

# ADVANCING RNN TRANSDUCER TECHNOLOGY FOR SPEECH RECOGNITION

George Saon, Zoltán Tüske, Daniel Bolanos and Brian Kingsbury

IBM Research AI, Yorktown Heights, USA

## ABSTRACT

We investigate a set of techniques for RNN Transducers (RNN-Ts) that were instrumental in lowering the word error rate on three different tasks (Switchboard 300 hours, conversational Spanish 780 hours and conversational Italian 900 hours). The techniques pertain to architectural changes, speaker adaptation, language model fusion, model combination and general training recipe. First, we introduce a novel multiplicative integration of the encoder and prediction network vectors in the joint network (as opposed to additive). Second, we discuss the applicability of i-vector speaker adaptation to RNN-Ts in conjunction with data perturbation. Third, we explore the effectiveness of the recently proposed density ratio language model fusion for these tasks. Last but not least, we describe the other components of our training recipe and their effect on recognition performance. We report a 5.9% and 12.5% word error rate on the Switchboard and CallHome test sets of the NIST Hub5 2000 evaluation and a 12.7% WER on the Mozilla CommonVoice Italian test set.

**Index Terms**— End-to-end ASR, recurrent neural network transducer, multiplicative integration

## 1. INTRODUCTION

End-to-end approaches directly map an acoustic feature sequence to a sequence of characters or even words without any conditional independence assumptions. Compared to traditional approaches which integrate various knowledge sources in a complex search algorithm, end-to-end methods resulted in a dramatic simplification of both training and decoding pipelines. This led to a rapidly evolving research landscape in end-to-end modeling for ASR with Recurrent Neural Network Transducers (RNN-T) [1] and attention-based models [2, 3] being the most prominent examples. Attention based models are excellent at handling non-monotonic alignment problems such as translation [4], whereas RNN-Ts are an ideal match for the left-to-right nature of speech [5–17].

Nowadays, end-to-end models can reach unprecedented levels of speech recognition performance in spite of, or maybe *because of*, the significantly simpler implementations. It has been shown that, given enough training data, end-to-end models are clearly able to outperform traditional approaches [18]. Nevertheless, data sparsity and overfitting are inherent problems for any direct sequence-to-sequence model and various approaches have been proposed to introduce useful variances and mitigate these issues [19–23].

After a short overview of RNN-T sequence modeling (section 2), section 3 investigates an architectural change and section 4 presents efficient training and decoding recipes for such models. The proposed techniques are then evaluated on three different languages (section 5), before conclusions are drawn (section 6). As will be shown, the consistent application of these methods will result in remarkable end-to-end model performance, even if only a few hundred hours of training data are available.

## 2. RNN-T MODEL DESCRIPTION

Borrowing some notations from [1], RNN-Ts model the conditional distribution  $p(\mathbf{y}|\mathbf{x})$  of an output sequence  $\mathbf{y} = (y_1, \dots, y_U) \in \mathcal{Y}^*$  of length  $U$  given an input sequence  $\mathbf{x} = (x_1, \dots, x_T) \in \mathcal{X}^*$  of length  $T$ . The elements of  $\mathbf{x}$  are typically continuous multidimensional vectors whereas the elements of  $\mathbf{y}$  belong to an output space which is typically discrete.  $p(\mathbf{y}|\mathbf{x})$  is expressed as a sum over all possible alignments  $\mathbf{a} = (a_1, \dots, a_{T+U})$  that are consistent with  $\mathbf{y}$ :

$$p(\mathbf{y}|\mathbf{x}) = \sum_{\mathbf{a} \in \mathcal{B}^{-1}(\mathbf{y})} p(\mathbf{a}|\mathbf{x}) \quad (1)$$

