



4-bit Conformer with Native Quantization Aware Training for Speech Recognition

Shaojin Ding, Phoenix Meadowlark, Yanzhang He, Lukasz Lew, Shivani Agrawal, Oleg Rybakov

Google LLC, USA

{shaojinding, meadowlark, yanzhanghe, lew, shivaniagraval, rybakov}@google.com

Abstract

Reducing the latency and model size has always been a significant research problem for live Automatic Speech Recognition (ASR) application scenarios. Along this direction, model quantization has become an increasingly popular approach to compress neural networks and reduce computation cost. Most of the existing practical ASR systems apply post-training 8-bit quantization. To achieve a higher compression rate without introducing additional performance regression, in this study, we propose to develop 4-bit ASR models with native quantization aware training, which leverages native integer operations to effectively optimize both training and inference. We conducted two experiments on state-of-the-art Conformer-based ASR models to evaluate our proposed quantization technique. First, we explored the impact of different precisions for both weight and activation quantization on the LibriSpeech dataset, and obtained a lossless 4-bit Conformer model with 7.7x size reduction compared to the float32 model. Following this, we for the first time investigated and revealed the viability of 4-bit quantization on a practical ASR system that is trained with large-scale datasets, and produced a lossless Conformer ASR model with mixed 4-bit and 8-bit weights that has 5x size reduction compared to the float32 model.

Index Terms: speech recognition, model quantization, 4-bit quantization

1. Introduction

With the fast growth of voice search and speech-interactive features, automatic speech recognition (ASR) [1, 2, 3, 4, 5] has become an essential component for user-interactive services and devices (e.g., search by voice functions in search engines and smartphones) over the years. Modern ASR applications are mostly developed based on an end-to-end model [6, 7, 8, 9, 10, 11], which has been shown to achieve significant recognition performance improvements compared to conventional hybrid systems [12] with a much smaller model size.

Improving latency and model size without compromising recognition quality has been an active research topic for years, as they benefit live ASR applications with both server-side and on-device models. Prior studies have explored the use of network pruning [13, 14, 15, 16], knowledge distillation [17], and model quantization [18, 19, 7, 20]. Among model quantization methods, post training quantization (PTQ) with int8 [21] is popular and easy to use (e.g. by just setting a flag during model conversion in TFLite [21]). It is successfully applied in multiple applications [7, 20]. One of the drawbacks of such technique is the potential performance degradation due to the loss of precision. Another limitation of PTQ is the lack of control over model quantization, e.g. it does not support int4 quantization or customized quantization of selected set of layers. That is why in this work we are focused on quantization aware training (QAT).

QAT has several flavors: “fake” [22] and native, discussed in [23] also called Accurate Quantized Training. It is called accurate because it uses native integer operations for executing quantized operations (e.g. matrix multiplications) and does not have any difference between accuracy during training and inference. Whereas “fake quantization” can have a numerical difference between training (with float operations) and inference (with integer operations) modes if float operation does not fit into 23 bits of mantissa during training.

In this work we use native QAT for ASR model, because it follows the approach of “what you train is what you serve”. It can use native integer operations during training (if hardware supports it or else it will use int emulated in float32, A.k.a. “fake quantization”). With native integer operations, there is no numerical difference between forward propagation of training and inference. We conduct extensive experiments on LibriSpeech and large-scale Voice Search datasets with a state-of-the-art Conformer ASR model [24] to evaluate this quantization approach, and systematically analyze how different quantization precisions affect the compute cost-accuracy tradeoff with the main focus on int4¹ quantization. Our main contributions are outlined as below:

- We leverage native QAT approach for ASR model, instead of running “fake quantization” that is mostly used in previous studies [25, 26, 27]. It allows us to run the model on both cloud (e.g. on TPU) and mobile applications, whereas “fake” QAT is used mostly for mobile applications and needs special conversion done by TFLite [21].
- We minimize the number of operations used for quantization; so that training time with native QAT is increased only by 7% (on TPU that supports float operations only) in comparison to training of a float model. If native QAT is executed on hardware that supports integer operations, then we expect more speed up in training time.
- We demonstrate that *Large* Conformer ASR model on LibriSpeech data has minimal or no accuracy loss with int4 weights and float32 activations quantization, or int4 weights and int8 activations quantization (as expected it also works well with int8 weights and activation quantization). For the first time we evaluate native 4-bit QAT on a real production model trained on large scale data. Despite of no accuracy regression on public Librispeech data, we observe that native QAT introduces regression on the model trained with large data sets. More importantly, we investigate several strategies to mitigate the regression and affect the cost-accuracy trade off, bringing new insights to both model quantization and ASR modeling studies.

