

# ACCELERATOR-AWARE TRAINING FOR TRANSDUCER-BASED SPEECH RECOGNITION

*Suhaila M. Shakiah, Rupak Vignesh Swaminathan, Hieu Duy Nguyen, Raviteja Chinta<sup>†</sup>,  
Tariq Afzal<sup>†</sup>, Nathan Susanj, Athanasios Mouchtaris, Grant P. Strimel, Ariya Rastrow*

Alexa Machine Learning, Amazon.com, USA; <sup>†</sup>Hardware Compute Group, Amazon.com, USA

## ABSTRACT

Machine learning model weights and activations are represented in full-precision during training. This leads to performance degradation in runtime when deployed on neural network accelerator (NNA) chips, which leverage highly parallelized fixed-point arithmetic to improve runtime memory and latency. In this work, we replicate the NNA operators during the training phase, accounting for the degradation due to low-precision inference on the NNA in back-propagation. Our proposed method efficiently emulates NNA operations, thus foregoing the need to transfer quantization error-prone data to the Central Processing Unit (CPU), ultimately reducing the user perceived latency (UPL). We apply our approach to Recurrent Neural Network-Transducer (RNN-T), an attractive architecture for on-device streaming speech recognition tasks. We train and evaluate models on 270K hours of English data and show a 5-7% improvement in engine latency while saving up to 10% relative degradation in WER.

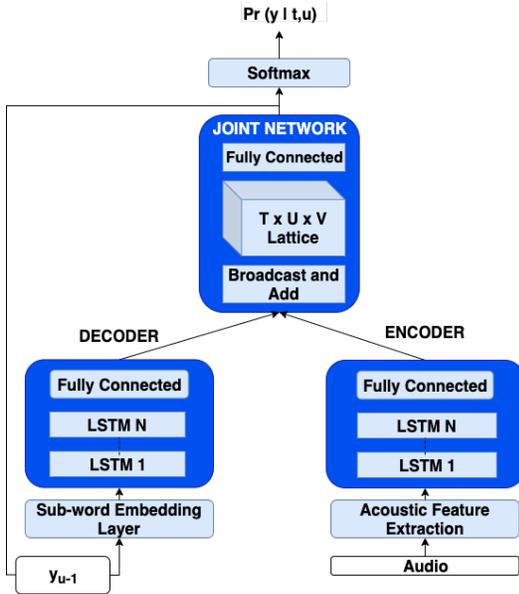
**Index Terms**— Accelerator-aware training, model compression, automatic speech recognition (ASR), recurrent neural network transducer (RNN-T).

## 1. INTRODUCTION

The task of transcribing an audio to the corresponding text transcriptions constitutes the Automatic Speech Recognition (ASR) component of voice assistants such as Alexa, Google Home or Siri. ASR solutions have evolved from traditional hybrid Deep Neural Network (DNN) - Hidden Markov Model (HMM) systems to modern end-to-end neural architectures including various transducer-based systems [1–5]. While many of these approaches have demonstrated high accuracy, an important differentiation comes from their streaming capability, which reduces the user’s perceived latency (UPL). A streaming ASR system is able to start transcribing the audio even before the user has finished speaking an utterance, i.e., the system does not require future context to inform the transcription results. Among streaming architectures, the Recurrent Neural Network-Transducer (RNN-T) model (see Figure 1) [6] stands out as an all-neural, end-to-end (E2E) method with both low latency and high accuracy, and it is widely adopted in modern speech recognition systems.

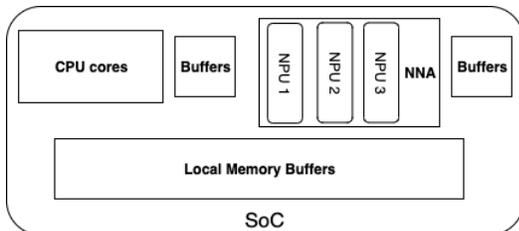
Apart from the current trend of moving to E2E architectures for machine learning (ML) applications, leading ML solution providers also utilize on-device processing to improve user experience and reduce UPL. As a result, hardware consisting of NNA chips have been gradually deployed to support on-device computer vision, natural language understanding (NLU), and ASR tasks. However, this comes with additional challenges when moving from floating point to fixed point operations supported by NNAs. Since computers are finite state machines, real numbers are represented and manipulated in floating point format in computer memory. Floating point numbers are characterized by a mantissa, an exponent and the sign bit, which enable the representation of a wide range of values with a floating decimal point. For example, a 32-bit floating point (FP-32) number can represent a total of  $2^{32}$  unique values with a higher precision than their fixed point counterparts. In the context of performing machine learning on the edge, floating point computations are time consuming and memory expensive, especially for deep learning models used for E2E speech recognition which involve millions (and sometimes billions) of multiply-and-accumulate (MAC) operations for a single inference cycle. To address the time and memory complexity, neural accelerator chips thus employ fixed point operations, in which each value is normally represented by a reduced number of bits.