The elements of  $\mathbf{a}$  belong to the augmented vocabulary  $\bar{\mathcal{Y}} = \mathcal{Y} \cup \{\phi\}$  where  $\phi$  (called BLANK) denotes the null output. The mapping  $\mathcal{B} : \bar{\mathcal{Y}}^* \rightarrow \mathcal{Y}^*$  is defined by  $\mathcal{B}(\mathbf{a}) = \mathbf{y}$ . For example, if  $\mathbf{y} = (C, A, T)$  and  $\mathbf{x} = (x_1, x_2, x_3, x_4)$  valid alignments include  $(\phi, C, \phi, A, \phi, T, \phi)$ ,  $(\phi, \phi, \phi, \phi, C, A, T)$ ,  $(C, A, T, \phi, \phi, \phi, \phi)$ , etc. Furthermore,  $p(\mathbf{a}|\mathbf{x})$  can be factorized as follows:

$$p(\mathbf{a}|\mathbf{x}) = p(\mathbf{a}|\mathbf{h}) \triangleq \prod_{i=1}^{T+U} p(a_i|h_{t_i}, \mathcal{B}(a_1, \dots, a_{i-1})) = \prod_{i=1}^{T+U} p(a_i|h_{t_i}, y_0, \dots, y_{u_{i-1}}) = \prod_{i=1}^{T+U} p(a_i|h_{t_i}, g_{u_i}) \quad (2)$$

where:  $\mathbf{h} = (h_1, \dots, h_T) = \text{Encoder}(\mathbf{x})$  is an embedding of the input sequence computed by an *encoder network*,  $\mathbf{g} = (g_1, \dots, g_U)$  is an embedding of the output sequence computed by a *prediction network* via the recursion  $g_u = \text{Prediction}(g_{u-1}, y_{u-1})$  (with the convention  $g_0 = \mathbf{0}, y_0 = \phi$ ), and  $p(a_i|h_{t_i}, g_{u_i})$  is the predictive output distribution over  $\bar{\mathcal{Y}}$  computed by a *joint network* which is commonly implemented as:

$$p(\cdot|h_t, g_u) = \text{softmax}[\mathbf{W}^{\text{out}} \tanh(\mathbf{W}^{\text{enc}} h_t + \mathbf{W}^{\text{pred}} g_u + b)] \quad (3)$$

$\mathbf{W}^{\text{enc}} \in \mathbb{R}^{J \times E}$ ,  $\mathbf{W}^{\text{pred}} \in \mathbb{R}^{J \times P}$  are linear projections that map  $h_t \in \mathbb{R}^E$  and  $g_u \in \mathbb{R}^P$  to a joint subspace which, after addition and  $\tanh$ , is mapped to the output space via  $\mathbf{W}^{\text{out}} \in \mathbb{R}^{|\bar{\mathcal{Y}}| \times J}$ . RNN-Ts are typically trained to minimize the negative log-likelihood (NLL) loss  $-\log p(\mathbf{y}|\mathbf{x})$ . From (1), calculating the sum by direct enumeration is intractable because the number of all possible alignments of length  $T + U$  is  $|\mathcal{B}^{-1}(\mathbf{y})| = \binom{T+U}{U} = \frac{(T+U)!}{T!U!} \geq (1 + \frac{T}{U})^U$  where  $\binom{n}{k}$  denotes the binomial coefficient *n choose k*. Luckily, the factorization in (2) allows for an efficient forward-backward algorithm with  $T \times U$  complexity for both loss and gradient computation [1, 24].

### 3. MULTIPLICATIVE INTEGRATION OF ENCODER AND PREDICTION NETWORK OUTPUTS

Here, we discuss a simple change to the joint network equation (3) and its implications:

$$p(\cdot|h_t, g_u) = \text{softmax}[\mathbf{W}^{out} \tanh(\mathbf{W}^{enc} h_t \odot \mathbf{W}^{pred} g_u + b)] \quad (4)$$

where  $\odot$  denotes elementwise multiplication (or Hadamard product). This modification was inspired by the work of [25] where the authors use multiplicative integration (MI) in the context of a recurrent neural network for fusing the information from the hidden state vector and the input vector. The advantages of MI over additive integration for fusing different information sources have also been mentioned in [26] and are summarized below.