Relation to prior work. A number of previous studies have also explored the idea of ASR model quantization [25, 26, 27].

¹We interchangeably use 4-bit and int4 (8-bit and int8) hereafter.

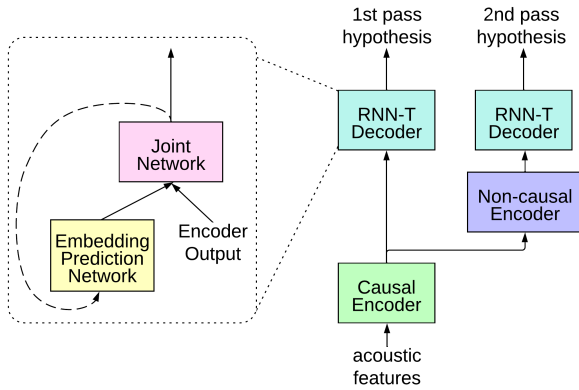


Figure 1: The architecture for the multi-pass ASR model.

Our work differs from them in several aspects. First and foremost, most previous studies conduct experiments only on relatively small datasets (e.g., LibriSpeech, Switchboard, and Call-Home). The models trained on these datasets are usually heavily over-parameterized, which makes it easy to obtain a loss-less 8/4-bit quantized model. By contrast, we for the first time examine 4-bit ASR quantization on a model trained on a large scale dataset ($\sim 400k$ hours) and showed the behavior difference. Second, all three studies use “fake quantization” that can have a numerical difference between training and inference. However, in this study, we consider native QAT, which has several advantages: 1) the same model representation can be used for both mobile and cloud model training and inference; 2) if hardware supports integer operations then we can expect model training speed up. Lastly, [26, 27] only explore 8- or 6-bit quantizations. Although [25] also explored 4-bit quantization that is similar to ours, there are several differences: 1) we focus on a recent state-of-the-art Conformer-based model, but [25] uses LSTM ASR; 2) we train model with native QAT from scratch, whereas in [25] they first train float model and then fine-tune it with QAT, which complicates training pipeline and increases total training time; 3) to minimize the impact of native QAT on total training time we reduced the number of operations in quantization function and do not have clamp with learnable boundaries as in [25].

2. Method

2.1. ASR model architectures

We use the state-of-the-art Conformer Transducer [24] backbones in this study. We considered slightly different architectures for the LibriSpeech experiments and the production experiments based on a real-world on-device ASR system to guarantee fair comparisons to the corresponding baselines. For LibriSpeech, we follow the architecture proposed in [24], which has a single encoder with different layers for *Small* (10M parameters) and *Large* (118M) models. The decoder is a standard RNN-Transducer decoder with 1-layer LSTM.

Alternatively, for model in practical applications, we follow the multi-pass architecture proposed in [28] and shown on Figure 1. It has a first pass causal encoder (47M parameters), followed by a second pass non-causal encoder (60M). This multi-pass model unifies the streaming and non-streaming ASRs, where the causal encoder uses only left context and produces partial results with minimal latency, and the non-causal can provide more accurate hypothesis by using both left and

```

1 def quantize(input, axis):
2     max_val = tf.math.reduce_max(tf.math.abs(
3         input), axis=axis, keepdims=True)
4     scale = tf.divide(max_val, 127.0)
5     scale = tf.stop_gradient(scale)
6     scaled_input = tf.math.divide_no_nan(input,
7         scale)
8     return tf.cast(
9         tf.math.floor(scaled_input + 0.5), tf.
10        int8)

```

Figure 2: int8 native quantization in TensorFlow.

right context. Unlike the original architecture, we use a separated decoder (4.4M) for each encoder. The decoder has an embedding prediction network [29] and a 1-layer fully-connected joint network. The embedding prediction network further reduces the decoder model size without accuracy degradation.