It is worth noting that neural ML models are highly sensitive to such reductions in the precision of weights and activations, and this effect is even more pronounced in recurrent architectures since the errors accumulate across multiple time steps. To address these problems, while also being cognizant of the low latency, power and memory requirements for on-device systems, the on-device chips adopt a hybrid architecture which includes general-purpose central processing units (CPUs) as well as specialized neural processing unit (NPU) cores, integrated into a single System-On-Chip (SoC) design (see Figure 2). In order to achieve the best performance, certain computations are performed in highly efficient NPUs whereas others requiring a higher precision are computed on CPUs. This trade-off results in additional on-chip computing and data transfer latency between the NPU and CPU, creating bottlenecks during inference. For example, on-device ASR models need to use the CPU to compute tanh and sigmoid activations, which are not only quantization error-prone



**Fig. 1:** Model diagram of Recurrent Neural Network Transducer (RNN-T).

on NPU, but also computation and memory intensive. In this work, we aim to reduce the latency incurred due to moving data between the NPU and CPU to perform non-linear activations. Using an ASR task and RNN-T architecture as a proof of concept, we show that we are able to improve model inference speed by 20% on-device with negligible, less than 1% relative, accuracy degradation by performing accelerator-aware training (AAT), thus making the models more robust to the NNA and the hardware activation functions.



**Fig. 2:** A representative SoC showing CPU and NPU cores with local memory and data buffers

The rest of the paper is structured as follows. Section 2 describes the related work for latency reduction in RNN-T speech models, followed by an overview of quantization in neural networks and NNAs in Section 3, as well as details of accelerator quantization schemes in Section 4. Section 5 describes our proposed accelerator-aware training technique. We then dive into the experimental details in Section 6 and the performance and latency results in Section 7. Section 8 concludes our paper with some remarks on the advantages of the proposed approach.

## 2. RELATED WORK

There is a large body of work focussing on improving the Word Error Rate (WER) and runtime latency for on-device RNN-T models. Architectural modifications to recurrent neural networks, such as CIFG-LSTM [7] and Simple Recurrent Units [8] have been used in RNN-T saving 30-40% compute while having negligible impact on the WER. Knowledge distillation techniques specific to RNN-T have also been studied [9–11]. Employing sparsity-based pruning on the weight matrices of LSTMs have been studied for both structured and unstructured sparse matrices [5, 12]. In [13], the authors proposed an additional regularization term in the RNN-T loss to penalize blank token prediction so that the model emits the labels faster, leading to a reduction in latency while maintaining the WER. In [14], the authors introduced a bifocal encoder architecture for RNN-T to improve streaming latency, where the low entropy wakeword segment of an utterance is processed by a small encoder allowing the larger encoder to catch up during decoding of the rest of the utterance as the frames become available. As a more general approach to switching between encoders of different compute capacity, [15] proposed an arbitrator network that could dynamically choose the encoder network on a per-frame basis.

In addition to modified architectures and loss functions, there have also been many advances in training methods that emulate inference; quantization-aware training (QAT) and sparse pruning methods have been proposed to improve on-device runtime latency without hurting accuracy [16–21]. In [22], the authors propose an accelerator-aware neural design where the architecture search space is explored to optimize performance on NNA. More relevant to this work is [23], in which the authors describe a quantization strategy for RNN-T based speech recognition systems using 16-bit activations. This work focuses on accelerators that uses 8-bit activations, demonstrating on-par performance with unquantized baselines even for models with smaller number of parameters, which are more susceptible to quantization.