**Higher-order interactions** MI allows for second-order interactions between the elements of  $h_t$  and  $g_u$  which facilitates modeling of more complex dependencies between the acoustic and LM embeddings. Importantly, the shift to second order is achieved with no increase in the number of parameters or computational complexity. In theory, feedforward neural nets with sufficient capacity and training data are universal function approximators; therefore, a joint network with  $(h_t, g_u)$  inputs and multiple layers should be optimal. However the memory requirements are prohibitive because the input tensor size for such a network is  $N \times T \times U \times (E + P)$  and the memory for the output tensors is  $N \times T \times U \times H \times L$  where  $N$  is the batchsize,  $H$  is the number of hidden units and  $L$  is the number of layers. Therefore, for pure training efficiency reasons, we are constrained to use shallow (single layer) joint networks, in which case the functional form of the layer becomes important.

**Generalizes additive integration** Indeed, when used in conjunction with biases, the additive terms appear in the expansion

$$\begin{aligned} (\mathbf{W}^{enc} h_t + b^{enc}) \odot (\mathbf{W}^{pred} g_u + b^{pred}) &= \mathbf{W}^{enc} h_t \odot \mathbf{W}^{pred} g_u \\ &+ b^{enc} \odot \mathbf{W}^{pred} g_u + b^{pred} \odot \mathbf{W}^{enc} h_t + b^{enc} \odot b^{pred} \end{aligned} \quad (5)$$

**Scaling and gating effect** By multiplying  $\tilde{h}_t := \mathbf{W}^{enc} h_t$  with  $\tilde{g}_u := \mathbf{W}^{pred} g_u$ , one embedding has a scaling effect on the other embedding. In particular, unlike for additive interactions, in MI the gradient with respect to one component is gated by the other component and vice versa. Concretely, let us denote by  $\mathcal{L}(\mathbf{x}, \mathbf{y}; \theta) = -\log p(\mathbf{y}|\mathbf{x}; \theta)$  the NLL loss for one utterance for an RNN-T with parameters  $\theta$ . The partial derivatives of the loss with respect to  $\tilde{h}_t$  and  $\tilde{g}_u$  are:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \tilde{h}_t} &= \tilde{g}_u \odot \frac{\partial \mathcal{L}}{\partial (\tilde{h}_t \odot \tilde{g}_u)} \\ \frac{\partial \mathcal{L}}{\partial \tilde{g}_u} &= \tilde{h}_t \odot \frac{\partial \mathcal{L}}{\partial (\tilde{h}_t \odot \tilde{g}_u)} \end{aligned} \quad (6)$$

These arguments suggest that multiplicative integration is a good candidate for fusing the two different information streams coming from the encoder and the prediction network. We expect multiplicative RNN-Ts to be, if not better, then at least complementary to additive RNN-Ts which should be beneficial for model combination.

### 4. TRAINING AND DECODING RECIPE

In this section, we discuss some of the techniques that are part of our RNN-T training and decoding recipe. While none of the techniques presented here are novel, the goal is to show their effectiveness for RNN-Ts in particular. The techniques can be roughly grouped into: (i) data augmentation/perturbation and model regularization, (ii) speaker adaptation and (iii) external language model fusion.

**Speed and tempo perturbation** [19] changes the rate of speech in the interval  $[0.9, 1.1]$  with or without altering the pitch or timbre of the speaker. This technique generates additional replicas of the training data depending on the number of speed and tempo perturbation values.

**Sequence noise injection** [20] adds, with a given probability, the downsampled spectra of randomly selected training utterances to the spectrum of the current training utterance. This technique does not increase the amount of training data per epoch.

**SpecAugment** [21] masks the spectrum of a training utterance with a random number of blocks of random size in both time and frequency.

**DropConnect** [22] zeros out entries randomly in the LSTM hidden-to-hidden transition matrices.

**Switchout** [23] randomly replaces labels in the output sequence with labels drawn uniformly from the output vocabulary.

**I-vector speaker adaptation** [27] appends a speaker identity vector to the input features coming from a given speaker. When used in conjunction with speed and tempo perturbation, the perturbed audio recordings of a speaker are considered to be different speakers for the purpose of universal background model (UBM) training, total variability matrix training, and i-vector extraction.

**Alignment-length synchronous decoding** [17] is a beam search technique with the property that all competing hypotheses within the beam have the same alignment length. It has been shown to be faster than time-synchronous search for the same accuracy.

