IMPROVING RNN TRANSDUCER BASED ASR WITH AUXILIARY TASKS

Chunxi Liu, Frank Zhang, Duc Le, Suyoun Kim, Yatharth Saraf, Geoffrey Zweig

Facebook AI, USA

{chunxiliu, frankz, duchoangle, suyounkim, ysaraf, gzweig}@fb.com

ABSTRACT

End-to-end automatic speech recognition (ASR) models with a single neural network have recently demonstrated state-of-the-art results compared to conventional hybrid speech recognizers. Specifically, recurrent neural network transducer (RNN-T) has shown competitive ASR performance on various benchmarks. In this work, we examine ways in which RNN-T can achieve better ASR accuracy via performing auxiliary tasks. We propose (i) using the same auxiliary task as primary RNN-T ASR task, and (ii) performing context-dependent graphemic state prediction as in conventional hybrid modeling. In transcribing social media videos with varying training data size, we first evaluate the streaming ASR performance on three languages: Romanian, Turkish and German. We find that both proposed methods provide consistent improvements. Next, we observe that both auxiliary tasks demonstrate efficacy in learning deep transformer encoders for RNN-T criterion, thus achieving competitive results -2.0%/4.2% WER on LibriSpeech test-clean/other - as compared to prior top performing models.

Index Terms— recurrent neural network transducer, speech recognition, auxiliary learning

1. INTRODUCTION

Building conventional hidden Markov model (HMM) based hybrid automatic speech recognition (ASR) systems include multiple engineered steps like bootstrapping, decision tree clustering of context-dependent phonetic/graphemic states [1], acoustic and language model training, etc. End-to-end ASR models [2, 3, 4, 5] use neural networks to transduce audio into word sequences, and can be learned in a single step from scratch. Specifically, recurrent neural network transducer (RNN-T) originally presented in [2] also referred to as sequence transducer - has been shown preferable on numerous applications. For example, the model size of RNN-T is much more compact than conventional hybrid models, being favorable as an on-device recognizer [6, 7, 8]. It also has been demonstrated as a high-performing streaming model in extensive benchmarks [9, 10, 11, 12]. Such recent success has motivated the efforts to improve RNN-T from various aspects, e.g. model pretraining [13, 14], generalization ability on long-form audios [15], training algorithms [7, 16], speech enhancement [17], etc.

In this work, we make an attempt on improving RNN-T via auxiliary learning, which aims to improve the generalization ability of a primary task by training on additional auxiliary tasks [18, 19]. While multitask learning [20] may aim to improve the performance of multiple tasks simultaneously, auxiliary learning selectively serves to assist the primary task and only the primary task performance is in focus. Auxiliary learning has been studied extensively in reinforcement learning [18, 21], where pseudo-reward functions are designed to enable the main policy to be learned more efficiently. In the context of attention-based sequence-to-sequence (seq2seq) ASR models, [22, 23] show that learning encoders with auxiliary tasks of predicting phonemes or context-dependent phonetic HMM states (i.e. senones [24]) can improve the primary ASR word error rate (WER). [25] shows that using auxiliary syntactic and semantic tasks can improve the main low-resource machine translation task.

In this paper, we consider the application of auxiliary tasks to RNN-T based ASR. First, we design an auxiliary task to be the same ASR task, where the transducer encoder forks from an intermediate encoder layer, and both the primary branch and auxiliary branch perform ASR tasks. Note that in this way, both primary and auxiliary branches can provide posterior distributions over output labels - characters or wordpieces. Inspired by the prior works [26, 27, 28], we exploit a symmetric Kullback-Leibler (KL) divergence loss between the output posterior distributions of primary and auxiliary branches, along with the standard RNN-T loss. Such mutual KL divergence loss is expected to implicitly penalize the inconsistent gradients from the primary and auxiliary losses with respect to their shared parameters, and relieve the optimization inconsistency across tasks [28]. Overall, the knowledge distilled from auxiliary tasks help a model learn better representations shared between primary and auxiliary branches, by enabling the model to find a more robust (flatter) minima and to better generalize to test data [26].

Secondly, we propose an alternative auxiliary task of predicting context-dependent graphemic states, also referred to as chenones [29], as in standard HMM-based hybrid modeling. Similar to the auxiliary senone classification for improving attention-based seq2seq model [22, 23], we exploit chenone prediction for improving RNN-T without relying on language-specific phonemic lexicon. HMM-based graphemic hybrid ASR systems have been shown to achieve comparable performance to phonetic lexicon based approaches [30, 31, 29], and still demonstrate state-of-the-art results on common benchmarks when compared to end-to-end models [32]. In this paper, we examine if the context-dependent graphemic knowledge - from a decision tree clustering of tri-grapheme HMM states - can be complementary to the character or wordpiece (i.e. subword unit) modeling used in end-to-end ASR [13], and if the auxiliary chenone prediction task provides an avenue of distilling such context-dependent graphemic knowledge into RNN-T training by providing additional discriminative information.

To evaluate our proposed methods, we first use streamable ASR models on a challenging task of transcribing social media videos, in both low-resource (training data size ~160 hours) and medium-resource (~3K hours) conditions. Next, on LibriSpeech, we consider the application of auxiliary tasks to the sequence transducers built with deep transformer encoders.

2. MODELING APPROACHES

In this section we begin with a review of RNN-T based ASR, as originally presented in [2]. Then we present our proposed auxiliary RNN-T task. Lastly, we describe the auxiliary context-dependent graphemic state prediction task.