## 3. OVERVIEW OF QUANTIZATION

In this section, we provide a brief overview of the quantization schemes, notations and definitions used in this paper. Quantization is the process of converting floating point values to a smaller set of discrete fixed point values, effectively reducing the number of bits used to represent the numbers. We will use the  $Q$ -notation to define the parameters of a signed fixed point number. A fixed point number denoted by  $Q_{m,n}$  has the following properties:

- Total bit width is  $m + n$ , including  $m$  integer and sign bits plus  $n$  fractional bits.
- Consider a signed representation, i.e., the most significant bit (MSB) is 1 for negative values. The minimum

possible fractional value that can be represented,  $f_{min}$ , is  $-2^{m-1}$ , while the maximum representable fractional value,  $f_{max}$ , is  $2^{m-1} - 2^{-n}$ .

- The resolution is  $2^{-n}$ , yielding a maximum quantization error of  $2^{-(n+1)}$  between a number and its quantized counterpart.

For example,  $Q3.2$  has 5 bits in total, 3 integer bits with the MSB indicating the sign bit, 2 fractional bits and can represent fractional numbers in the interval  $[-4, 3.75]$  with a resolution of 0.25, thus able to represent a total of  $2^5$  unique values.

### 3.1. Static Quantization

In static quantization, the input value  $R_i$ , is first clipped to be within the quantizable integer range, thus accumulating clipping error. Afterwards, the resulting FP-32 value is scaled and rounded to the nearest integer, which is then scaled back to its floating point equivalent. This ensures that the numbers, during inference on-device, will have the lowest quantization (rounding) errors. For a  $Q_{m.n}$  fixed point quantization scheme, the quantized floating point equivalent,  $R_q$ , of an input FP-32 value  $R_i$  is given by,

$$R_q = R(C(R_i, f_{min}, f_{max}) 2^n) 2^{-n} \quad (1)$$

in which  $R$  denotes *Round* function, which rounds the values to the nearest integer or towards zero depending on the implementation and  $C$  denotes *Clip* function which clips the values at  $f_{min}$  and  $f_{max}$ . In this work, we round the weights to their nearest value and round the inputs and hidden states towards zero as implemented on the accelerator.

### 3.2. Dynamic Quantization

In dynamic quantization, there is an additional degree of freedom, in which the input values can be scaled with a scaling factor, so that the range fits into the quantizable range of integers. This scaling factor needs to be a power of 2 and is calculated dynamically for every sample and time step for the NNA inference. Simply put, the scaling factor,  $S$ , is the smallest power of 2 that scales the incoming tensor and fits all values within the quantizable range of the given target  $Q_{m.n}$  scheme. A clipping error is incurred if a larger scaling factor is required. After the quantization, the outputs are scaled back up. The modified equation is as follows:

$$R_q = S \times R\left(C\left(\frac{R_i}{S}, f_{min}, f_{max}\right) 2^n\right) 2^{-n} \quad (2)$$

where  $S \in \{1, 2, 4, 8, 16\}$ , such that  $f_{min} \leq \frac{R_i}{S} \leq f_{max}$ .

## 4. NEURAL NETWORK ACCELERATOR

In this section, we briefly discuss how neural network weights and activations are quantized in our NNA experiments. The number of bits for the quantization of weights and activations, the type of quantizations performed, supported layers, the specific data paths between the various on-chip components, etc., differ from one chip version to another. In our experiments, we use the following schemes:

- **LSTM Cells:** The hidden states are statically quantized to  $Q1.7$  in all LSTM cells in the encoder and decoder.
- The inputs to the first LSTM layer in the encoder and decoder are dynamically quantized to  $Q1.7$  (because it directly follows a CPU data path, hence will not incur additional latencies to compute dynamic scaling factors). The inputs to all other LSTM layers in the encoder and decoder are statically quantized to  $Q1.7$
- The sigmoid and tanh activations for calculating the gate values  $i$ ,  $f$ ,  $g$  and  $o$  in the LSTM cells are non-uniformly quantized 8-bit values as shown in Figs. 3a and 3b. As sigmoid and tanh are non-linear, non-uniform quantizations are more suitable to reduce the means-squared error (MSE) than uniform quantization.
- **Dense Layers:** The inputs to all dense layers in the encoder, decoder and joint network are dynamically quantized to  $Q1.7$ .
- **Weights:** All weights in the embedding layer, encoder, decoder and joint network are statically quantized to  $Q1.7$

