

---

# Quant-Noisier: Second-Order Quantization Noise

Model Compression

---

Sergio Charles, William Ellsworth, Lyron Co Ting Keh

Department of Computer Science, Stanford University

{scharles1, willells, lyronctk}@stanford.edu

## Abstract

Applying quantization to a neural network after training often leads to severe performance regressions. Quantization Aware Training (QAT), proposed by Jacob et al., addresses this problem by applying simulated training-time quantization for the model to learn robustness to inference-time quantization. One key drawback of this approach is that quantization functions induce biased gradient flow through the network during back-propagation, thus preventing the network from best fitting to the learning task. Fan et al. addressed this issue by proposing *Quant-Noise*, in which they apply simulated quantization to a fixed proportion, called the "quantization noise rate", of parameters during training rather than all of them which sets a new state-of-the-art for quantization. Our study, *Quant-Noisier*, builds upon their technique by exploring a variable noise rate instead of a fixed one, which we term second-order noise. We craft three candidate functions to vary noise rate during training and evaluate the variants with 3 datasets and 3 quantization schemes for each dataset. Surprisingly, we observe negative results on our hand-crafted candidate functions, and somewhat positive results on our ablation study of randomly varying the noise rate during training. Future work should determine whether these results hold for more datasets and quantization schemes, as well as investigating other schemes for varying the noise rate during training.

## 1 Introduction

Modern deep learning architectures are getting larger, with recent models spanning many billions of parameters [1]. By contrast, many NLP and SLP use-cases involve deployment to embedded devices (e.g. voice assistants, IoT devices), where massive models are prohibitively memory intensive and computationally expensive.

Model compression techniques aim to address this limitation by creating compact representations of models that can attain the same level of performance as their larger counterparts. For instance, pruning methods do so by removing extraneous weights and activations to produce sparse architectures with reduced parameter counts [2]. Our work is an investigation into another such compression technique: quantization. Rather than removing entire units as pruning does, scalar quantization (used in this work) shrinks networks by reducing the bit-widths of their parameters [3]. In addition to a reduced memory footprint, computations can, likewise, be sped up with the appropriate hardware accelerators [4]. Throughout this paper, we will refer to scalar quantization to  $n$ -bit integers as  $\text{int}_n$  quantization, e.g. " $\text{int}_8$  quantization."

We aim to improve upon Quant-Noise, a state-of-the-art (SotA) quantization method developed by Facebook AI Research (FAIR) [3]. Quant-Noise falls under the umbrella of methods that simulate quantization at train-time, so the model learns "robustness" to quantization. A core challenge to training-time quantization is that quantization functions are non-differentiable. The functions

necessary to approximate gradients for these non-differentiable operations, such as the straight through estimator [5], inherently produce biased gradients during backpropagation. Quant-Noise [3] addresses this limitation by quantizing only a random subset of the total network during each training step, allowing some unbiased gradients to flow while still teaching the model to be resilient to the distortion. Note that Quant-Noise selects this subset using a Bernoulli trial for each parameter, where the probability input (labelled the *noise rate*) is constant across all examples and epochs.

We suggest that using a varying noise rate can lead to further robustness to quantization. We term this "second-order noise," since we are adding noise to the noise rate itself.

We propose several novel formulations where the quantization noise rate changes throughout training. We evaluate these methods on three tasks (RTE, MRPC, and HarperValleyBank), with three quantization schemes apiece (`int8`, `int4`, and `int1`). We use three random seeds for most of our studies; counting each (task, method, quantization scheme, seed) tuple as a single experiment, we perform a total of 129 experiments in this study.

Our contributions can be summarized under two categories:

### Research

- *The "Random Jitter" form of Quant-Noisier outperforms Quant-Noise.* We identify a method in the Quant-Noisier family, which we call "random jitter," that outperforms Quant-Noise on two out of three quantization schemes for all three tasks.
- *Quantization improves performance on a small speech dataset.* For a small speech dataset, HarperValleyBank, `int4` and `int8` quantization (using either post-training quantization or training-time quantization noise) actually outperforms the uncompressed baseline. Based on follow-up analyses on our speech dataset and prior work, we suspect that quantization has a regularizing effect for this dataset.

### Engineering

- *New Quant-Noise-enabled modules.* We enabled Quant-Noise for layers in RoBERTa and CTC architectures; this involved custom adaptations of stacked LSTMs and objective functions. We plan to open a pull request to FAIR's sequence modeling toolkit (FairSeq).
- *New quantization schemes.* We added support for `int4` and `int1` (previous lowest was `int8`) quantization, which involved low-level PyTorch modifications. A pull request to integrate a subset of these changes to FairSeq can be found at <https://github.com/pytorch/fairseq/pull/3370>.
- *Quant-Noisier.* We added functionality to vary the quantization noise rate per example as a function of previous loss, converged uncompressed loss, or schedule; these involved modifications to the base modules and the trainer that coordinates them.