**Density ratio (DR) LM fusion** [13] is a shallow fusion technique that combines two language models: an external LM trained on a target domain corpus and a language model trained on the acoustic transcripts (source domain) only. The latter is used to subtract the effect of the intrinsic LM given by the prediction network (idea further developed in [14]). Decoding using DR fusion is done according to:

$$\mathbf{y}^* = \underset{\mathbf{y} \in \mathcal{Y}^*}{\text{argmax}} \{ \log p(\mathbf{y}|\mathbf{x}) - \mu \log p^{src}(\mathbf{y}) + \lambda \log p^{ext}(\mathbf{y}) + \rho |\mathbf{y}| \} \quad (7)$$

where  $\mu, \lambda, \rho$  are the weights corresponding to the source LM  $p^{src}$ , external LM  $p^{ext}$ , and label length reward  $|\mathbf{y}|$ .

## 5. EXPERIMENTS AND RESULTS

We investigate the effectiveness of the proposed techniques on one public corpus (English conversational telephone speech 300 hours) and two internal datasets (Spanish and Italian conversational speech 780 hours and 900 hours, respectively).

### 5.1. Experiments on Switchboard 300 hours

The first set of experiments was conducted on the Switchboard speech corpus which contains 300 hours of English telephone conversations between two strangers on a preassigned topic. The acoustic data segmentation and transcript preparation are done according

to the Kaldi `s5c` recipe [28]. We report results on the commonly used Hub5 2000 Switchboard and CallHome test sets as well as Hub5 2001 and RT'03, which are processed according to the LDC segmentation and scored using Kaldi scoring for measuring WER.

We extract 40-dimensional speaker independent log-Mel filterbank features every 10 ms which are mean and variance normalized per conversation side. The features are augmented with  $\Delta$  and  $\Delta\Delta$  coefficients and every two consecutive frames are stacked and every second frame is skipped resulting in 240-dimensional vectors extracted every 20 ms. These features are augmented with 100-dimensional i-vectors that are extracted using a 2048 40-dimensional diagonal covariance Gaussian mixture UBM trained on speaker independent PLP features transformed with LDA and a semi-tied covariance transform.

Speed and tempo perturbation is applied to each conversation side with values in  $\{0.9, 1.1\}$  for both speed and tempo separately resulting in 4 additional training data replicas which, together with the original data, amounts to 1500 hours of training data per epoch. For sequence noise injection, we add, with probability 0.8, to the spectrum of each training utterance the spectrum of one random utterance of similar length scaled by a factor of 0.4. For SpecAugment we used the settings published in [21]. Lastly, in label switchout, for a sequence of length  $U$ , we first sample  $\hat{n} \in \{0, \dots, U\}$  with  $p(\hat{n}) \propto e^{-\hat{n}/\tau}$  and then we replace, with probability  $\hat{n}/U$ , the true characters with random characters for each position in the sequence. We set the temperature to  $\tau = 10$  in our experiments.

The architecture of the trained models is as follows. The encoder contains 6 bidirectional LSTM layers with 640 cells per layer per direction and is initialized with a network trained with CTC based on [29] similar to [5, 6, 8, 15]. The prediction network is a single unidirectional LSTM layer with only 768 cells (this size has been found to be optimal after external LM fusion). The joint network projects the 1280-dimensional stacked encoder vectors from the last layer and the 768-dimensional prediction net embedding to 256 dimensions and combines the projected vectors using either  $+$  or  $\odot$ . After the application of hyperbolic tangent, the output is projected to 46 logits followed by a softmax layer corresponding to 45 characters plus BLANK. All models have 57M parameters.

Optimizer	LR policy	Batch size	SWB	CH
Momentum SGD	const+decay	256	9.6	17.5
Momentum SGD	const+decay	64	9.4	17.5
AdamW	const+decay	256	9.4	17.8
AdamW	const+decay	64	8.5	16.7
AdamW	OneCycleLR	64	8.1	16.5

**Table 1.** Effect of optimizer, batch size and learning rate schedule on recognition performance for Switchboard 300 hours (Hub5'00 test set).