### 4.1. Hidden States and Inputs Quantization

The NNA uses a symmetric linear quantization scheme to map 32-bit floating point numbers to 8-bit integers, rounded toward zero. Considering our previous example of  $Q3.2$ , the range of representable fractional values is  $[-4, 3.75]$ , which on the hardware is represented as integers in the range  $[-16, +15]$ . Any value less than  $-4$  or greater than  $3.75$  is clipped to these limiting values, thus accumulating clipping error during static quantization. For dynamic quantization, the allowable dynamic scales are 1, 2, 4, and 16, which are calculated on the CPU and used to scale the original values into the quantizable range. The outputs are uniformly quantized with equal step sizes and quantization intervals as illustrated in Fig. 3c.

### 4.2. Non-linear Activation Quantization

The hyperbolic tangent and sigmoid functions are the commonly used non-linear activation functions in neural networks and are integral components in learning long range memory. They are expensive to compute even on CPUs, and thus are

approximated on the NNA hardware through careful digital circuit design and error analysis. An efficient hardware implementation of activation functions is required to meet the performance, area, power and cost targets of neural accelerators. Multiple designs and approximation algorithms have been proposed to balance this trade-off [24–26]. Usually, a combination of linear interpolators, shifting operations, look-up tables and multiplexers are used in the digital circuit to approximate the activation values. We are interested in translating the on-device operators into the model training workflow to train accelerator-aware models.

For the accelerator considered in this work, activation functions are implemented as a piece-wise linear approximator, which gives non-uniformly quantized 8-bit values as outputs. The design, modeling and analysis of these digital circuits are outside the scope of this work (please refer to [24–26]). The approximation algorithm is optimized to balance the trade-off between accuracy and on-chip area, and leverages the fact that mathematically, the tanh and sigmoid functions are shifted and scaled versions of one another. It also reduces quantization error by carefully choosing non-uniform quantization centers to pack more quantization bins at parts of the functions with higher gradients. Despite the careful design, it is not possible to circumvent the information loss due to quantization. Given that the activation functions play a key role in learning long-term memory, as we show later in the results, switching from high precision to 8-bit values leads to a large performance degradation. To alleviate this performance degradation, these values need to be computed on the CPU, which requires constant transfer of data on the chip to the CPU in every time step of the LSTM cell to perform activation functions, and back to the NNA to perform MAC operations, incurring additional processing latency. A major contribution of this paper is to incorporate the bit-exact hardware operators into the training workflow with meaningful gradient backpropagation and further fine-tune the model to non-uniformly quantized activations, yielding 5-7% latency reduction and negligible, less than 1% relative WER degradation.

## 5. ACCELERATOR-AWARE TRAINING

We propose a two-stage method to address the errors propagated by the activation and intermediate quantizations in the model. Before diving into the implementation details, we first discuss the motivation for choosing such an approach.

The sigmoid and tanh functions saturate approximately at  $\pm 7$  and  $\pm 4$  respectively (Figs. 3a and 3b). However, in the process of quantizing the activations, a large clipping and precision error is accumulated between  $\pm 7$  and  $\pm 4$ , for the activations when compared to the FP-32 (full precision) values. This is because a significant portion of the inputs to the activation functions in the encoder and decoder LSTMs of a fully trained RNN-T model lie in this error-prone range. Due

to the large number of values present in the saturation range of these functions, the vanishing gradient problem hinders learning using backpropagated gradients, which is already affected due to the addition of quantization. To overcome this issue, we add an activity regularization loss on the inputs to the activation functions as described in Section 5.1.

For tuning the inputs, hidden states and activation function values into the quantization levels, we perform a bit-exact quantization in the forward pass and approximate the gradients in the backward pass using straight-through estimation. This is detailed for the various components in Section 5.2.