## 2 Related Work

Model compression methods reduce the memory requirements for neural networks. Pruning and knowledge distillation are forms of compression that reduce the number of network weights. Pruning is a method that removes weights based on their network-level importance. For instance, magnitude pruning introduced in [6] removes weights with low magnitude. Knowledge distillation, as first named in [7], begins with a large pre-trained teacher model, and trains a compact student model by exposing it to raw predictions generated by the teacher model on an unlabeled dataset.

Quantization, by contrast, minimizes the number of bits needed to represent each weight. Standard neural networks use 32-bit floating point precision; scalar quantization replaces these weights with an  $N$  bit representation. Applying quantization methods like `int1`, `int4`, or `int8` as a post-training step causes errors to accumulate at each layer in the model, leading to significantly worse performance. Quantization Aware Training (QAT) is a method proposed in [8] that applies simulated quantization *during training*, so the network learns robustness to the quantization transformations. However, quantization operators tend to be stepwise functions and thus have null gradient; therefore, during back-propagation, gradients are approximated by a straight through estimator (STE) [9]. The STE is a simple idea: treat the stepwise function as if it were the identity function [9]. As the original

authors acknowledge, this estimator is "clearly...biased" [9]. In high compression regimes like `int4`, the error due to biased gradient flow from STE is severe.

Quant-Noise seeks to reduce this error [3]. Quant-Noise only applies simulated quantization to a random subset (at a rate called the "noise rate") of weights during training, allowing some unbiased gradient flow in backpropagation, improving upon the QAT method. Note that the noise rate is fixed across all examples and all epochs, and is tuned empirically by the authors in [3]. However, it is reasonable to think that a model might benefit from using different noise rates for different examples or different epochs. For instance, a model might benefit from using harder, or earlier, examples to learn the task's objective and using easier, or later, examples to learn robustness to quantization.

### 3 Approach

#### 3.1 Quant-Noise and Quant-Noisier

Let  $p$  be the noise rate—the proportion of weights that undergo simulated quantization during training. Formally, we partition a weight matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$  in  $r \times s$  blocks  $\mathbf{b}_{ij}(k, e)$  at location  $(i, j)$  for a given training example  $k$  for a given epoch  $e$ . Let  $\varphi$  be a noise function which simulates quantization during training and let  $p_{ke}$  be the noise rate for example  $k$  at epoch  $e$ . In the case of scalar quantization, each block is simply a single weight. A1 in the Appendix has additional information on scalar quantization.

Then, for a forward pass, the authors apply the following function to each block  $\mathbf{b}_{ij}(k, e)$  in the weight matrix.

$$\psi(\mathbf{b}_{ij}(k, e)) = \begin{cases} \varphi(\mathbf{b}_{ij}(k, e)) & \text{with probability } p_{ke} = p \\ \mathbf{b}_{ij}(k, e) & \text{otherwise.} \end{cases} \quad (1)$$

Notice that all examples, at every epoch, share the same  $p$  value. We hypothesize that "harder" examples should propagate more unbiased gradient flow through the network in order for the model to better learn these examples. Thus, it would be beneficial to set a lower  $p_{ke}$  for these examples so fewer computational nodes require bias-inducing STE. Conversely, "easier" examples should have a higher value of  $p_{ke}$  to teach the network to be robust to quantization. In other words, our training procedure will leverage hard examples—in the active learning sense—to teach the network how best reach its objective and leverage easy examples to maintain performance—in the QAT sense—after compression.

With this in mind, we propose the following (novel) modification to Quant-Noise, which we call Quant-Noisier:

$$\psi(\mathbf{b}_{ij}(k, e)) = \begin{cases} \varphi(\mathbf{b}_{ij}(k, e)) & \text{with probability } p_{ke} = \pi(k, e) \\ \mathbf{b}_{ij}(k, e) & \text{otherwise} \end{cases} \quad (2)$$

where  $\pi(k, e)$  is some function that takes in example  $k$  and an epoch  $e$ . For Quant-Noise baselines, we set  $p_{ke}$  to be a constant value, which we call  $\hat{p}$ . For our adaptive methods, we experiment with three variants of  $\pi(k, e)$ .

**Variante 1.**  $\pi_{\text{hard1}}(k, e) = \hat{p} - \lambda h(k, e)$ , where  $h(k, e)$  quantifies the *hardness* of the example for the quantized model and  $\lambda$  is a positive scaling factor. In these experiments, we set  $h(k, e)$  to the loss  $\mathcal{L}_{k, e-1}$  for example  $k$  on the previous epoch  $e - 1$  that is min-max scaled to the range  $[-1, 1]$  using loss for all examples at epoch  $e - 1$ . We add  $\hat{p}$  so the distribution of  $\pi(k, e)$  is centered around the noise rate we use for Quant-Noise experiments.