The models were trained in Pytorch on V100 GPUs for 20 epochs using SGD variants that differ in optimizer, learning rate schedule and batch size. In Table 1 we look at the effect on performance of changing the learning strategy. For momentum SGD, the learning rate was set to 0.01 and decays geometrically by a factor of 0.7 every epoch after epoch 10. For AdamW, the maximum learning rate is set to  $5e-4$  and the OneCycleLR policy [30] consists in a linear warmup phase from  $5e-5$  to  $5e-4$  over the first 6 epochs followed by a linear annealing phase to 0 for the next 14 epochs. As can be seen, AdamW with OneCycleLR scheduling and a batch size of 64 appears to be optimal in this experiment.

Model	No ext LM			With ext LM		
	SWB	CH	Avg	SWB	CH	Avg
Baseline	7.9	15.7	11.8	6.4	13.4	9.9
No Switchout	8.1	15.5	11.8	6.3	13.1	9.7
No Seq. noise	8.4	16.1	12.2	6.6	13.8	10.2
No i-vectors	8.1	16.0	12.0	6.6	13.9	10.3
No CTC init.	8.3	16.3	12.3	6.6	14.1	10.4
No Density ratio	7.9	15.7	11.8	6.9	14.4	10.7
No DropConnect	8.2	16.7	12.5	7.0	14.8	10.9
No SpecAugment	8.8	18.1	13.5	6.9	15.5	11.2
No Speed/tempo	10.0	18.3	14.2	7.6	15.2	11.4

**Table 2.** Ablation study on Switchboard 300 hours (Hub5'00 test set).

Next, we perform an ablation study on the final recipe to tease out the importance of the individual components. We show results in Table 2 for models with multiplicative integration with and without external language model fusion. The source LM is a one-layer LSTM with 768 cells (5M parameters) trained on the Switchboard 300 hours character-level transcripts (15.4M tokens) whereas the target LM is a two-layer LSTM with 2048 cells per layer (84M parameters) trained on the Switchboard+Fisher character-level transcripts (126M tokens).

Surprisingly, deactivating switchout from the final recipe actually improves recognition performance after external LM fusion, which was not the case in prior experiments where this technique was marginally helpful. Also, another unexpected finding was the large gain due to density ratio LM fusion. We attribute this to optimizing the other elements of the training recipe around this technique (e.g. reducing the size of the prediction network).

ext LM	Model	Hub5'00		Hub5'01			RT'03	
		swb	ch	swb	s2p3	s2p4	swb	fsh
no	$+$	8.0	15.6	8.7	11.4	16.3	18.0	11.4
	$\odot$	8.1	15.5	8.5	11.7	15.9	18.5	11.8
	Comb.	7.5	14.3	7.9	10.8	15.2	17.1	10.7
yes	$+$	6.4	13.4	7.0	9.2	13.4	15.0	9.2
	$\odot$	6.3	13.1	7.1	9.4	13.6	15.4	9.5
	Comb.	<b>5.9</b>	<b>12.5</b>	<b>6.6</b>	<b>8.6</b>	<b>12.6</b>	<b>14.1</b>	<b>8.6</b>

**Table 3.** Recognition results for additive, multiplicative and combined RNN-Ts on Switchboard 300 hours (Hub5'00, Hub5'01, RT'03 test sets).

In the next experiment, we compare RNN-Ts with additive versus multiplicative integration in the joint network. Based on the previous findings, we train these models without switchout. We also perform log-linear model combination with density ratio LM fusion according to:

$$\mathbf{y}^* = \underset{\mathbf{y} \in \mathcal{H}_+(\mathbf{x}) \cup \mathcal{H}_\odot(\mathbf{x})}{\operatorname{argmax}} \left\{ \alpha \log p(\mathbf{y}|\mathbf{x}; \theta_+) + \beta \log p(\mathbf{y}|\mathbf{x}; \theta_\odot) - \mu \log p^{\text{src}}(\mathbf{y}) + \lambda \log p^{\text{ext}}(\mathbf{y}) + \rho |\mathbf{y}| \right\} \quad (8)$$

where  $\mathcal{H}_+(\mathbf{x}), \mathcal{H}_\odot(\mathbf{x})$  are the  $n$ -best hypotheses generated by the additive and multiplicative RNN-Ts. Concretely, we rescore the union of the top 32 hypotheses from each model with  $\alpha = \beta = \mu =$

0.5,  $\lambda = 0.7$  and  $\rho = 0.2$ . The results from Table 3 show that multiplicative RNN-Ts are comparable on Hub5'00 and Hub5'01 but slightly worse on RT'03 and that model combination significantly improves the recognition performance across all test sets.

