Engineer, GTRI  
andy.stark@gtri.gatech.edu

# Possible Approaches

- ~~Modify the provided architecture [Speed]~~
- **Reduce the quantization with Brevitas [Speed]**
  - Four bits for both weights and activations
- **Prune weights [Speed]**
  - L1 unstructured Iterative Magnitude Pruning (IMP)
  - Prune when accuracy threshold reached
- **Adjust training paradigm [Accuracy]**
  - Learning Rate Scheduler → Reduce LR on Plateau

# Methods

## Architecture



## IMP (Simplified)

```
for num_prune_iterations:
    for num_epochs:
        train model
        test model
        if model_accuracy > 0.56:
            save model weights
            prune 20% of weights
            break
        else:
            lr_scheduler.step
```



Alex Renda  
@alex\_renda\_

- 1) Train to completion.
- 2) Globally prune the 20% of weights with the lowest magnitudes.
- 3) Retrain with learning rate rewinding for the original training time.
- 4) Iteratively repeat steps 2 and 3 until the desired sparsity is reached.

That's it.

11:03 AM · Mar 10, 2020 · Twitter Web App

6 Retweets 2 Quote Tweets 36 Likes

## Compression Summary

Quantity	Original*	Final
Bit Ops	807,699,904	24,436,576
Weight Bits	1,244,936	68,072
Compression	1x	9.313x
Sparsity	0%	89.26%

\*Values from provided code

## Notes

- Did not reset LR scheduler
  - May have reduced our end performance
- Sparsity % =  $1 - (0.8 \wedge 10)$
- Compression =  $1 / (0.8 \wedge 10)$

# Final Results

## Inference Cost Score:

- 0.042467

## Overall Test Accuracy:

- 0.5625