**Variante 2.**  $\pi_{\text{hard2}}(k, e) = \hat{p} - \gamma g(k)$ , where  $g(k)$  quantifies the hardness of the example for the *unquantized* model and  $\gamma$  is another positive scaling factor. This variant is analogous to the previous, with the difference in measuring hardness with the unquantized model instead. That is, the hardness is measured by the loss of the uncompressed model  $g(k) = \mathcal{L}_k^{\text{uncompressed}}(\theta)$  which is only a function of example  $k$ ; that is, the hardness is constant across epochs. We normalize  $g(k)$  using loss for all examples for the unquantized model at convergence.

**Variante 3.**  $\pi_{\text{sched}}(k, e) = \hat{p} - \alpha z(e)$ , where  $z(e)$  quantifies the current position in training and  $\alpha$  is another positive scaling factor. This approach is based on scheduled learning rates where the

optimal learning rate early on in training is distinct from the optimal learning rate in later stages [10]. We hypothesize that quantization noise rate benefits from a similar schedule. Low values of  $p$  in the beginning let the network learn the task at hand. Then, when the model has had the chance to reasonably fit to the training data, it can begin developing resilience to quantization. In these experiments, we choose  $z(e)$  to be a linear function such that  $\pi_{\text{sched}}$  takes on a mean value of  $\hat{p}$  over all epochs. More precisely, we set  $z(e) = 1 - \frac{2e}{m}$ , where  $m$  is set number of epochs the model will train for.

Finally, in what was intended as an ablation study, we propose "random jitter," which adds a random amount of noise to the noise rate, drawing from a uniform distribution  $\pi \sim \hat{p} - \mathcal{U}(-\beta, \beta)$ , where  $\beta$  is a tunable hyperparameter.

We built off of the original codebase for Quant-Noise, available at [11]. Each of the variants for  $\pi(k, e)$  were custom implementations. This involved functionality to vary the quantization noise rate per example, per epoch, and per schedule using modifications to PyTorch modules, the trainer that coordinates them, and other functions called by the trainer. Two of our compression regimes (`int1` and `int4`) also involved additions to FAIR’s quantization operators and corresponding low-level PyTorch code. A pull request (found at <https://github.com/pytorch/fairseq/pull/3370>) has been opened to merge a small subset of our additions to FairSeq [11].

### 3.2 Model Architectures

For our NLP experiments, we fine-tune a RoBERTa base (available at [12]) with an added classification head. For each of RTE and MRPC, we establish baselines with uncompressed models, post-training quantization, and Quant-Noise. Each baseline is measured using `int8`, `int4`, and `int1` schemes.

For our SLP experiments, we train (from scratch) an RNN-based network tasked on a multi-task objective. We adopt the training framework proposed in [13] and open-sourced at [14] that jointly optimizes CTC loss on the task of Automatic Speech Recognition (ASR) and Cross-Entropy (CE) loss on the auxiliary task of intent classification. Henceforth, we refer to this architecture as the CTC model. Note that our reported downstream performance is measured with respect to the task of intent classification since its validation curve converged 3x faster than the curve for ASR. The CTC model is composed of a recurrent layer, a stacked LSTM with *depth* = 2, and a softmax layer. Unlike many of the layers required for RoBERTa, LSTMs are not supported off-the-shelf by Quant-Noise. In addition, PyTorch also does not release its LSTM implementation since much of it is in C for optimization purposes. Enabling quantization noise with LSTMs for our experiments required adapting and integrating a native PyTorch LSTM frame from the GitHub repository located at [15]. We plan to open a pull request for this work soon.

## 4 Experiments

### 4.1 Data

For NLP branch of experiments, we focus on binary classification problems of (1) detecting textual entailment and (2) determining semantic equivalence. The former involves determining whether a given text fragment is entailed by another text fragment [16]; we use the RTE dataset (2.5k fragment pairs). The latter involves determining whether a pair of text fragments have semantic equivalence; we use the MRPC dataset (5.8k sentence pairs [17]). Both datasets are drawn from the GLUE benchmark [18]. For our SLP branch of experiments, we focus on the task of intent classification. Our chosen dataset, HarperValleyBank (1.4k customer-agent conversations) [13], frames this as an 8-class (e.g. *order checks*, *transfer money*) classification problem. We take as input Mel-frequency cepstral coefficients (MFCCs), which are encoded representations of raw audio signal. MFCC feature generation is the standard process of encoding audio in speech.

### 4.2 Evaluation method

We consider both the task-specific evaluation metrics (accuracy for all three tasks considered) and the compression ratio (*original\_model\_size* / *quantized\_model\_size*).

### 4.3 Experimental details

**NLP** We slightly modify the hyperparameters suggested in [11] under RoBERTa fine-tuning. First, we remove settings related to fp16. Second, due to memory constraints, we were only able to train with a batch size of 4 instead of 16, but we set our parameter update frequency to 4 (instead of the default 1) so that gradient updates would still occur at a batch size of  $4 * 4 = 16$ . All of our experiments use  $\hat{p} = 0.5$ , a reasonable default choice, as indicated by [11] and [3]. We also set  $\lambda = 0.125$  for  $\pi_{\text{hard1}}$ ,  $\gamma = 0.125$  for  $\pi_{\text{hard2}}$ ,  $\alpha = 0.4$  for  $\pi_{\text{sched}}$ , and  $= 0.125$ . We train for 10 epochs.