For comparison, we include in Table 4 the top performing single model systems from the literature on the Switchboard 300 hours corpus. The numbers should be compared with 6.3/13.1 (row 5 from Table 3). As can be seen, the proposed modeling techniques exhibit excellent performance on this task.

System	Type	ext. LM	SWB	CH
Park et al.'19 [21]	Att. enc-dec	LSTM	6.8	14.1
Irie et al.'19 [31]	Hybrid	Transf.	6.7	12.9
Hadian et al.'18 [32]	LF-MMI	RNN	7.5	14.6
Tüske et al.'20 [33]	Att. enc-dec	LSTM	6.4	12.5

**Table 4.** Single model performance for existing systems on Switchboard 300 hours (Hub5'00 test set).

## 5.2. Experiments on conversational Spanish 780 hours

We also investigated the effectiveness of the proposed techniques on an internal dataset which consists of 780 hours of Spanish call center data collected using a balanced distribution of Castilian and other Central and South American dialects from roughly 4000 speakers. The test set on which we report results comes from the call center domain and contains 4 hours of speech, 113 speakers and 31.6k words.

The model architecture, training and decoding recipes are identical to the ones from the previous subsection (except for a 40-character output layer) with the difference that we do not use i-vector speaker adaptation because of the small amount of data per speaker. The external language model is a two layer LSTM with 2048 cells per layer (84M parameters) trained on the character-level acoustic transcripts (36M characters) as well as additional text data from the customer care domain (173M characters). The source LM is a single-layer LSTM with 1024 cells (8.7M parameters).

Model/technique	Training data	WER
Initial experiment	250 h	34.8
+ DropConnect+seq. noise	250 h	27.6
+ Speed/tempo	780 h	25.0
+ CTC encoder pretraining	780 h	23.6
+ Multiplicative integration	780 h	22.7
+ SpecAugment	780 h	21.9
+ AdamW+OneCycleLR	780 h	20.8
+ Shallow LM fusion	780 h	20.3
+ Density ratio LM fusion	780 h	20.0

**Table 5.** Cumulative effect of proposed techniques on conversational Spanish 780 hours (internal test set).

In Table 5, we look at the impact of the various techniques on word error rate for bidirectional RNN-Ts. The first 4 models have additive integration and the first 6 models are trained with momentum SGD and a constant+decay learning rate schedule for 20 epochs. Unlike in the previous experiments, here we observe a significant gain from multiplicative integration (0.9% absolute) and a smaller gain from density ratio fusion (0.3% absolute).

## 5.3. Experiments on conversational Italian 900 hours

The last set of experiments was conducted on an internal Italian corpus of conversational speech that was collected using scripted dialogs from several domains such as banking, insurance, telco, retail to name a few. The data has a balanced demographic and dialectal distribution with a large fraction of the speakers speaking standard dialect (as spoken in the capital or news broadcasts). We report results on the Mozilla CommonVoice<sup>1</sup> Italian test set which has 1.2 hours of audio, 764 utterances and 6.6K reference words.

Similar to the Spanish setup, we trained character-level RNN-Ts (53 outputs) without i-vector speaker adaptation using the recipe described in section 4 with specifics from subsection 5.1. In addition, we also train RNN-Ts with unidirectional LSTM encoders with 6 layers and 1024 cells per layer on stacked log-Mel features augmented with  $\Delta$ ,  $\Delta\Delta$  with 5 frames lookahead as proposed in [10]. The language model configuration for density ratio fusion is as follows: the source LM has one LSTM layer with 1024 units (8.7M parameters) and is trained on the character-level acoustic transcripts (41.7M characters) whereas the external LM is a two layer 1024 cells/layer LSTM (21.3M parameters) trained on the acoustic transcripts and 10% of Italian Wikipedia data (306M characters).