### 5.1. Activity Regularization

In Stage I of training, we initialize the model with random weights and train it from scratch with an activity regularization loss on the outputs  $z_0, z_1, z_2$ , and  $z_3$  of the LSTM cell [27] to alleviate the quantization errors as discussed above. The additional loss term aims to restrict the range of the inputs to the activations functions by penalizing the values proportional to their distance outside the allowable range  $[z_{min}, z_{max}]$ . We achieve this by using a regularization loss,  $L_{activity}$ , defined by a shifted ReLU function.

$$L_{activity} = ReLU(z + z_{min}) + ReLU(z - z_{max}) \quad (3)$$

$$L_{total} = L_{model} + \lambda L_{activity} \quad (4)$$

where the Rectified Linear Unit,  $ReLU(x) = \max(0, x)$  and  $z$  is the concatenated array  $[z_0, z_1, z_2, z_3]$ . For values of  $z$  outside the allowable range, the loss adds a proportional penalizing term, thus restricting their values. The loss is differentiable with respect to the input except at the range extrema. However, the derivative at these extrema is approximated to be 0, making it a loss function compatible for backpropagation. We train the model with  $L_{activity}$  added as an additional term in the total loss  $L_{total}$  (see Eq. 4), where  $\lambda$  is used to weigh the activity regularization term.

### 5.2. Quantization of Inputs, Hidden States and Activations

In Stage II of training, we initialize the model with the trained weights from stage I and apply accelerator-aware training for the linear quantization of inputs and hidden states, and the non-linear quantization of the tanh and sigmoid activation functions.

#### 5.2.1. Bit-exact Quantization Replication during Forward Pass

- For the linear quantization of inputs and hidden states, we use Eq. 1 in the forward pass.

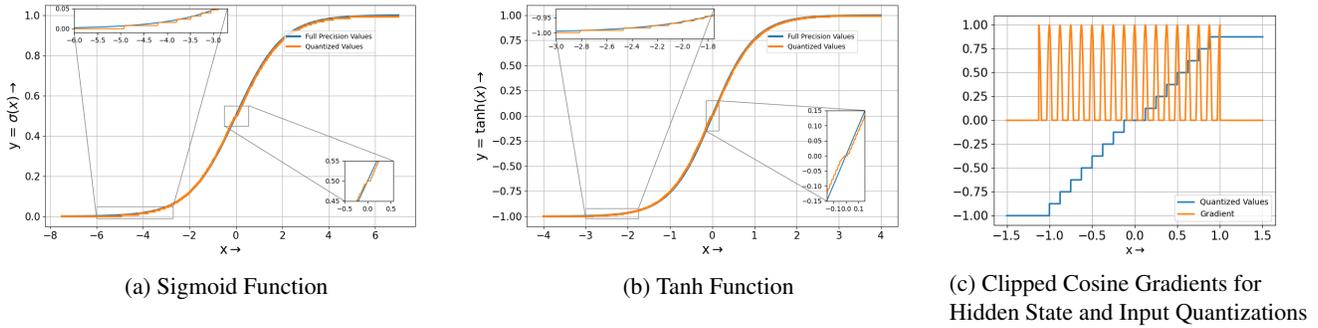


Fig. 3: Quantization of Sigmoid, Tanh, Hidden States and Inputs

---

**Algorithm 1** QAT for Inputs, Hidden States and Activations

---

**For Inputs and Hidden States**

*Forward pass:* Bit-exact linear quantization of input tensor  $X$

a.  $Y = Q(X)$

*Backward pass:* Clipped Cosine Gradients

a. Frequency,  $f_q = 1$

b.  $\frac{\partial Y}{\partial X} = \text{clip}(\cos(2\pi f_q X), \min = 0, \max = 1)$

{Gradients are clipped to 0 at the quantization bin midpoints}

**For Non-linear Activations**

*Forward pass:* Bit-exact non-linear quantization of input tensor  $X$

a.  $Y = Q_{act}(X)$

*Backward pass:* tanh and sigmoid gradients

a. For tanh,  $\frac{\partial Y}{\partial X} = 1 - \tanh^2(X)$

b. For sigmoid,  $\frac{\partial Y}{\partial X} = \text{sigmoid}(X) * (1 - \text{sigmoid}(X))$

---

- For the non-linear activation functions, the on-device linear interpolator in Fig. 3 is replicated during the forward pass. We denote this bit-exact quantization function implemented in the training workflow as  $Q_{act}$  (see Algorithm 1). Although we demonstrate it here for the *sigmoid* and *tanh* functions for a particular hardware implementation, this method can be generally applied to any other similar linearly or non-linearly quantized functions.