2.2. Native quantization aware training

Standard approach for QAT in TensorFlow(TF) [30] is based on “fake” QAT and relies on `tf.quantization.fake_quant_*` operations. In this case, researchers have to use these operations during training and server-side inference. For on-device models, they use TFLite for converting fake quantization operations to integer operations and then run output model with TFLite on mobile phone. Even though this approach provides end-to-end user experience, it also has some disadvantages: 1) It requires additional conversion step of fake quantization operations to integer. 2) Existing API `tf.quantization.fake_quant_with_min_max_vars_per_channel` supports per-channel min and max values estimation over the last dimension only. However, there are use cases when channel dimension is not the last one. In this case we have to permute dimensions of the input tensor to make channel dimension the last one, then we “fake” quantize the tensor with `tf.quantization.fake_quant_with_min_max_vars_per_channel` and after that permute dimensions back to the original order of the input tensor. These additional permutation operations can increase training time in comparison to the non-quantized model training time. To address all above we propose to use native QAT and with native tf operations, also discussed in [23]; called Accurate Quantized Training. It allows us to follow an approach called “what you train is what you serve”. As a result we can use the same TF model for training and inference on both mobile and cloud TPU applications with opportunity to speed up model training if hardware supports integer operations (e.g. matrix multiplications).

One of the methods of 4-bit ASR model quantization is based on QAT fine-tuning [25]. It has several steps: train float model, then use it for QAT fine tuning. This approach always increases total training time (it includes float model training and QAT fine tuning). By contrast, we train our model from scratch in quantization aware mode and show that training time is increased by only 7% (on hardware with float operations). We achieve it by using native QAT based on dynamic quantization shown in Figure 2 (for int8 use case). It estimates max values over axis which will be quantized (it supports channel-wise quantization). Then it computes *scale* by dividing max values with 127 for int8 (or by 7 for int4). After that, the input tensor is quantized by dividing with *scale* and casting it to int. De-quantization will be done by multiplying a tensor with *scale*. To reduce overall computations we do not use “zero point” as in [22] and assume that input tensor values distribution is sym-

metrical. Even though it is a strong assumption, in Section 4.1, we empirically show that model weights can be quantized by native 4-bit QAT without accuracy degradation.

3. Experimental setups

3.1. Datasets

We conducted experiments on both LibriSpeech and an internal large-scale datasets. LibriSpeech training set contains 960 hours of speech, where 460 hours of them are “clean” speech and the other 500 hours are “noisy” speech. The testing set also consists of a “clean” and a “noisy” subset.

When running experiments with the large-scale datasets, we train the models with a training set [31, 20] consisting of $\sim 400k$ hours English audio-text pairs from multiple domains, such as YouTube and anonymized voice search traffic. YouTube were transcribed in a semi-supervised fashion [32]. All other domains are anonymized and hand-transcribed. During evaluations, we use the Voice Search (VS) test set that contains around 12k voice search utterances, each having an average length of 5.5 seconds. Our data handling abides by *Google AI Principles* [33].

3.2. Conformer model with Librispeech data

The LibriSpeech conformer model [24] uses a frontend of 80-dimensional log Mel-filterbank energies, extracted from 25ms window and 10ms shift. The *Small(S)* and *Large(L)* variants have 16, 17 layers, with a dimensionality of 144, 512, respectively. The *Small* model has 4 attention heads in self-attention layers, while the *Large* model has 8 heads. The kernel size of the depthwise convolutions is set to 32. The LSTM layer in decoder has 640 units.