In Table 6, we compare additive versus multiplicative models with bidirectional and unidirectional encoders with and without external LM fusion (regular and density ratio). Based on these results, three observations can be made. First, multiplicative integration significantly outperforms the additive counterpart for unidirectional models and after shallow LM fusion (for both). Second, there is a severe degradation in recognition performance from using unidirectional encoders which can be mitigated with techniques from [12, 34]. Third, density ratio LM fusion significantly outperforms regular shallow fusion which we attribute to the mismatched training and testing conditions.

	Bidirectional		Unidirectional	
	+	⊖	+	⊖
No external LM	17.8	17.6	28.0	26.2
Shallow fusion	15.2	13.9	22.9	21.4
Density ratio fusion	13.6	12.7	20.9	18.6

**Table 6.** Various comparisons on conversational Italian 900 hours (Mozilla CommonVoice test set).

## 6. CONCLUSION

The contribution of this paper is twofold. First, we have introduced a simple yet effective modification of the joint network whereby we combine the encoder and prediction network embeddings using multiplicative integration instead of additive. MI outperforms additive integration on two out of three tasks and shows good model complementarity on the Switchboard 300 hours corpus. Second, we have shown that careful choice of optimizer and learning rate scheduling, data augmentation/perturbation and model regularization, speaker adaptation, external language model fusion, and model combination can lead to excellent recognition performance when applied to models with a very simple architecture and character outputs. Future work will look at alternative neural transducer architectures (e.g. [35]) and training criteria (e.g. [36]) and attempt to simplify the training recipe without sacrificing recognition performance.