5.2.2. *Backpropagation Using Meaningful Gradients*

Since quantization is not a differentiable process, we employ the straight-through estimation method [28] for the forward pass and backpropagate meaningful gradients through the various quantization nodes in the backward pass.

- For the inputs and hidden states, we use clipped cosine gradients as illustrated in Fig. 3c. Instead of passing

a unity gradient throughout the range of the inputs at the quantization node as proposed in [28], we propose to backpropagate through the quantization node a periodic clipped cosine gradient to disincentivize values to occupy unwanted areas that lie outside the quantized levels. The clipped cosine gradients drive the values into surrounding quantization bins, while leaving the values around the center of the bins unchanged.

- For the activations, we use the full precision *sigmoid* and *tanh* gradients.

**6. EXPERIMENTAL SETUP**

We conduct our experiments under two setups A and B. In setup A, we train a small RNN-T model with 5 layers of an encoder and 2 layers of a decoder with an additive and feedforward joint network. The number of hidden units in the LSTM cell for setup A is 640, yielding a model with 26M parameters. Setup B is a larger RNN-T variant with the same number of encoder, decoder and joint network layers, in which the LSTMs have 1024 hidden units, yielding a model with 66M parameters. For both setups, the baseline models are trained for a total of 600k steps, with 5k steps of warmup to a learning rate (LR) of  $5e-4$  which was held constant for 150k steps, followed by an exponential decay to a learning rate of  $1e-5$  for the remaining steps. The Accelerator-aware training (AAT) models are trained in Stage I for 500k steps with  $\lambda = 2$  (see Eq. 4), and the same LR schedule as the baseline. This was followed by 5k-10k steps of training in Stage II with a constant learning rate of  $1e-4$ . The acoustic features are 64-dimensional Log Mel Filterbank Energies (LFBE) with a window size of 25 ms and 10 ms overlap between frames.

All models are trained with absolute cosine quantization aware regularization for 8-bit on-device quantization of weights [16]. To train models under both setups, we use an in-house collection of a far-field English training dataset with 270k hours of audio. For evaluation, we use 6 test sets: 3

for Setup A, and 3 for Setup B. Since Setup A has a smaller model, we evaluate it on a smaller number of intents. The number of utterances in each test set is available in Table 1.

## 7. RESULTS AND DISCUSSION

In Table 1, we compare the relative WER reduction (WERR) performance between the baseline and AAT models. Note that a positive WERR number indicates a WER improvement, whereas negative numbers signify WER degradation. Furthermore, in order to evaluate the effectiveness of each proposed stage, we present WER for each stage in Table 1.1. Here, Q denotes quantized WER numbers. Model A-I and A-2 represent the model in Setup A after Stage 1 and Stage II trainings, respectively. All WERR numbers are computed between the quantized (Q) version of the respective model and the un-quantized version of the baseline model. For example, the WERR of quantized Model A-II with respect to the un-quantized Baseline-A is +0.3%.

**Table 1:** Relative WER Reductions (WERR) for Setup A and B. WERR numbers are computed between the quantized (Q) version of the respective model and the un-quantized version of the baseline (i.e. Baseline (UN-Q) = 0.0).

**Table 1.1: Setup A**

Datasets (Num. Utts)	WERR (%)		
	Baseline (Q)	A-I (Q)	A-II (Q)
A-D1 (42688)	-9.9	-5.8	+0.3
A-D2 (26226)	-6.0	-3.0	+0.4
A-D3 (46208)	-8.0	-3.9	+0.8
<b>Stage I</b>	✗	✓	✓
<b>Stage II</b>	✗	✗	✓

**Table 1.2: Setup B**

Datasets (Num. Utts)	WERR (%)	
	Baseline (Q)	B-I (Q)
B-D1 (155936)	-5.4	-0.4
B-D2 (46530)	-7.2	-0.6
B-D3 (20279)	-3.0	+1.3
<b>Stage I</b>	✗	✓
<b>Stage II</b>	✗	✓

As shown in the Table 1, quantizing the inputs, hidden states and activations for the baseline models without the proposed AAT approach leads to 6-10% relative WER degradation for Setup A, and 3-7% for Setup B across our testing datasets. This degradation is reduced to be within 1% relative to the baseline performance after the two-stage AAT. In particular, as expected, we observed a larger degradation for the model in Setup A than the model in Setup B since the former has smaller number of parameters and is more susceptible to quantization than the latter. It is observed that by

using the two-stage approach, we can effectively reduce the WER gap between quantized and unquantized versions for small and large models alike.