2.1. RNN-T

ASR can be formulated as a sequence-to-sequence problem. Each speech utterance is parameterized as an input acoustic feature vector sequence $\mathbf{x} = \{\mathbf{x}_1 \dots \mathbf{x}_T\} = \mathbf{x}_{1:T}$, where $\mathbf{x}_t \in \mathbb{R}^d$ and T is the number of frames. Denote a grapheme set or a wordpiece inventory as \mathcal{Y} , and the corresponding output sequence of length U as $\mathbf{y} = \{y_1 \dots y_U\} = \mathbf{y}_{1:U}$, where $y_u \in \mathcal{Y}$.

We define $\overline{\mathcal{Y}}$ as $\mathcal{Y} \cup \{\emptyset\}$, where \emptyset is the blank label. Denote $\overline{\mathcal{Y}}^*$ as the set of all sequences over output space $\overline{\mathcal{Y}}$, and the element $\mathbf{a} \in \overline{\mathcal{Y}}^*$ as an alignment sequence. Then we have the posterior probability as:

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

where $\mathcal{B}: \bar{\mathcal{Y}}^* \to \mathcal{Y}^*$ is a function that removes blank symbols from an alignment **a**. RNN-T model parameterizes the alignment probability $P(\mathbf{a}|\mathbf{x})$ and computes it with an encoder network (i.e. transcription network in [2]), a prediction network and a joint network. The encoder performs a mapping operation, denoted as f^{enc} , which converts **x** into another sequence of representations $\mathbf{h}_{1:T}^{\text{enc}} = {\mathbf{h}_{1}^{\text{enc}} \dots \mathbf{h}_{T}^{\text{enc}}}$:

$$\mathbf{h}_{1:T}^{\text{enc}} = f^{\text{enc}}(\mathbf{x}) \tag{2}$$

A prediction network f^{pred} , based on RNN or its variants, takes both its state vector and the previous non-blank output label y_{u-1} predicted by the model, to produce the new representation $\mathbf{h}_{u}^{\text{pred}}$:

$$\mathbf{h}_{1:u}^{\text{pred}} = f^{\text{pred}}(y_{0:(u-1)}) \tag{3}$$

where u is output label index and $y_0 = \emptyset$. The joint network f^{join} is a feed-forward network that combines encoder output $\mathbf{h}_t^{\text{enc}}$ and prediction network output $\mathbf{h}_u^{\text{pred}}$ to compute logits $\mathbf{z}_{t,u}$:

$$\mathbf{z}_{t,u} = f^{\text{join}}(\mathbf{h}_t^{\text{enc}}, \mathbf{h}_u^{\text{pred}})$$
(4)

$$P(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) = \operatorname{Softmax}(\mathbf{z}_{t,u})$$
(5)

such that the logits go through a softmax function and produce a posterior¹ distribution of the next output label over $\overline{\mathcal{Y}}$. Finally, the RNN-T loss function is then the negative log posterior as in Eq. 1:

$$\mathcal{L}^{\text{RNN-T}}(\theta) = -\log P(\mathbf{y}|\mathbf{x})$$
(6)

where θ denotes the model parameters. Note that the encoder is analogous to an acoustic model, and the combination of prediction network and joint network can be seen as a decoder.

2.2. Auxiliary sequence transducer modeling

The RNN-T decoder can be viewed as a RNN language model. The RNN takes both its state vector and y_{u-1} to predict y_u , so implicitly predicting y_u is conditioned on the whole label history $y_1 \dots y_{u-1}$ as in Eq. 5. Since the label history can be very informative in predicting



Fig. 1. Illustration of the proposed auxiliary RNN-T and KL divergence criteria. For the auxiliary criteria, decoder is shown in a dashed box when it is used by the auxiliary branch to compute the logits (Eq. 7) in the forward pass, while the decoder is not updated in the backward pass.

the next output label, we conjecture that the posterior entropy over $\bar{\mathcal{Y}}$ computed by Eq. 5 may be excessively reduced, resulting in encoder undertraining. In other words, if the decoder has played a major role in predicting each y_u by such teacher forcing procedure, which can still result in a reasonable training loss, the encoder may underfit the input **x**, and the resulting generalization can be worse than a model with an adequately trained encoder.

Additionally, gradient flow [34] through a deep neural network architecture is difficult in general, due to the gradient vanishing/exploding problem at lower layers. Although we could add shortcut connections [35, 36] across encoder layers that would help gradient flow through the encoder, it does not address the encoder undertraining problem – if the posterior of Eq. 5 has been peaked at the true label due to the strong cue from previous label history.

2.2.1. Auxiliary RNN-T criterion

An alternative proposal to increase the gradient signal is based on connecting auxiliary classifiers to intermediate layers directly [37]. In this work, to address encoder underfitting and provide the encoder with more backward gradients, we take the approach of connecting an auxiliary branch to an intermediate encoder layer and applying the same RNN-T loss function.

As in Figure 1, given an *L*-layer encoder network, denote $\mathbf{h}^{\text{enc},l}$ as the hidden activations of an intermediate layer *l*, where $1 \leq l < L$. $\mathbf{h}^{\text{enc},l}$ goes through a one-hidden-layer multi-layer perceptron (MLP), parameterized by ϕ^l , and use the same decoder to compute the logits of auxiliary branch:

$$\mathbf{z}_{t,u}^{\text{aux},l} = f^{\text{join}}(\text{MLP}(\mathbf{h}_t^{\text{enc},l}), \mathbf{h}_u^{\text{pred}})$$
(7)

$$P_{\text{aux},l}(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) = \text{Softmax}(\mathbf{z}_{t,u}^{\text{aux},l})$$
(8)

such that we can apply another RNN-T objective function to this auxiliary branch, and the overall objective function becomes:

$$\mathcal{L}(\theta, \phi) = \mathcal{L}^{\text{RNN-T}}(\theta) + \lambda_{\text{aux}} \mathcal{L}^{\text{RNN-T}}(\theta_{\text{shared}}^{\text{enc}}, \phi)$$