3.3. Conformer model with large-scale data

The model with large-scale data uses a 128-dimensional log Mel-filterbank energies as the frontend feature, in which the 4 contiguous frames are stacked, and the stacked sequence is sub-sampled by a factor of 3. We use causal convolution for all layers with a kernel size of 15. The causal conformer encoder has 7 conformer layers (first 3 layers have no self-attention), which has 23-frame left context per layer and no right context to strictly prevent the model from using future inputs. The non-causal encoder 6 conformer layers, with additional 30-frame right context across 6 layers that processes 900ms speech from the future. All the self-attention layers have 8 heads. Each separate RNN-T decoder is comprised of an 320-dimensional embedding prediction network and a 384-dimensional fully-connected joint network.

Evaluations of production models are running on an on-device inference pipeline, where we convert the TensorFlow graphs to TFLite format, with the corresponding quantization approach. Additionally, we did not use any language model in our experiments, as this is orthogonal to the end-to-end model performance. When applying native QAT to the backbone models, we only quantize the encoders since the decoders are relatively small. We do not quantize convolutional kernels for the same reason.

4. Results

We conduct two sets of experiments to evaluate the proposed native QAT on ASR models. First, we examine different quan-

Table 1: Results of our proposed int8/4 QAT on Conformer Large(L) and Small(S) models with the baseline approach [27] on LibriSpeech test-clean and test-other subsets. Please see Section 4.1 for the meanings of the method abbreviations.

Conformer (L)			
Method	test-clean	test-other	Model size (MB)
Float	2.0	4.4	472
I8W	2.0	4.5	118
I4W	2.0	4.4	61
I8WA	2.0	4.5	118
I4WI8A	2.1	4.4	61
I4WA	3.1	8.2	61
FakeI4W	2.0	4.6	61
Conformer (S)			
Float	2.5	6.1	40
I8W	2.5	6.0	10
I4W	2.7	6.3	6
I8WA	2.5	6.0	10
I4WI8A	2.8	6.6	6
I4WA	5.0	12.1	6
FakeI4W	2.9	6.9	6
Baselines			
[34]I8WA	6.9	N/A	8
[35]I8W	8.7	22.3	60
[35]I6W	8.9	22.8	45
[27]I8W	2.7	6.9	123
[27]I6W8A	3.6	8.2	92

tization precision and compared with five baselines on LibriSpeech dataset, demonstrating the validity and the effectiveness of native QAT in conformer based ASR models. Second, we explore applying native QAT to a practical ASR models that is trained on large-scale data and the corresponding optimizations.

4.1. Experiments on LibriSpeech

We experiment with conformer *Large* and *Small* models to examine the behaviors of QAT with different model sizes. In terms of QAT, we consider 7 quantization configurations to understand the impacts of 8/4-bit quantization on weights only or both weights and activations:

- *Float*: float32 weight, float32 activation (baseline)
- *I8W*: int8 weight, float32 activation
- *I4W*: int4 weight, float32 activation
- *I8WA*: int8 weight, int8 activation
- *I4WI8A*: int4 weight, int8 activation
- *I4WA*: int4 weight, int4 activation
- *FakeI4W*: fake int4 weights, float32 activation

Quantizing the weights alone can reduce the model size, but we still need to convert the weights to float during inference. By contrast, quantizing the activations will benefit in both model size and run-time efficiency, with integer matrix multiplication.

We train models with different combinations of weight and activation quantization using native QAT and “fake” QAT approaches, and reported their WERs in Table 1. Here, we first analyze the impact of quantizing only model weights to 8/4-bit. In 8-bit cases (*I8W*), we see no regression in either *Large* or

Table 2: Results of applying int8/4 QAT to production ASR model. We use the actual TensorFlow Lite file size of the models with corresponding quantization approaches for model size, measure in megabyte (MB). PTQ refers to post-training quantization. From model E4 to E13, we apply in4 QAT to the listed layers and int8 PTQ for the remaining layers.