**Table 2:** Normalized Latency Measurements for Setup A and B. All values are normalized with respect to the baseline latencies at p50, for the respective setups.

Setups	Statistic	Baseline		AAT	
		EL	UPL	EL	UPL
Setup A	<b>p50</b>	1.0	1.0	0.94	0.95
	<b>p90</b>	1.17	1.49	1.11	1.35
	<b>p99</b>	1.34	7.21	1.28	7.21
Setup B	<b>p50</b>	1.0	1.0	0.94	0.94
	<b>p90</b>	1.53	1.43	1.42	1.35
	<b>p99</b>	2.41	1.90	2.16	2.36

To demonstrate the latency gains, we provide the ASR engine latency (EL) and UPL numbers in Table 2. Here, EL measures the time elapsed between the user completing the utterance and the ASR recognition result being available in the ASR engine. UPL denotes the time elapsed between the user completing the utterance, and when the system responds with an appropriate dialogue / action. Thus, UPL includes EL plus the other server and device side processing required to fulfill the user’s request. We conduct on-device tests with 500 utterances for both baseline and AAT models in setups A and B, and report both EL and UPL numbers at the 50th, 90th and 99th percentiles. As before, positive relative numbers signify reductions in latency whereas negative numbers signify degradation. We see that by compiling the AAT models, we can get approximately 4 - 9% EL and UPL reductions at p50. AAT enables this by allowing us to shift to a more quantized, faster NNA data path during inference with negligible degradation in performance.

## 8. CONCLUSIONS

In this work, we introduce a two-stage AAT approach to alleviate the performance degradation due to performing post-training model conversion and quantized inference on the NNA. In particular, we incorporate the NNA quantization operators into model training and apply regularization as well as cosine gradients at the quantization nodes to reduce the performance gap between the in-training and post-training models. It is observed that with our proposed approach, there is little to no WER performance degradation between the unquantized baseline models and the quantized AAT counterparts. Compared to the original quantized baseline, the proposed two-stage quantized AAT models save 3-10% relative WER degradation. Furthermore, AAT enables NPU/hardware execution of activation functions in runtime, leading to 5-7% relative latency reductions.