= $-\log P(\mathbf{y}|\mathbf{x}) - \lambda_{\text{aux}} \log P_{\text{aux},l}(\mathbf{y}|\mathbf{x})$ (9)

where θ denotes the parameters of primary branch including the whole encoder and decoder, and $\theta_{\text{shared}}^{\text{snc}}$ denotes the encoder layers 1 - *l* shared

¹Note that, the posterior distribution in Eq. 5 can also be written as $P(y_u|\mathbf{x}_{1:T}, y_{1:(u-1)})$, if the encoder uses global/infinite context, like a BLSTM or non-streaming transformer network [32, 33].

by primary and auxiliary branches, and λ_{aux} is a weighting parameter. Note that the auxiliary branch requires a decoder to compute \mathbf{h}_{u}^{pred} and then $\mathbf{z}_{t,u}^{aux,l}$. Instead of adding another decoder specifically for auxiliary branch, we propose to share the primary decoder during the forward pass; however, we do not update the decoder parameters if the gradients are back propagated from the auxiliary RNN-T loss. Because the auxiliary loss is to address the encoder underfitting issue, decoder is not explicitly learned to fit the auxiliary objective function.

Note that for the auxiliary model, we connect a nonlinear MLP (Eq. 7) – rather than a single linear layer – to the intermediate encoder layer l. Since the lower encoder layers are focused on feature extraction rather than the meaningful final label prediction, directly encouraging discrimination in the low-level representations is suboptimal. This is similar to the primary branch, where additional encoder layers of l + 1 to L are added on top of layer l; thus, adding the MLP allows for a similar coarse-to-fine architecture, and the shared encoder layers play a more consistent role for both branches.

Finally, when an encoder has a large network depth, we can apply such criterion to multiple encoder layers. Denote Φ as a set of encoder layer indices that are connected with each auxiliary criterion, and $\Phi \subseteq \{1 \dots L-1\}$. Denote *I* as a binary indicator function as

$$I(l) = \begin{cases} 1, & l \in \Phi\\ 0, & l \notin \Phi \end{cases}$$
(10)

where $1 \leq l < L$. Then Eq. 9 becomes

$$\mathcal{L}(\theta, \phi^{\Phi}) = \mathcal{L}^{\text{RNN-T}}(\theta) + \lambda_{\text{aux}} \mathcal{L}^{\text{RNN-T}}(\theta^{\text{enc}}_{\text{shared}}, \phi^{\Phi})$$

= $-\log P(\mathbf{y}|\mathbf{x}) - \lambda_{\text{aux}} \sum_{l=1}^{L-1} I(l) \log P_{\text{aux},l}(\mathbf{y}|\mathbf{x})$ (11)

2.2.2. Auxiliary symmetric KL divergence criterion

Further, prior works [26, 27, 28] show that aligning the pairwise posterior distributions of multiple (sub)networks in a mutual learning strategy achieves better performance than learning independently. Thus, other than the supervised learning objective function, i.e. RNN-T criterion, we also explore an additional symmetric KL divergence criterion between the output posterior distributions of both branches:

$$\mathcal{L}^{\mathrm{KL}}(\theta, \phi^{\Phi}) = \frac{1}{T} \sum_{t=1}^{T} \sum_{l=1}^{L-1} \frac{1}{U} \sum_{u=1}^{U} I(l) \begin{bmatrix} \\ D_{\mathrm{KL}}(P(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) \| P_{\mathrm{aux},l}(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) \\ + D_{\mathrm{KL}}(P_{\mathrm{aux},l}(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) \| P(y_u | \mathbf{x}_{1:t}, y_{1:(u-1)}) \end{bmatrix}$$
(12)

where the posteriors are given by Eq. 5 and 8. Such KL divergence criterion can also guide the auxiliary branch with the supervision signals from primary branch, as a knowledge distillation procedure. As analyzed in [28], the gradients of multiple loss functions can be counteractive, and such KL loss penalizes the inconsistent gradients with respect to their shared parameters. Thus, the training objective can be:

$$\mathcal{L}(\theta, \phi^{\Phi}) = \mathcal{L}^{\text{RNN-T}}(\theta) + \lambda_{\text{aux}} \mathcal{L}^{\text{KL}}(\theta, \phi^{\Phi})$$
(13)

However, the direct application of RNN-T criterion to the auxiliary model can still be useful, since the auxiliary branch thus always contributes meaningful gradients before the primary model outputs are informative. Therefore, the overall training objective becomes:

$$\mathcal{L}(\theta, \phi^{\Phi}) = \mathcal{L}^{\text{RNN-T}}(\theta) + \lambda_{\text{aux}}(\mathcal{L}^{\text{RNN-T}}(\theta^{\text{enc}}_{\text{shared}}, \phi^{\Phi}) + \mathcal{L}^{\text{KL}}(\theta, \phi^{\Phi}))$$
(14)



Fig. 2. Illustration of the proposed auxiliary context-dependent graphemic state prediction task.

Finally, after training, we discard the auxiliary branch and there is no additional computation overhead for ASR decoding.

2.3. Auxiliary context-dependent graphemic state prediction

In an HMM-based phonetic hybrid ASR system, the triphone HMM states are tied via traditional decision tree clustering [1]. Such a set of tied triphone HMM states – also referred to as context-dependent phonetic states or senones [24] – are used as the output units for the neural network based acoustic model. To further remove the need of a pronunciation lexicon, context-dependent graphemic hybrid models have been developed, and the tri-grapheme HMM states are tied instead. Accordingly, the neural network output units become tied tri-grapheme states, i.e. chenones [29], and the training criterion is cross entropy (CE) loss in conventional hybrid CE models.