**SLP** We largely maintain the same hyperparameter settings as HarperValleyBank’s baselines (found at [14]). Data augmentation is also performed (reasoning explained in *Results*) with time and frequency masks applied with a probability of 0.5 for each training example at every epoch. See Figure 6 for an example. Note that we used SpecAugment to apply these augmentations [19]. Another deviation from the original model configuration is that we train for 40 epochs instead of 200. As with our NLP experiments, we set  $p = 0.5$  when evaluating Quant-Noise.

### 4.4 Results

First, a general note on figures: dashed lines indicate performance of the uncompressed model; error bars represent standard deviations based on the three seeds for each trial.

### 4.5 NLP Results

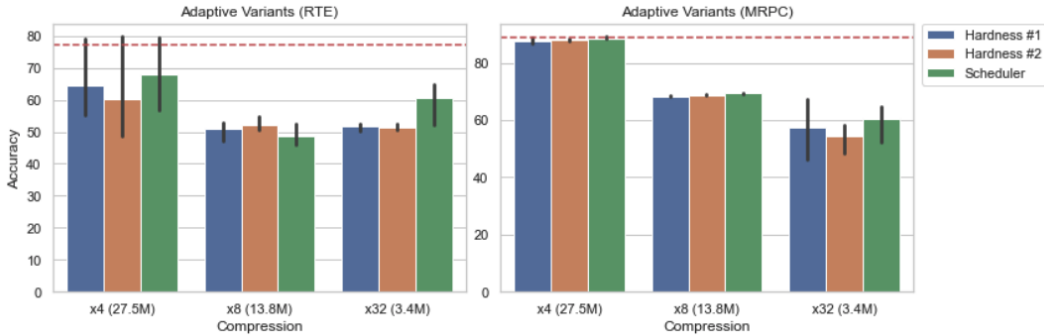


Figure 1: We compare performance of our three adaptive variants. The scheduler variant appears to perform best.

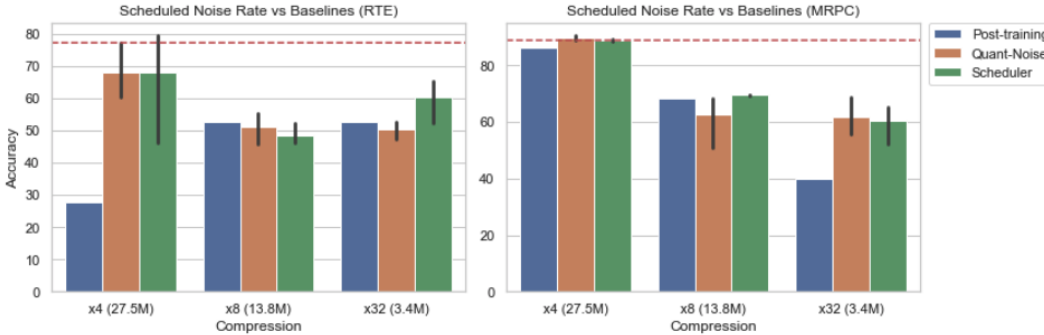


Figure 2: We compare the performance of our scheduler against Quant-Noise and the post-training quantization baseline. The scheduler variant is roughly even with Quant-Noise and outperforms post-training quantization.

We display our results in Table 2. We report average performance across three training runs, each with a different seed, with the exception of the uncompressed model and post-training quantization, for which we only train one model.

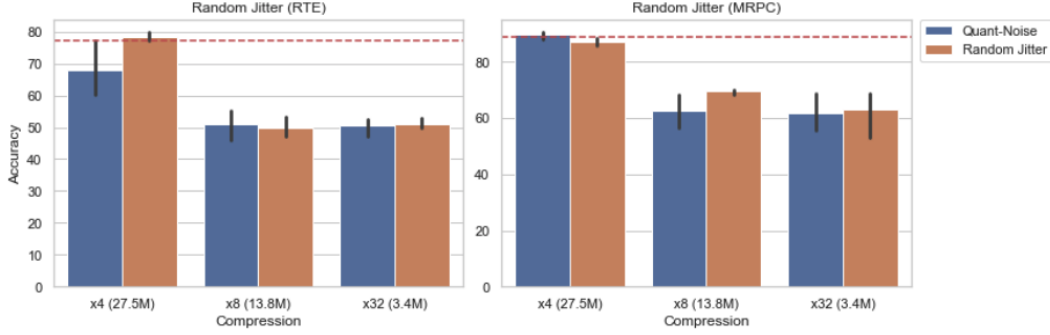


Figure 3: We compare the performance of "random jitter," our best method from the Quant-Noisier family, against Quant-Noise. Random jitter outperforms Quant-Noise.