Exp	Model	Voice Search WER		Model size (MB)
		1st pass	2nd pass	
B0	float32 model	8.0	5.7	460
E0	int8 PTQ	7.9	5.8	115
E1	int8 QAT	7.9	5.8	115
E2	int4 PTQ	>100.0	>100.0	61
E3	int4 QAT	8.5	6.2	61
E4	int4 QAT causal	8.5	5.8	91
E5	int4 QAT non-causal	7.9	6.2	84
E6	int4 except first& last	8.5	6.1	78
E7	int4 except self-atten	8.1	5.9	94
Quantizing different number of layers in both encoders				
E8	int4 first layer	7.9	5.8	108
E9	int4 first 2 layer	7.9	5.8	101
E10	int4 first 3 layer	7.9	5.9	94
E11	int4 first 4 layer	8.0	6.1	84
E12	int4 first 5 layer	8.1	6.2	75
E13	int4 first 6 layer	8.5	6.2	66

Small model. However, in 4-bit cases (*I4W*), although the *Large* model can still retain the float performance, the *Small* model has introduced 0.2 WER increase compared to the baseline due to the limited capacity of the *Small* model. Additionally, we also investigate the possibility of quantizing both weights and activations. With *I8WA*, the *Large* and *Small* models can still preserve the float performance. As we quantized the model more aggressively with *I4WI8A* and *I4WA*, we see more serious performance loss. To summarize, the most light-weight configurations without WER loss of *Large* model is to have int4 weights and int8 activations. For *Small* model, the most light-weight configuration without WER loss is to use int8 for both weights and activations. Compared against the baseline models, our most light-weight *I4WI8A* Conformer (S) model (6MB) has already significantly outperformed all the baselines while having more aggressive quantization scheme, demonstrating our state-of-the-art performance.

“Fake” QAT with 4-bit weights has the same accuracy as the float model on *test-clean* and little accuracy reduction on *test-other* as shown on Table 1. We benchmarked “fake” QAT *Fake4W* and native QAT *I4W* with [36] on TPU and observed that training with native QAT is 7% slower than float model. We explain it by additional operations shown in Figure 2 and by the fact that we used hardware with float operations support only (we expect speedup on hardware with integer operations support). We also observed that training with *Fake4W* is 6% slower than native QAT.

4.2. Exploring the limit of 4-bit quantization on large-scale data

Although we can easily obtain int4 quantized model that has no performance regression in LibriSpeech, we still see performance drops when applying native int4 QAT to the model trained with large-scale data, even only with weight quantization. Results are shown in Table 2. Comparing to float32 models (*B0*), the int8 models (*E0* and *E1*) have similar results. Unless otherwise indicated, we compare the following 4-bit models with *E0*, since int8 PTQ has been widely used in real-world ap-

plications. For int4 models, with PTQ (*E2*), the model cannot produce a reasonable performance anymore. By contrast, native QAT (*E3*) has additional 0.6/0.4 degradations compared to int8 models, which again confirms the effectiveness of native QAT. The reason of the different conclusions between LibriSpeech and production models could be that LibriSpeech models are highly over-parameterized (LibriSpeech training set: 960 hours; large-scale training set: 400k hours), and therefore, the models still have enough capacities with 4-bit weights. By contrast, the production model is trained with an extremely large amount of data, so the model is not overfitting anymore. This can also be verified through the int8 results – when applying native QAT to int8 model, we did not see any performance improvement over PTQ.

Additionally, with the multi-pass encoder architecture, we explore only quantizing either the first pass causal encoder or the second pass non-causal encoder. We keep other layers to be int8 PTQ hereafter. As shown in Table 2 (*E4* and *E5*), this keeps the performance of the non-quantized pass unchanged, which can be an intermediate solution for the scenarios where either the first or second pass is more important (e.g., we are more tolerant to regressions on the first pass if only using it for partial hypothesis).

We further investigate the possibility of having no performance loss in either pass with three strategies. First, we keep the first and the last layer to be int8, as they tend to be most sensitive to quantization [23, 37]. However, as shown in Table 2, this observation does not hold for our model (*E6*), and the performance remains to be the same as quantizing the entire causal encoder to int4. Second, we keep the self-attention layers to be int8 since they are the most essential layers in conformer. With this, we obtain a model (*E7*) with 0.2/0.1 WER regression for the two passes compared to the int8 model. Lastly, we experiment with quantizing different number of layers from bottom to the top in both passes (*E8* to *E13*). When quantizing the first 3 layers of both encoders (*E10*), we only see 0.1 performance loss on the 2nd pass. When quantizing additional layers, the model starts to have more serious degradations on both passes. These results indicate that lower layers are more robust to quantization, as the upper layers can mitigate the performance loss. It is also worthwhile to experiment with quantizing different number of layers in the two decoders, which we will continue investigating in future work.