While RNN-T uses context-independent graphemes or wordpieces as output units, adding the chenone prediction supervision to encoder layers can transfer complementary tri-grapheme knowledge, encouraging diverse and discriminative encoder representations. Then we can apply such CE criterion to multiple encoder layers. Similarly, given an *L*-layer encoder, denote Φ as a set of encoder layer indices that are connected to chenone prediction, and $\Phi \subseteq \{1 \dots L\}$. Denote *I* as a binary indicator function, and $1 \le l \le L$. As in Figure 2, if I(l) = 1, $\mathbf{h}^{\text{enc},l}$ goes through a one-hidden-layer multi-layer perceptron (MLP)², parameterized by ϕ^l , and then a softmax function to provide a posterior distribution over chenone label set S:

$$P(s_t | \mathbf{h}_t^{\text{enc}, l}) = \text{Softmax}(\text{MLP}(\mathbf{h}_t^{\text{enc}, l}))$$
(15)

where $s_t \in S$, and the auxiliary CE loss is

$$\mathcal{L}^{\text{CE}}(\theta_{\text{shared}}^{\text{enc}}, \phi^{\Phi}) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{l=1}^{L} I(l) \log P(s_t | \mathbf{h}_t^{\text{enc}, l})$$
(16)

The overall training objective is:

$$\mathcal{L}(\theta, \phi^{\Phi}) = \mathcal{L}^{\text{RNN-T}}(\theta) + \lambda_{\text{ce}} \mathcal{L}^{\text{CE}}(\theta_{\text{shared}}^{\text{enc}}, \phi^{\Phi})$$
(17)

where λ_{ce} is a tunable weighting parameter.

 $^{^{2}}$ Note that we use a linear layer rather than a MLP for the topmost/*L*th layer, since the top encoder layer has been designed for final label prediction.

Table 1. The amounts of audio data in hours.

Train	Valid	Test	
		clean	noisy
161	5.2	5.1	10.2
3.1K	13.6	21.2	23.4
3.2K	13.8	24.5	24.0
	Train 161 3.1K 3.2K	Train Valid 161 5.2 3.1K 13.6 3.2K 13.8	Train Valid Technologie 161 5.2 5.1 3.1K 13.6 21.2 3.2K 13.8 24.5

3. EXPERIMENTS

3.1. Experimental setup

3.1.1. Data

We first evaluate our proposed approaches on our in-house Romanian, Turkish and German video datasets, which are sampled from public social media videos and de-identified before transcription. These videos contain a diverse range of acoustic conditions, speakers, accents and topics. The test sets for each language are composed of *clean* and *noisy* categories, with *noisy* category being more acoustically challenging than *clean*. The dataset sizes are shown in Table 1. Moreover, we also perform evaluations on the public LibriSpeech dataset [38].

Input acoustic features are 80-dimensional log-mel filterbank coefficients with 25 ms window size, and we apply mean and variance normalization. We apply the frequency and time masking as in the policy LD from SpecAugment [39], with p = 0.2 on video datasets and p = 1.0 on LibriSpeech. We perform speed perturbation [40] of the training data, and produce three versions of each audio with speed factors 0.9, 1.0 and 1.1. The training data size is thus tripled. For the low-resource Romanian, we further apply another 2-fold data augmentation based on additive noise as in [41], and the training data size is thus 6 times the size of original train set.

3.1.2. System implementation details

For each video language, RNN-T output labels consist of a blank label and 255 wordpieces generated by the unigram language model algorithm from SentencePiece toolkit [42]. To provide chenone labels (Section 2.3), forced alignments are generated via a graphemic hybrid model [29] for each language, and the number of unique chenone labels range from 7104 to 9272.

For video datasets, we build each RNN-T encoder based on latency-controlled bidirectional long short-term memory (LC-BLSTM) network [43]. Each encoder is a 5-layer LC-BLSTM network with 800 hidden units in each layer and direction, and dropout 0.3. Two subsampling layers with stride 2 are applied after first and second LC-BLSTM layer. The prediction network is a 2-layer LSTM of 160 hidden units for Romanian, and 512 units for Turkish and German, with dropout 0.3. Each joint network has 1024 hidden units, and a softmax layer of 256 units for blank and wordpieces. For all neural network implementation, we use an in-house extension of PyTorch-based *fairseq* [44] toolkit. All experiments use multi-GPU and mixed precision training supported in *fairseq*, Adam optimizer [45], and tri-stage [39] learning rate schedule with peak learning rate $4e^{-4}$.

For LibriSpeech, we experiment with two VGG transformer encoders of 24 and 36 layers as in [32], except that we use three VGG blocks with stride 2 in the first two blocks and 1 in the third block. Each transformer layer has an embedding dimension 512 and attention heads 8; feed-forward network (FFN) size is 2048 for 24-layer transformer, and 3072 for the 36-layer. Wordpiece size is 1000 for

Table 2. WER results on Romanian dataset. λ_{aux} is used in Eq. 9, 13 and 14. "aux" and "kl" loss denote the auxiliary RNN-T (Section 2.2.1) and KL divergence criterion (Section 2.2.2) respectively. "crosslingual pretrain" denotes the encoder pretrained from a high-resource Spanish RNN-T. WERR (%) is the unweighted average of the respective relative WER reductions on clean and noisy test sets.