As corroborated by Figure 1, we find that of the 3 different adaptive variants, the scheduled adaptive Quant-Noise variant outperforms other adaptive variants for 2 of the 3 quantization schemes on RTE. Namely, the schedule adaptive variant attains an accuracy of 67.99, 48.49, and 60.46 for `int8`, `int4`, and `int1`, respectively. Likewise, on MRPC, the adaptive scheduler variant consistently outperforms other adaptive variants for all quantization schemes.

Furthermore, as shown in Figure 2, we compare the adaptive scheduler variant to our baselines of Quant-Noise and post-training quantization. Notably, the adaptive scheduler variant outperforms most post-training quantization and some Quant-Noise baselines, attaining a higher average accuracy than Quant-Noise on 3 of the 6 datapoints for RTE and MRPC.

Most surprisingly, we found that fine-tuning RoBERTa base using Quant-Noise + random jitter is better than expected; this outperforms the Quant-Noise baseline on 4 of 6 (task, quantization scheme) pairs, as shown in Figure 3. In particular, as illustrated in 1, random jitter on RTE substantially outperforms Quant-Noise for RTE on `int8` and MRPC on `int4` (78.46% vs. 67.87%; 69.53% vs. 62.50%, respectively).

The introduction of second-order noise with random jitter leads to improvements upon Quant-Noise for `int1`, `int4`, and `int8` quantization schemes. By jittering the noise rate, all examples have the same expected noise rate, but the amount of noise varies randomly throughout training. More precisely, that for an example  $x_k$  sampled from the training data  $\mathcal{D}$ , the expected noise rate of any example is  $\mathbb{E}_{x_k \sim \mathcal{D}}[\pi(k)] = \hat{p}$ . The additional layer of random noise appears to help the model learn further robustness to quantization. This is a surprising result and merits further investigation to determine precisely why this second-order noise is helpful.

Table 1: Quant-Noise versus Random Jitter with `int1`, `int4`, and `int8` quantization schemes.

Quantization	Scheme/compression rate on RTE			Scheme/compression rate on MRPC			Scheme/compression rate on HVB		
	<code>int8</code> / <code>x 4</code>	<code>int4</code> / <code>x 8</code>	<code>int1</code> / <code>x 32</code>	<code>int8</code> / <code>x 4</code>	<code>int4</code> / <code>x 8</code>	<code>int1</code> / <code>x 32</code>	<code>int8</code> / <code>x 4</code>	<code>int4</code> / <code>x 8</code>	<code>int1</code> / <code>x 32</code>
Random jitter	<b>78.46</b>	49.82	<b>51.02</b>	86.93	<b>69.53</b>	<b>62.83</b>	<b>45.1</b>	41.1	<b>17.4</b>
Quant-Noise	67.87	<b>51.14</b>	50.54	<b>89.46</b>	62.50	61.85	42.5	<b>42.6</b>	14.1

#### 4.6 SLP Results

Our group’s custom implementation of the CTC model with quantization noise attains an accuracy score of 39.6% (see Figure 7), which is within reasonable range of the baseline reported in the original HVB study [13]. The subsequent baselines we establish for *Post-training Quantization* and *Quant-Noise* have unexpected results. One would expect the former to lead to heavy regressions in performance while the latter is shielded from much of it. Figure 7 displays a contradictory trend: the two methods perform equally as well for all compression levels and, even more remarkably, improves on uncompressed accuracy with 4x and 8x scalar quantization. This is a surprising result

and indicates that quantization can actually have a helpful effect on performance for this dataset; we explore this more in *SLP Analysis*.

Figure 4 illustrates the results of a follow-up experiment with data augmentation added. Stronger regularization removed the unexpected positive correlation between accuracy and compression ratios. `int8` quantization (4x compression) produced a slight drop in accuracy, similar to what is observed in [3]. `int4` quantization (8x compression) continues to perform just as well as the `int8` model. The lack of a performance gap between post-training quantization and Quant-Noise is still evident even with `int4` quantization where [3] observed significant gains with using quantization noise.

We again observe that adding second-order noise in the form of random jitter (the best performing form during the NLP experiments) improves upon the original Quant-Noise method. It does so for both `int8` ( $\times 4$ ) and `int1` ( $\times 32$ ) quantization, as can be seen in Table 1.