<sup>1</sup><https://commonvoice.mozilla.org>

## References

- [1] A. Graves, “Sequence transduction with recurrent neural networks,” *arXiv preprint arXiv:1211.3711*, 2012.
- [2] D. Bahdanau, J. Chorowski, D. Serdyuk, et al., “End-to-end attention-based large vocabulary speech recognition,” in *ICASSP*. IEEE, 2016, pp. 4945–4949.
- [3] W. Chan, N. Jaitly, Q. Le, and O. Vinyals, “Listen, attend and spell: A neural network for large vocabulary conversational speech recognition,” in *ICASSP*. IEEE, 2016, pp. 4960–4964.
- [4] Y. Wu, M. Schuster, Z. Chen, et al., “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *arXiv preprint arXiv:1609.08144*, 2016.
- [5] A. Graves, A.-r. Mohamed, and G. Hinton, “Speech recognition with deep recurrent neural networks,” in *ICASSP*. IEEE, 2013, pp. 6645–6649.
- [6] K. Rao, H. Sak, and R. Prabhavalkar, “Exploring architectures, data and units for streaming end-to-end speech recognition with RNN-transducer,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2017, pp. 193–199.
- [7] E. Battenberg, J. Chen, R. Child, et al., “Exploring neural transducers for end-to-end speech recognition,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2017, pp. 206–213.
- [8] S. Wang, P. Zhou, W. Chen, et al., “Exploring RNN-Transducer for Chinese speech recognition,” *arXiv preprint arXiv:1811.05097*, 2018.
- [9] Y. He, T. N. Sainath, R. Prabhavalkar, et al., “Streaming end-to-end speech recognition for mobile devices,” in *ICASSP*. IEEE, 2019, pp. 6381–6385.
- [10] J. Li, R. Zhao, H. Hu, and Y. Gong, “Improving RNN transducer modeling for end-to-end speech recognition,” *arXiv preprint arXiv:1909.12415*, 2019.
- [11] A. Tripathi, H. Lu, H. Sak, and H. Soltau, “Monotonic recurrent neural network transducer and decoding strategies,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 944–948.
- [12] M. Jain, K. Schubert, J. Mahadeokar, et al., “RNN-T for latency controlled ASR with improved beam search,” *arXiv preprint arXiv:1911.01629*, 2019.
- [13] E. McDermott, H. Sak, and E. Variani, “A density ratio approach to language model fusion in end-to-end automatic speech recognition,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 434–441.
- [14] E. Variani, D. Rybach, C. Allauzen, and M. Riley, “Hybrid autoregressive transducer (HAT),” in *ICASSP*. IEEE, 2020, pp. 6139–6143.
- [15] A. Zeyer, A. Merboldt, R. Schlüter, and H. Ney, “A new training pipeline for an improved neural transducer,” *arXiv preprint arXiv:2005.09319*, 2020.
- [16] C.-C. Chiu, A. Narayanan, W. Han, et al., “RNN-T models fail to generalize to out-of-domain audio: Causes and solutions,” *arXiv preprint arXiv:2005.03271*, 2020.
- [17] G. Saon, Z. Tüske, and K. Audhkhasi, “Alignment-length synchronous decoding for RNN transducer,” in *ICASSP*. IEEE, 2020, pp. 7804–7808.
- [18] C.-C. Chiu, T. N. Sainath, Y. Wu, et al., “State-of-the-art speech recognition with sequence-to-sequence models,” in *ICASSP*. IEEE, 2018, pp. 4774–4778.
- [19] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, “Audio augmentation for speech recognition,” in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.
- [20] G. Saon, Z. Tüske, K. Audhkhasi, and B. Kingsbury, “Sequence noise injected training for end-to-end speech recognition,” in *ICASSP*. IEEE, 2019.
- [21] D. S. Park, W. Chan, Y. Zhang, et al., “SpecAugment: A simple data augmentation method for automatic speech recognition,” *Proc. Interspeech 2019*, pp. 2613–2617, 2019.
- [22] L. Wan, M. Zeiler, S. Zhang, et al., “Regularization of neural networks using DropConnect,” in *International conference on machine learning*, 2013, pp. 1058–1066.
- [23] X. Wang, H. Pham, Z. Dai, and G. Neubig, “Switchout: an efficient data augmentation algorithm for neural machine translation,” *arXiv preprint arXiv:1808.07512*, 2018.
- [24] T. Bagby, K. Rao, and K. C. Sim, “Efficient implementation of recurrent neural network transducer in TensorFlow,” in *Spoken Language Technology Workshop (SLT)*. IEEE, 2018, pp. 506–512.
- [25] Y. Wu, S. Zhang, Y. Zhang, et al., “On multiplicative integration with recurrent neural networks,” in *Advances in neural information processing systems*, 2016, pp. 2856–2864.
- [26] S. M. Jayakumar, W. M. Czarnecki, J. Menick, et al., “Multiplicative interactions and where to find them,” in *International Conference on Learning Representations*, 2019.
- [27] G. Saon, H. Soltau, D. Nahamoo, and M. Picheny, “Speaker adaptation of neural network acoustic models using i-vectors,” in *Proc. ASRU*, 2013.
- [28] D. Povey, A. Ghoshal, G. Boulianne, et al., “The Kaldi speech recognition toolkit,” in *IEEE workshop on automatic speech recognition and understanding*. IEEE Signal Processing Society, 2011.
- [29] K. Audhkhasi, G. Saon, Z. Tüske, et al., “Forget a bit to learn better: Soft forgetting for CTC-based automatic speech recognition,” in *INTERSPEECH*, 2019, pp. 2618–2622.
- [30] L. N. Smith and N. Topin, “Super-convergence: Very fast training of neural networks using large learning rates,” in *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*. International Society for Optics and Photonics, 2019, vol. 11006, p. 1100612.
- [31] K. Irie, A. Zeyer, R. Schlüter, and H. Ney, “Training language models for long-span cross-sentence evaluation,” in *Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 419–426.
- [32] H. Hadian, H. Sameti, D. Povey, and S. Khudanpur, “End-to-end speech recognition using lattice-free MMI,” in *Proceedings of INTER-SPEECH*, 2018.
- [33] Z. Tüske, G. Saon, K. Audhkhasi, and B. Kingsbury, “Single headed attention based sequence-to-sequence model for state-of-the-art results on Switchboard-300,” in *INTERSPEECH*, 2020.
- [34] G. Kurata and G. Saon, “Knowledge distillation from offline to streaming RNN transducer for end-to-end speech recognition,” in *INTERSPEECH*, 2020.
- [35] C.-F. Yeh, J. Mahadeokar, K. Kalgaonkar, et al., “Transformer-transducer: End-to-end speech recognition with self-attention,” *arXiv preprint arXiv:1910.12977*, 2019.
- [36] C. Weng, C. Yu, J. Cui, et al., “Minimum Bayes risk training of RNN-Transducer for end-to-end speech recognition,” *arXiv preprint arXiv:1911.12487*, 2019.