Model	λ_{aux}	valid	clean	noisy	WERR
baseline	-	24.0	20.5	22.0	-
+ aux loss	0.1	23.2			
	0.3	22.8	19.6	21.0	4.5%
	0.6	23.1			
+ kl loss	0.3	22.9			
	0.6	22.6	19.3	20.6	6.1%
	0.9	22.7			
+ aux + kl loss	0.3	22.5	19.1	20.8	6.1%
+ crosslingual pretrain		19.4	15.9	17.6	21.2%
+ aux + kl loss	0.3	18.9	15.7	17.2	22.6%

Table 3. WER results on Romanian. λ_{ce} is used in Eq. 17. "ce pretrain" denotes encoder pretraining from graphemic hybrid CE model. "ce loss" denotes auxiliary chenone prediction objective function (Section 2.3). "mid" denotes connecting CE loss to the 3rd (middle) encoder layer, and "top" denotes connecting CE loss to the 5th (topmost) encoder layer.

Model	$\lambda_{ m ce}$	valid	clean	noisy	WERR
baseline	-	24.0	20.5	22.0	_
+ ce pretrain	-	22.8	19.3	20.9	5.4%
	0.3	23.2			
+ ce loss, top	0.6	22.9	19.8	21.2	3.5%
	0.9	23.1			
	0.3	22.3			
+ ce loss, mid	0.6	22.0	18.5	20.3	8.7%
	0.9	22.0			
+ ce pretrain, ce loss, mid	0.6	21.4	17.9	19.6	11.8%
+ ce pretrain, ce loss, mid, top	0.6	21.2	17.8	19.5	12.3%

the 24-layer, and 2048 for the 36-layer, resulting in total model parameters of 83.3M and 160.3M respectively.

3.2. Auxiliary RNN-T modeling results on video datasets

We first perform experimental evaluations on the low-resource language Romanian, and obtain the optimal λ_{aux} in Eq. 9, 13 and 14. ASR word error rate (WER) results are shown in Table 2. For both *clean* and *noisy* test sets, we first compute the relative WER reduction (WERR) over respective baseline as a percentage, and then take the unweighted average of two percentages, which we refer to as an average WERR.

As shown in Table 2, for auxiliary RNN-T loss (Eq. 9), we vary λ_{aux} over {0.1, 0.3, 0.6}, and observe 0.3 gives the lowest WER on valid set. So we proceed with $\lambda_{aux} = 0.3$ to decode the clean and noisy test sets, and see an average WERR 4.5%. Similarly for the symmetric KL divergence loss (Eq. 13), we vary λ_{aux} over {0.3, 0.6, 0.9}; we find $\lambda_{aux} = 0.6$ works best and provides an average WERR 6.1%. When combining the two objectives with $\lambda_{aux} = 0.3$ (Eq. 14), we find it also gives an average WERR 6.1%, which is better than using auxiliary RNN-T loss on its own.

For the low-resource scenario, one approach to address the lack

	Turkish			German		
Model	clean	noisy	WERR	clean	noisy	WERR
baseline	17.1	18.9	_	11.6	13.0	_
+ aux loss	16.8	18.8	1.1%	11.3	12.6	2.8%
+ kl loss	16.7	18.8	1.4%	11.5	12.8	1.2%
+ aux + kl loss	16.4	18.5	3.1%	11.3	12.6	2.8%
+ crosslingual pretrain	16.6	18.6	2.3%	11.4	12.8	1.6%
+ aux + kl loss	16.1	18.1	5.0%	11.3	12.4	3.6%
+ ce pretrain	16.8	18.9	0.9%	11.5	12.8	1.2%
+ ce loss, mid	16.5	18.4	3.1%	11.3	12.5	3.2%
+ ce loss, mid, top	16.3	18.2	4.2%	11.2	12.3	4.4%

Table 4. WER results on Turkish and German, with $\lambda_{aux} = 0.3$ and $\lambda_{ce} = 0.6$.

of resources are to make use of data from high-resource languages. We thus perform crosslingual pretraining experiments with a Spanish RNN-T model trained on 7K hours. We use the Spanish encoder as the pretrained encoder for Romanian, and proceed with RNN-T training as before, which provides substantial improvements as in Table 2. While on top of crosslingual pretraining, adding auxiliary RNN-T and KL divergence loss provides moderate gain.

We use the optimal λ_{aux} found in each condition and evaluate the performance on Turkish and German. As shown in Table 4, the proposed combination of auxiliary RNN-T and KL divergence loss provides consistent improvements, which is also better than using each individually. We use the same Spanish encoder for crosslingual pretraining, and the improvements are much less due to the increased training data size. Along with the proposed auxiliary RNN-T modeling, they combine to produce noticeable gains.

3.3. Auxiliary chenone prediction results on video datasets

Since we build graphemic hybrid systems to provide chenone labels, we can additionally use the hybrid model as pretrained encoder for RNN-T. As shown in [13, 14], pretraining RNN-T encoder with connectionist temporal classification (CTC) or hybrid CE criterion can improve performance, and we also find CE pretraining produces an average 5.4% WERR on the low-resource Romanian as in Table 3.

For the medium-resource Turkish and German (i.e. training data size of ~3K hours), we initially find pretraining with hybrid CE model can provide 2 - 4% improvements with a relatively small training mini-batch size. However, after optimizing the memory cost by mixed precision training and function merging [7], RNN-T training can enable larger mini-batch size, and we only observe minor improvements 0.9 - 1.2% in Table 4.

Then we experiment with λ_{ce} (Eq. 17) on Romanian. Given each 5-layer LC-BLSTM encoder, we also examine connecting chenone prediction to 3rd (middle) layer or 5th (topmost) layer. As in Table 3, $\lambda_{ce} = 0.6$ works best in each case. While attaching chenone prediction to middle layer performs better than top layer, they combine to provide further improvements on top of CE pretraining.

We continue to evaluate the Turkish and German performance with $\lambda_{ce} = 0.6$. As in Table 4, training on both middle and top layers for auxiliary chenone prediction outperforms training on each alone, and produces noticeable improvements of 4.2 - 4.4% when combined with CE pretraining.