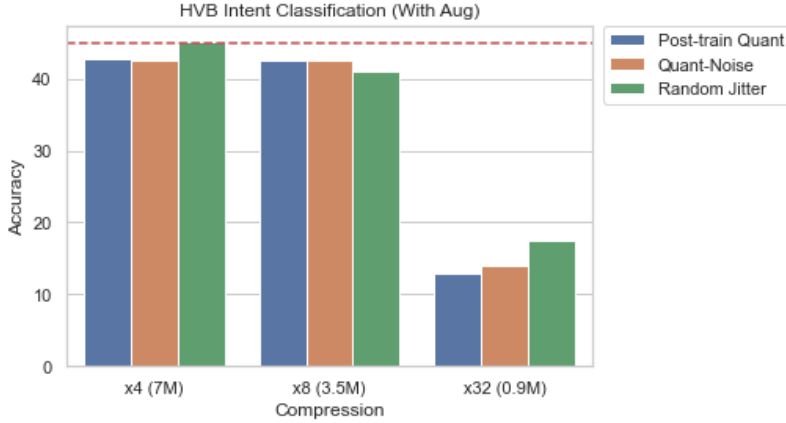


Figure 4: Repeat *Post-training Quantization* vs *Quant-Noise* experiment with time and frequency mask augmentations.

## 5 Analysis

### 5.1 NLP Analysis

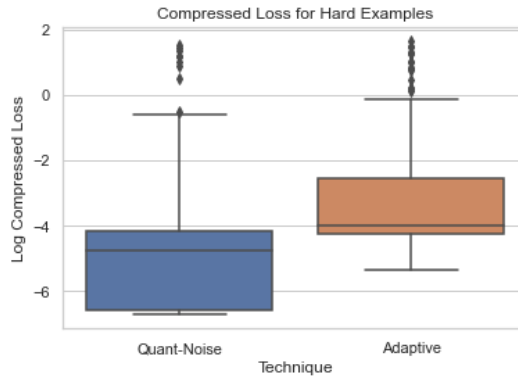


Figure 5: We compare loss on the hardest training examples in RTE (90th percentile of loss of the uncompressed model) between a Quant-Noise model and an adaptive model (variant 2).

One of the underlying assumptions for our adaptive variants was that unbiased gradient flow on an example would help the model better learn that example. Is this borne out by our results? In this section, we conduct qualitative analyses to answer this question.

In Figure 8, we display a scatterplot with the loss of the uncompressed examples on RTE on the x-axis, and the loss of the adaptive (variant 2) model and the Quant-Noise model on the y-axis. Both y-axis

models use `int8` quantization, and we plot losses on training examples. Our scatterplot indicates the opposite of what we expected; loss does not appear lower for the "harder" (higher uncompressed loss) examples on our adaptive variant 2 model as compared to Quant-Noise.

We also create a boxplot comparing the loss of the Quant-Noise model against the loss of the adaptive (variant 2) model on the "hardest" (highest uncompressed loss) training examples. Our scatterplot indicates the opposite of what we expected; loss is higher for the hardest examples on our adaptive variant 2 model as compared to Quant-Noise.

## 5.2 SLP Analysis

A significant portion of the speech evaluations are aimed at testing our group’s custom implementation of the CTC model with quantization noise. Results displayed in Figure 7 under-perform relative to the scores reported in the original HVB study (47.8% accuracy) [13]. Meeting with one of its authors confirmed that this was expected since there were additional steps not included in the released code. Incorporating data augmentation as in Figure 4 closes this gap by attaining 45% accuracy.

We then move to understand the impact of established quantization methods on model behavior for our problem setting since Quant-Noise is untested on both speech data and recurrent models. First, we explain the counter-intuitive positive correlation of compression level and accuracy in our baseline. Diagnostics on the train and validation loss curves reveal that models at all compression levels are overfitting to the small dataset. This points an explanation based on quantization’s regularizing effect: perhaps gains with higher compression ratios was due to lowering variance? Replicating the baselines with stronger regularization in the form of data augmentation confirmed this hypothesis by removing the positive correlation. Notable approaches, such as one done by Wu and Flierl [20] and another by Hirose et al. [21], explicitly use quantization as regularization mechanisms. Our results in this study indicate that such methods are likely to help CTC-based models generalize as well.

Attaining explainable results after data augmentation allowed us to reliably test *random jitter*, the best second-order noise method from the NLP trials. *Random jitter* outperforming *Quant-Noise* on two compression schemes provides additional empirical evidence that the former produces consistent gains. This is especially strong considering that speech is a different domain that operates on waveforms that are quite different from text-based sequences. In addition, HarperValleyBank is the smallest dataset we test, continuing support for *Random jitter* across dataset sizes as well.

Likely reasons for the small performance gap between *Post-training Quantization* and the baselines are that (1) intent classification is straight-forward enough of a task such that even aggressively quantized networks can effectively learn it and (2) quantization noise is not much more effective than post-training quantization in low-resource regimes. The first reason is tested by mirroring the experiment from Figure 4 and evaluating ASR word error rate instead of intent classification. All performance trends persisted, thus disproving the first explanation. The second, on the other hand, is supported by the fact that quantization noise brings about significant gains with experiments involving Wikitext-103 (100 million tokens), ImageNet (1.2 million images), and RTE/MRPC (small datasets, but used with pretraining) while it does not on the smaller HarperValleyBank dataset. Insights from our NLP runs also supports this hypothesis: Quant-Noise is a high-variance method. One of the well-studied properties of high-variance methods is that they require more data to outperform less expressive alternatives [22]. Though this is a likely explanation with a theoretical base, future work can provide empirical proof by using our framework to evaluate ASR with a larger dataset such as LibriSpeech [23].