## 9. REFERENCES

- [1] Liang Lu, Changliang Liu, Jinyu Li, and Yifan Gong, “Exploring transformers for large-scale speech recognition,” 2020.
- [2] Yongqiang Wang, Abdelrahman Mohamed, Due Le, Chunxi Liu, Alex Xiao, Jay Mahadeokar, Hongzhao Huang, Andros Tjandra, Xiaohui Zhang, Frank Zhang, and et al., “Transformer-based acoustic modeling for hybrid speech recognition,” *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2020.
- [3] Linhao Dong, Shuang Xu, and Bo 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)*, 2018, pp. 5884–5888.
- [4] Martin Radfar, Rohit Barnwal, Rupak Vignesh Swaminathan, Feng-Ju Chang, Grant P Strimel, Nathan Susanj, and Athanasios Mouchtaris, “Convtrnn-t: Convolutional augmented recurrent neural network transducers for streaming speech recognition,” *Proceedings of the Annual Conference of the International Speech Communication Association (Interspeech)*, 2022.
- [5] Yuan Shangguan, Jian Li, Qiao Liang, Raziell Alvarez, and Ian McGraw, “Optimizing speech recognition for the edge,” *arXiv preprint arXiv:1909.12408*, 2019.
- [6] A Graves, “Sequence transduction with recurrent neural networks,” in *Proceedings of the International Conference on Machine Learning (ICML)*, 2012.
- [7] Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber, “Lstm: A search space odyssey,” *IEEE transactions on neural networks and learning systems*, vol. 28, no. 10, pp. 2222–2232, 2016.
- [8] Tao Lei, Yu Zhang, Sida I Wang, Hui Dai, and Yoav Artzi, “Simple recurrent units for highly parallelizable recurrence,” *arXiv preprint arXiv:1709.02755*, 2017.
- [9] Sankaran Panchapagesan, Daniel S Park, Chung-Cheng Chiu, Yuan Shangguan, Qiao Liang, and Alexander Gruenstein, “Efficient knowledge distillation for rnn-transducer models,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5639–5643.
- [10] Rupak Vignesh Swaminathan, Brian King, Grant P Strimel, Jasha Droppo, and Athanasios Mouchtaris, “Codert: Distilling encoder representations with co-learning for transducer-based speech recognition,” *Proceedings of the Annual Conference of the International Speech Communication Association (Interspeech)*, 2021.
- [11] Varun Nagaraja, Yangyang Shi, Ganesh Venkatesh, Ozlem Kalinli, Michael L Seltzer, and Vikas Chandra, “Collaborative training of acoustic encoders for speech recognition,” *arXiv preprint arXiv:2106.08960*, 2021.
- [12] Kai Zhen, Hieu Duy Nguyen, Feng-Ju Chang, Athanasios Mouchtaris, and Ariya Rastrow, “Sparsification via compressed sensing for automatic speech recognition,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 6009–6013.
- [13] Jiahui Yu, Chung-Cheng Chiu, Bo Li, Shuo yiin Chang, Tara N. Sainath, Yanzhang He, Arun Narayanan, Wei Han, Anmol Gulati, Yonghui Wu, and Ruoming Pang, “Fastemit: Low-latency streaming asr with sequence-level emission regularization,” 2021.
- [14] Jon Macoskey, Grant P. Strimel, and Ariya Rastrow, “Bifocal neural asr: Exploiting keyword spotting for inference optimization,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5999–6003.
- [15] Jonathan Macoskey, Grant P Strimel, Jinru Su, and Ariya Rastrow, “Amortized neural networks for low-latency speech recognition,” *arXiv preprint arXiv:2108.01553*, 2021.
- [16] H D Nguyen, A Alexandridis, and A Mouchtaris, “Quantization aware training with absolute-cosine regularization for automatic speech recognition,” in *Proceedings of the Annual Conference of the International Speech Communication Association (INTERSPEECH)*, 2020.
- [17] Raziell Alvarez, Rohit Prabhavalkar, and Anton Bakhtin, “On the efficient representation and execution of deep acoustic models,” 2016.
- [18] 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.
- [19] Song Han, Huizi Mao, and William J. Dally, “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding,” 2016.

- [20] M. Zhu and S. Gupta, “To prune, or not to prune: exploring the efficacy of pruning for model compression,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [21] K Zhen, H D Nguyen, F.-J. Chang, A Mouchtaris, and A Rastrow, “Sparsification via compressed sensing for automatic speech recognition,” in *Proceedings of the International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021.
- [22] Hyeong-Ju Kang, “Accelerator-aware pruning for convolutional neural networks,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 7, pp. 2093–2103, 2019.
- [23] Jian Li and Raziq Alvarez, “On the quantization of recurrent neural networks,” *arXiv preprint arXiv:2101.05453*, 2021.
- [24] Ashkan Hosseinzadeh Namin, Karl Leboeuf, Roberto Muscedere, Huapeng Wu, and Majid Ahmadi, “Efficient hardware implementation of the hyperbolic tangent sigmoid function,” in *2009 IEEE International Symposium on Circuits and Systems*, 2009, pp. 2117–2120.
- [25] Che-Wei Lin and Jeen-Shing Wang, “A digital circuit design of hyperbolic tangent sigmoid function for neural networks,” in *2008 IEEE International Symposium on Circuits and Systems*, 2008, pp. 856–859.
- [26] Karl Leboeuf, Ashkan Hosseinzadeh Namin, Roberto Muscedere, Huapeng Wu, and Majid Ahmadi, “High speed vlsi implementation of the hyperbolic tangent sigmoid function,” in *2008 Third International Conference on Convergence and Hybrid Information Technology*, 2008, vol. 1, pp. 1070–1073.
- [27] Sepp Hochreiter and Jürgen Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [28] Yoshua Bengio, Nicholas Léonard, and Aaron Courville, “Estimating or propagating gradients through stochastic neurons for conditional computation,” 2013.