Consequently, we find that model *E10* has the best trade-offs, mostly retaining the performance of int8 or float32 models. In terms of model size, our int4 model has 81% of int8 model size and only 20% of the original float32 model size, which can significantly reduce memory, latency, and power consumption for on-device and server applications.

5. Conclusions

In this paper, we proposed a novel approach based on native QAT to establish state-of-the-art 4-bit quantized ASR model. Through an experiment on LibriSpeech, we analyzed the impact of different quantization configurations for both weights and activations, validating the effectiveness of QAT in building 4-bit quantized ASR models. Following this, we for the first time applied it to a production ASR with large-scale data. Observing a performance regression on production ASR models, we proposed three strategies to alleviate the regression and thoroughly discussed the cost-accuracy tradeoffs.

6. References

- [1] D. Wang, X. Wang, and S. Lv, "An overview of end-to-end automatic speech recognition," *Symmetry*, vol. 11, no. 8, p. 1018, 2019.
- [2] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Sathesh, S. Sengupta, A. Coates *et al.*, "Deep speech: Scaling up end-to-end speech recognition," *arXiv preprint arXiv:1412.5567*, 2014.
- [3] A. Graves, "Sequence transduction with recurrent neural networks," *arXiv preprint arXiv:1211.3711*, 2012.
- [4] J. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in *ICONIP*, 2015, pp. 577–585.
- [5] L. Dong, S. Xu, and B. Xu, "Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5884–5888.
- [6] J. Li, Y. Wu, Y. Gaur *et al.*, "On the Comparison of Popular End-to-End Models for Large Scale Speech Recognition," in *Proc. Interspeech*, 2020.
- [7] Y. He, T. N. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, D. Zhao, D. Rybach, A. Kannan, Y. Wu, R. Pang *et al.*, "Streaming end-to-end speech recognition for mobile devices," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6381–6385.
- [8] C.-C. Chiu, T. N. Sainath, Y. Wu *et al.*, "State-of-the-art Speech Recognition With Sequence-to-Sequence Models," in *Proc. ICASSP*, 2018.
- [9] S. Kim, T. Hori, and S. Watanabe, "Joint CTC-attention based end-to-end speech recognition using multi-task learning," in *Proc. ICASSP*, 2017.
- [10] J. Li, R. Zhao, H. Hu, and Y. Gong, "Improving RNN transducer modeling for end-to-end speech recognition," in *Proc. ASRU*, 2019.
- [11] A. Zeyer, A. Merboldt, R. Schlüter, and H. Ney, "A new training pipeline for an improved neural transducer," in *Proc. Interspeech*, 2020.
- [12] G. Pundak and T. N. Sainath, "Lower frame rate neural network acoustic models," in *Proc. Interspeech*, 2016.
- [13] R. Takeda, K. Nakadaï, and K. Komatani, "Node pruning based on entropy of weights and node activity for small-footprint acoustic model based on deep neural networks," in *INTERSPEECH*, 2017, pp. 1636–1640.
- [14] Y. Shangguan, J. Li, Q. Liang, R. Alvarez, and I. McGraw, "Optimizing speech recognition for the edge," *arXiv preprint arXiv:1909.12408*, 2019.
- [15] D. Gao, X. He, Z. Zhou, Y. Tong, K. Xu, and L. Thiele, "Rethinking pruning for accelerating deep inference at the edge," in *SIGKDD*, 2020, pp. 155–164.
- [16] S. Ding, T. Chen, and Z. Wang, "Audio lottery: Speech recognition made ultra-lightweight, noise-robust, and transferable," in *International Conference on Learning Representations*, 2021.
- [17] C. Li, L. Zhu, S. Xu, P. Gao, and B. Xu, "Compression of acoustic model via knowledge distillation and pruning," in *ICPR*. IEEE, 2018, pp. 2785–2790.