 Table 5.
 WER results on LibriSpeech, with 24-layer transformer

 encoder and 83M total model parameters.

Model	test-clean	WERR	test-other	WERR
baseline	2.77	_	6.60	-
+ aux + kl loss	2.48	10.6%	5.62	14.8%
+ ce loss	2.42	12.6%	5.75	12.9%
+ aux + kl + ce loss	2.31	16.5%	5.26	20.3%

 Table 6. Comparison of our models (with 36-layer transformer encoder and 160M total model parameters) with recently published best results on LibriSpeech.

Model	w/o	LM	w/ LM		
	test-clean	test-other	test-clean	test-other	
LAS					
LSTM [46]	2.6	6.0	2.2	5.2	
Hybrid					
Transformer [32]	2.6	5.6	2.3	4.9	
CTC					
Transformer [47]	2.3	4.8	2.1	4.2	
Sequence Transducer					
Transformer [33]	2.4	5.6	2.0	4.6	
Conformer [48]	2.1	4.3	1.9	3.9	
Transformer (Ours)	2.2	4.7	2.0	4.2	

3.4. Results on LibriSpeech with transformer encoders

While we use streamable 5-layer LC-BLSTM encoders on video datasets above, we experiment with 24/36-layer transformer encoders instead on LibriSpeech. Given the much larger encoder depth, when evaluating the auxiliary RNN-T and KL divergence, we find it more effective to apply the loss at multiple layers. Thus for the 24-layer transformer, we apply it to the 6th, 12th and 18th encoder layers. As in Table 5, it provides about 11% and 15% WERR on each test set. When evaluating the auxiliary CE loss, we apply it at the middle (12th) and top (24th) layer again, which also produces substantial relative gains about 13%.

Additionally, we also attempt to apply both auxiliary tasks simultaneously, i.e., auxiliary RNN-T and KL divergence loss at 6th and 18th layers, and CE loss at 12th and 24th layers. In all cases, we use $\lambda_{aux} = 0.3$ and $\lambda_{ce} = 0.6$ found above (Section 3.2 and 3.3). As in Table 5, both auxiliary tasks combine to produce significant and complementary improvements. These performance gains are much larger than those on the video datasets with 5-layer LC-BLSTM encoder. We conjecture that transformer networks of increased depth suffer more from the encoder undertraining and gradient vanishing problem at lower layers (as discussed in Section 2.2), and auxiliary tasks play more effective roles in addressing it.

We proceed to increase transformer encoder depth from 24 to 36 layers, FFN size from 2048 to 3072, and wordpiece size from 1000 to 2048. We observe that without the auxiliary tasks, neither 24-layer transformer of FFN size 3072 nor 36-layer transformer of FFN 2048 is able to converge. Instead both can converge while using either of the two auxiliary tasks. Finally, the 36-layer transformer of FFN 3072 - which uses auxiliary RNN-T and KL divergence loss at 9th and 27th layers, and CE loss at 18th and 36th layers - produces our best results in Table 6. Auxiliary tasks thus provide an opening for learning deep encoder network, and the increased depth is central to accuracy gains.

We further perform first-pass shallow fusion [49] with an external language model (LM). We use a 4-layer LSTM LM with 4096 hidden units, and LM training data consists of LibriSpeech transcripts and text-only corpus (800M word tokens), tokenized with the 2048 wordpiece model. As in Table 6, we thus achieve competitive results compared to the prior top-performing models.

4. RELATED WORK

Attaching auxiliary objective functions to intermediate layers has been explored in various prior works. For improving image recognition, multiple auxiliary classifiers with squared hinge losses were used in [50], and CE objective functions used in [37, 51], while later [51] only reported limited performance gain. [28] made further progress by showing that, the gradients of multiple loss functions with respect to their shared parameters can counteract each other, and minimizing the symmetric KL divergence between the multiple classifier outputs can penalize such inconsistent gradients and provide more performance gains.

Similarly, for improving hybrid ASR models trained with CE criterion, [27] connected an intermediate layer directly with a linear projection layer to compute the logits over senones, and used an asymmetric KL divergence loss between the primary model output (i.e. senone posterior) and the auxiliary classifier output. While in our work, we found connecting a nonlinear MLP - rather than a single linear layer - to the intermediate layer is more effective, which disentangled low-level feature extraction from final wordpiece prediction. Also for improving ASR, [52] applied CTC or CE objective functions to multiple encoder layers, although without the cross-layer KL divergence loss. While CTC or hybrid senone/chenone models can directly produce posteriors over output labels, RNN-T requires a decoder to compute the output (wordpiece) posterior. Thus, in applying the auxiliary RNN-T or KL divergence loss, we specifically share the RNN-T decoder during the forward pass while keeping it intact from the backward pass (as discussed in Section 2.2.1).

Note that compared to CTC or hybrid models, attention-based seq2seq model is more similar to RNN-T, since both have a neural encoder and decoder. And for attention-based seq2seq model, using auxiliary senone labels has shown improved WERs in [22, 23], while recent work [53] showed contrary observations.

5. CONCLUSIONS

In this work, we propose the use of auxiliary tasks in improving RNN-T based ASR. We first benchmark the streamable LC-BLSTM encoder based performance on video datasets. Applying either auxiliary RNN-T or symmetric KL divergence objective function to intermediate encoder layers has been shown to improve ASR performance, and combining both is more effective than each on its own. Performing auxiliary chenone prediction also provides noticeable complementary gains on top of hybrid CE pretraining.