## 6 Conclusion

In this study, we developed novel methods for quantization, learning that adding random jitter to Quant-Noise can outperform Quant-Noise on three datasets. Our work has limitations: first, for our NLP tasks, we perform both model selection (i.e., selecting the best model checkpoint from training) and evaluation on the development set. This is standard practice for GLUE tasks, but if we had more time, it would be better to evaluate our best models on the GLUE test set. Similarly, for HarperValleyBank, we perform both model selection and evaluation on the development set; again, however, this is consistent with what is done in the original paper [13]. Second, for adaptive variants (1) and (2) we normalize our losses to minimum of  $-1$  and maximum of  $1$ . In both cases, we use



min-max scaling, but other methods, which may involve dividing by mean loss and clamping to  $[-1, 1]$ , might work better.

In the future, we will further investigate why our adaptive method led to increased loss on harder examples as compared to Quant-Noise. If we can find a reliable method for decreasing loss on harder examples, this may be a promising avenue for further performance gains in quantization. Additionally, we hope to measure not just compression ratios but also speedups in inference time. Finally, we could, rather than hand-writing functions for the noise rate, search over a large space of functions, perhaps using reinforcement learning (where reward is performance on the development set).

## References

- [1] T Brown and B et al. Mann. Language models are few-shot learners, 2020.
- [2] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks, 2018.
- [3] Angela Fan, Pierre Stock, Benjamin Graham, Edouard Grave, Remi Gribonval, Herve Jegou, and Armand Joulin. Training with quantization noise for extreme model compression, 2020.
- [4] R Ding and Z et al. Liu. Quantized deep neural networks for energy efficient hardware-based inference, 2018.
- [5] P Yin and J et al. Lyu. Understanding straight-through estimator in training activation quantized neural nets, 2019.
- [6] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks, 2015.
- [7] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015.
- [8] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference, 2017.
- [9] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation, 2013.
- [10] V Plagianakos and M Magoulas. Learning rate adaptation in stochastic gradient descent, 2016.
- [11] Fan et al. Fairseq github. *GitHub. Note: <https://github.com/pytorch/fairseq>*, 2020.
- [12] Yinhan Liu and Myle Ott et al. Roberta: A robustly optimized bert pretraining approach, 2019.
- [13] Mike Wu, Jonathan Nafziger, Anthony Scodary, and Andrew Maas. Harpervalleybank: A domain-specific spoken dialog corpus, 2020.
- [14] Wu et al. Harpervalleybank github. *GitHub. Note: <https://github.com/cricketclub/gridspace-stanford-harper-valley>*, 2020.
- [15] Daehwan Nam. Pytorch rnn util. *GitHub. Note: <https://github.com/daehwannam/pytorch-rnn-util>*, 2019.
- [16] ACL Team. Recognizing textual entailment. "[https://aclweb.org/aclwiki/Recognizing\\_Textual\\_Entailment](https://aclweb.org/aclwiki/Recognizing_Textual_Entailment)", 2013. [Online; Accessed 26-Feb-2021].
- [17] William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005.
- [18] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding, 2019.

- [19] Zach Caceres. Specaugment with pytorch. *GitHub*. *Note:* <https://github.com/zcaceres/specaugment>, 2019.
- [20] Zach Caceres. Quantization-based regularization for autoencoders. 2019.
- [21] K Hirose and K et al. Ando. Quantization error-based regularization in neural networks. 2017.
- [22] Jason Brownlee. How to reduce variance in a final machine learning model. 2018.
- [23] Panayotov et al. Librispeech: an asr corpus based on public domain audio books. 2015.

## A Appendix (optional)

### A.1 Scalar Quantization

Scalar quantization uses lower-bit representations of floating point weights. According to [3], for fixed-point `intn` quantization, we quantize weights by applying the following element-wise transform:

$$Q = (\text{round}(W/s + z) - z) \times s \quad (3)$$

where the *scale* is  $s = \frac{\max W - \min W}{2^n - 1}$  and the *bias* is  $z = \text{round}(\min W/s)$ . Note, the activations are also quantized at inference time. Following this compression scheme, the corresponding compression ratio is  $32/n$ .

### A.2 Supplementary Figures and Tables

Table 2: Fine-tuning RoBERTa on RTE and MRPC using Quant-Noise versus Random Jitter with `int1`, `int4`, and `int8` quantization schemes.