- [18] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [19] R. Alvarez, R. Prabhavalkar, and A. Bakhtin, "On the efficient representation and execution of deep acoustic models," *Interspeech 2016*, pp. 2746–2750, 2016.
- [20] T. N. Sainath, Y. He, B. Li, A. Narayanan, R. Pang, A. Bruguier, S.-y. Chang, W. Li, R. Alvarez, Z. Chen *et al.*, "A streaming on-device end-to-end model surpassing server-side conventional model quality and latency," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6059–6063.
- [21] [Online]. Available: https://www.tensorflow.org/lite/performance/post_training_quantization
- [22] B. Jacob, S. Kligys, B. Chen, M. Zhu, M. Tang, A. G. Howard, H. Adam, and D. Kalenichenko, "Quantization and training of neural networks for efficient integer-arithmetic-only inference," *CoRR*, vol. abs/1712.05877, 2017. [Online]. Available: <http://arxiv.org/abs/1712.05877>
- [23] A. Abdolrashidi, L. Wang, S. Agrawal, J. Malmaud, O. Rybakov, C. Lechner, and L. Lew, "Pareto-optimal quantized resnet is mostly 4-bit," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3091–3099.
- [24] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," *Proc. Interspeech 2020*, pp. 5036–5040, 2020.
- [25] A. Fasoli, C.-Y. Chen, M. Serrano, X. Sun, N. Wang, S. Venkataramani, G. Saon, X. Cui, B. Kingsbury, W. Zhang *et al.*, "4-bit quantization of lstm-based speech recognition models," *arXiv preprint arXiv:2108.12074*, 2021.
- [26] A. Bie, B. Venkitesh, J. Monteiro, M. Haidar, M. Rezagholizadeh *et al.*, "A simplified fully quantized transformer for end-to-end speech recognition," *arXiv preprint arXiv:1911.03604*, 2019.
- [27] S. Kim, A. Gholami, Z. Yao, N. Lee, P. Wang, A. Nrusimha, B. Zhai, T. Gao, M. W. Mahoney, and K. Keutzer, "Integer-only zero-shot quantization for efficient speech recognition," *arXiv preprint arXiv:2103.16827*, 2021.
- [28] T. N. Sainath, Y. He, A. Narayanan *et al.*, "An Efficient Streaming Non-Recurrent On-Device End-to-End Model with Improvements to Rare-Word Modeling," in *Proc. of Interspeech*, 2021.
- [29] R. Botros, T. Sainath, R. David, E. Guzman, W. Li, and Y. He, "Tied & reduced rnn-t decoder," in *Proc. Interspeech*, 2021.
- [30] [Online]. Available: <https://www.tensorflow.org/>
- [31] A. Narayanan, R. Prabhavalkar, C.-C. Chiu, D. Rybach, T. N. Sainath, and T. Strohmaier, "Recognizing long-form speech using streaming end-to-end models," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 920–927.
- [32] H. Liao, E. McDermott, and A. Senior, "Large Scale Deep Neural Network Acoustic Modeling with Semi-supervised Training Data for YouTube Video Transcription," in *Proc. ASRU*, 2013.
- [33] Google, "Artificial Intelligence at Google: Our Principles." [Online]. Available: <https://ai.google/principles/>
- [34] A. Prasad, P. Motlicek, and S. Madikeri, "Quantization of acoustic model parameters in automatic speech recognition framework," *arXiv preprint arXiv:2006.09054*, 2020.
- [35] H. D. Nguyen, A. Alexandridis, and A. Mouchtaris, "Quantization aware training with absolute-cosine regularization for automatic speech recognition," in *Interspeech*, 2020, pp. 3366–3370.
- [36] J. Shen, P. Nguyen, Y. Wu, Z. Chen *et al.*, "Lingvo: a modular and scalable framework for sequence-to-sequence modeling," 2019.
- [37] J. Choi, Z. Wang, S. Venkataramani, P. I.-J. Chuang, V. Srinivasan, and K. Gopalakrishnan, "Pact: Parameterized clipping activation for quantized neural networks," *arXiv preprint arXiv:1805.06085*, 2018.