Next, we demonstrate the efficacy of both auxiliary tasks in improving the transformer encoder based sequence transducer results on LibriSpeech. Both auxiliary tasks provide substantial and complementary gains, and we find that, critical to the convergence of learning deep transformer encoders is the application of auxiliary objective functions to multiple encoder layers. Lastly, to participate in the LibriSpeech benchmark challenge, we develop a 36-layer transformer encoder via both auxiliary tasks, which achieves a WER of 2.0% on test-clean, 4.2% on test-other.

6. REFERENCES

- Steve J Young, Julian J Odell, and Philip C Woodland, "Treebased state tying for high accuracy acoustic modelling," in *Proceedings of the workshop on Human Language Technology*. Association for Computational Linguistics, 1994, pp. 307–312.
- [2] Alex Graves, "Sequence transduction with recurrent neural networks," *arXiv preprint arXiv:1211.3711*, 2012.
- [3] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," in *Proc. ICASSP*, 2016.
- [4] Geoffrey Zweig, Chengzhu Yu, Jasha Droppo, and Andreas Stolcke, "Advances in all-neural speech recognition," in *Proc. ICASSP*, 2017.
- [5] Hasim Sak, Matt Shannon, Kanishka Rao, and Françoise Beaufays, "Recurrent neural aligner: An encoder-decoder neural network model for sequence to sequence mapping," in *Proc. Interspeech*, 2017.
- [6] Yanzhang He, Tara N Sainath, Rohit Prabhavalkar, Ian McGraw, Raziel Alvarez, Ding Zhao, David Rybach, Anjuli Kannan, Yonghui Wu, Ruoming Pang, et al., "Streaming end-to-end speech recognition for mobile devices," in *Proc. ICASSP*, 2019.
- [7] Jinyu Li, Rui Zhao, Hu Hu, and Yifan Gong, "Improving RNN transducer modeling for end-to-end speech recognition," in *Proc. ASRU*, 2019.
- [8] Tara N Sainath, Yanzhang He, Bo Li, Arun Narayanan, Ruoming Pang, Antoine Bruguier, Shuo-yiin Chang, Wei Li, Raziel Alvarez, Zhifeng Chen, et al., "A streaming on-device end-toend model surpassing server-side conventional model quality and latency," in *Proc. ICASSP*, 2020.
- [9] Eric Battenberg, Jitong Chen, Rewon Child, Adam Coates, Yashesh Gaur Yi Li, Hairong Liu, Sanjeev Satheesh, Anuroop Sriram, and Zhenyao Zhu, "Exploring neural transducers for end-to-end speech recognition," in *Proc. ASRU*, 2017.
- [10] Chung-Cheng Chiu, Wei Han, Yu Zhang, Ruoming Pang, Sergey Kishchenko, Patrick Nguyen, Arun Narayanan, Hank Liao, Shuyuan Zhang, Anjuli Kannan, et al., "A comparison of end-to-end models for long-form speech recognition," in *Proc. ASRU*, 2019.
- [11] Jinyu Li, Yu Wu, Yashesh Gaur, Chengyi Wang, Rui Zhao, and Shujie Liu, "On the comparison of popular end-to-end models for large scale speech recognition," *arXiv preprint arXiv:2005.14327*, 2020.
- [12] Xiaohui Zhang, Frank Zhang, Chunxi Liu, Kjell Schubert, Julian Chan, Pradyot Prakash, Jun Liu, Ching-feng Yeh, Fuchun Peng, Yatharth Saraf, and Geoffrey Zweig, "Benchmarking LF-MMI, CTC and RNN-T criteria for streaming ASR," in *Proc. SLT*, 2021.
- [13] Kanishka Rao, Haşim Sak, and Rohit Prabhavalkar, "Exploring architectures, data and units for streaming end-to-end speech recognition with rnn-transducer," in *Proc. ASRU*, 2017.
- [14] Hu Hu, Rui Zhao, Jinyu Li, Liang Lu, and Yifan Gong, "Exploring pre-training with alignments for rnn transducer based end-to-end speech recognition," in *Proc. ICASSP*, 2020.
- [15] Arun Narayanan, Rohit Prabhavalkar, Chung-Cheng Chiu, David Rybach, Tara N Sainath, and Trevor Strohman, "Recognizing long-form speech using streaming end-to-end models," in *Proc. ASRU*, 2019.