Quantization	Scheme/compression rate on RTE			Scheme/compression rate on MRPC		
	<code>int8</code> ×4	<code>int4</code> ×8	<code>int1</code> ×32	<code>int8</code> ×4	<code>int4</code> ×8	<code>int1</code> ×32
Uncompressed	77.26	77.26	77.26	88.73	88.73	88.73
Post-training Quantization	27.80	<b>52.71</b>	52.71	86.03	68.14	39.71
Quant-Noise	67.87	51.14	50.54	<b>89.46</b>	62.50	61.85
Adaptive noise variant 1	60.40	52.22	51.26	87.99	68.79	54.33
Adaptive noise variant 2	60.40	52.22	51.26	87.99	68.79	54.33
Adaptive noise variant 3	67.99	48.49	<b>60.46</b>	88.65	69.36	60.46
Random jitter	<b>78.46</b>	49.82	51.02	86.93	<b>69.53</b>	<b>62.83</b>

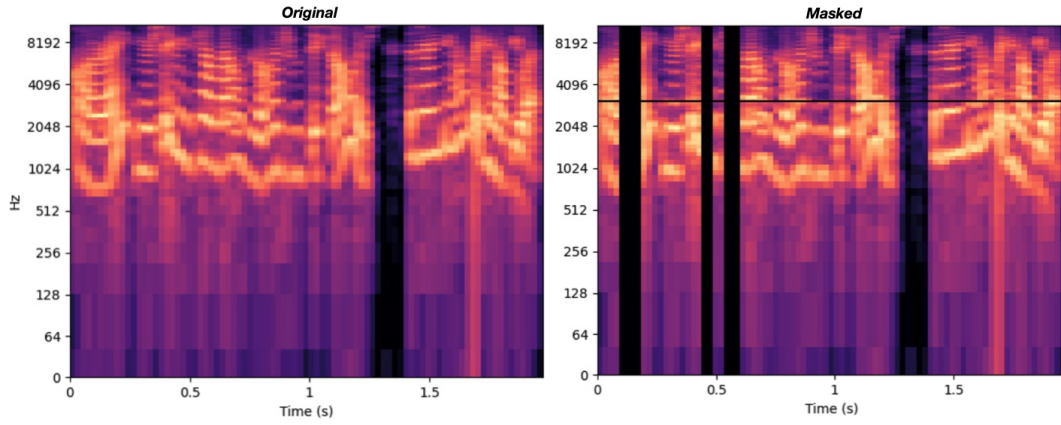


Figure 6: A sample spectrogram with time and frequency masks applied.

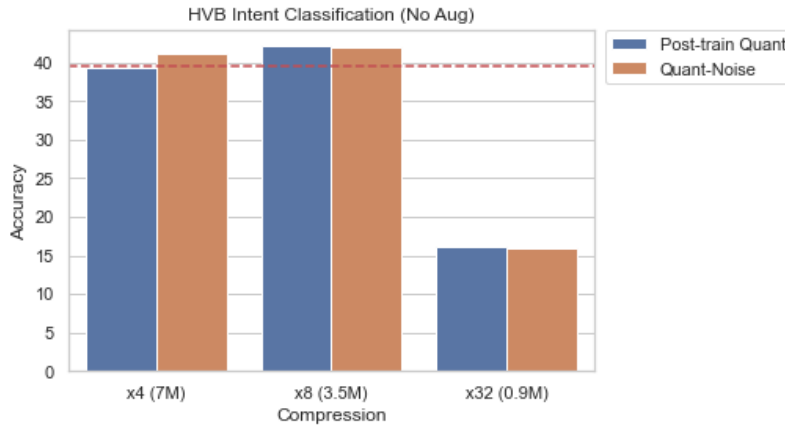


Figure 7: Establish baselines with *Post-training Quantization* to *Quant-Noise* at four levels of compression.

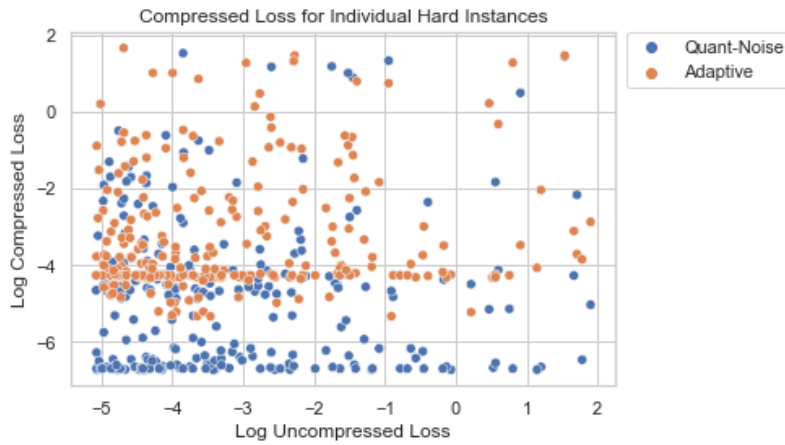


Figure 8: We plot loss on training examples for the uncompressed model against loss on these same examples for the Quant-Noise model and adaptive model (variant 2).