- [16] Chao Weng, Chengzhu Yu, Jia Cui, Chunlei Zhang, and Dong Yu, "Minimum Bayes risk training of RNN-transducer for endto-end speech recognition," *arXiv preprint arXiv:1911.12487*, 2019.
- [17] Ashutosh Pandey, Chunxi Liu, Yun Wang, and Yatharth Saraf, "Dual application of speech enhancement for automatic speech recognition," in *Proc. SLT*, 2021.
- [18] Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David Silver, and Koray Kavukcuoglu, "Reinforcement learning with unsupervised auxiliary tasks," *Proc. ICLR*, 2017.
- [19] Shikun Liu, Andrew Davison, and Edward Johns, "Selfsupervised generalisation with meta auxiliary learning," in *Proc. Advances in Neural Information Processing Systems*, 2019.
- [20] Rich Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [21] Yunshu Du, Wojciech M Czarnecki, Siddhant M Jayakumar, Razvan Pascanu, and Balaji Lakshminarayanan, "Adapting auxiliary losses using gradient similarity," *arXiv preprint arXiv:1812.02224*, 2018.
- [22] Shubham Toshniwal, Hao Tang, Liang Lu, and Karen Livescu, "Multitask learning with low-level auxiliary tasks for encoderdecoder based speech recognition," *Proc. Interspeech*, 2017.
- [23] Takafumi Moriya, Sei Ueno, Yusuke Shinohara, Marc Delcroix, Yoshikazu Yamaguchi, and Yushi Aono, "Multi-task learning with augmentation strategy for acoustic-to-word attention-based encoder-decoder speech recognition.," in *Proc. Interspeech*, 2018.
- [24] George E Dahl, Dong Yu, Li Deng, and Alex Acero, "Context-dependent pre-trained deep neural networks for largevocabulary speech recognition," *IEEE Transactions on audio, speech, and language processing*, vol. 20, no. 1, pp. 30–42, 2012.
- [25] Poorya Zaremoodi and Gholamreza Haffari, "Adaptively scheduled multitask learning: The case of low-resource neural machine translation," in *Proceedings of the 3rd Workshop on Neural Generation and Translation*, 2019, pp. 177–186.
- [26] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu, "Deep mutual learning," in *Proc. CVPR*, 2018.
- [27] Liang Lu, Eric Sun, and Yifan Gong, "Self-teaching networks," Proc. Interspeech, 2019.
- [28] Duo Li and Qifeng Chen, "Dynamic hierarchical mimicking towards consistent optimization objectives," in *Proc. CVPR*, 2020.
- [29] Duc Le, Xiaohui Zhang, Weiyi Zheng, Christian Fügen, et al., "From senones to chenones: tied context-dependent graphemes for hybrid speech recognition," *Proc. ASRU*, 2019.
- [30] Stephan Kanthak and Hermann Ney, "Context-dependent acoustic modeling using graphemes for large vocabulary speech recognition," in *Proc. ICASSP*, 2002.
- [31] Mark JF Gales, Kate M Knill, and Anton Ragni, "Unicodebased graphemic systems for limited resource languages," in *Proc. ICASSP*, 2015.
- [32] Yongqiang Wang, Abdelrahman Mohamed, Duc Le, Chunxi Liu, Alex Xiao, Jay Mahadeokar, Hongzhao Huang, Andros Tjandra, Xiaohui Zhang, Frank Zhang, et al., "Transformerbased acoustic modeling for hybrid speech recognition," in *Proc. ICASSP*, 2020.

- [33] Qian Zhang, Han Lu, Hasim Sak, Anshuman Tripathi, Erik McDermott, Stephen Koo, and Shankar Kumar, "Transformer transducer: A streamable speech recognition model with transformer encoders and RNN-T loss," in *Proc. ICASSP*, 2020.
- [34] Sepp Hochreiter, Yoshua Bengio, Paolo Frasconi, Jürgen Schmidhuber, et al., "Gradient flow in recurrent nets: the difficulty of learning long-term dependencies," 2001.
- [35] Rupesh K Srivastava, Klaus Greff, and Jürgen Schmidhuber, "Training very deep networks," in *Proc. NIPS*, 2015.
- [36] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proc. CVPR*, 2016.
- [37] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich, "Going deeper with convolutions," in *Proc. CVPR*, 2015.
- [38] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "LibriSpeech: an ASR corpus based on public domain audio books," in *Proc. ICASSP*, 2015.
- [39] Daniel S Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D Cubuk, and Quoc V Le, "SpecAugment: A simple data augmentation method for automatic speech recognition," *Proc. Interspeech*, 2019.
- [40] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and Sanjeev Khudanpur, "Audio augmentation for speech recognition," in *Proc. Interspeech*, 2015.
- [41] Chunxi Liu, Qiaochu Zhang, Xiaohui Zhang, Kritika Singh, Yatharth Saraf, and Geoffrey Zweig, "Multilingual graphemic hybrid ASR with massive data augmentation," in *Proceedings* of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL), 2020.
- [42] Taku Kudo and John Richardson, "Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing," *arXiv preprint arXiv:1808.06226*, 2018.
- [43] Yu Zhang, Guoguo Chen, Dong Yu, Kaisheng Yaco, Sanjeev Khudanpur, and James Glass, "Highway long short-term memory RNNs for distant speech recognition," in *Proc. ICASSP*, 2016.
- [44] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli, "fairseq: A fast, extensible toolkit for sequence modeling," in *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.
- [45] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.
- [46] Daniel S Park, Yu Zhang, Chung-Cheng Chiu, Youzheng Chen, Bo Li, William Chan, Quoc V Le, and Yonghui Wu, "SpecAugment on large scale datasets," in *Proc. ICASSP*, 2020.
- [47] Frank Zhang, Yongqiang Wang, Xiaohui Zhang, Chunxi Liu, Yatharth Saraf, and Geoffrey Zweig, "Faster, Simpler and More Accurate Hybrid ASR Systems Using Wordpieces," *Proc. Interspeech*, 2020.
- [48] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, et al., "Conformer: Convolutionaugmented transformer for speech recognition," arXiv preprint arXiv:2005.08100, 2020.

- [49] Anjuli Kannan, Yonghui Wu, Patrick Nguyen, Tara N Sainath, Zhijeng Chen, and Rohit Prabhavalkar, "An analysis of incorporating an external language model into a sequence-to-sequence model," in *Proc. ICASSP*, 2018.
- [50] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu, "Deeply-supervised nets," in *Artificial intelligence and statistics*, 2015, pp. 562–570.
- [51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna, "Rethinking the inception architecture for computer vision," in *Proc. CVPR*, 2016.
- [52] Andros Tjandra, Chunxi Liu, Frank Zhang, Xiaohui Zhang, Yongqiang Wang, Gabriel Synnaeve, Satoshi Nakamura, and Geoffrey Zweig, "Deja-vu: Double feature presentation and iterated loss in deep transformer networks," in *Proc. ICASSP*, 2020.
- [53] Hirofumi Inaguma, Yashesh Gaur, Liang Lu, Jinyu Li, and Yifan Gong, "Minimum latency training strategies for streaming sequence-to-sequence ASR," in *Proc. ICASSP*, 2020.