

# LEARNING WORD-LEVEL CONFIDENCE FOR SUBWORD END-TO-END ASR

David Qiu<sup>1</sup>, Qiuqia Li<sup>2\*</sup>, Yanzhang He<sup>1</sup>, Yu Zhang<sup>1</sup>, Bo Li<sup>1</sup>, Liangliang Cao<sup>1</sup>,  
Rohit Prabhavalkar<sup>1</sup>, Deepti Bhatia<sup>1</sup>, Wei Li<sup>1</sup>, Ke Hu<sup>1</sup>, Tara N. Sainath<sup>1</sup>, Ian McGraw<sup>1</sup>

<sup>1</sup> Google, LLC, USA, <sup>2</sup> University of Cambridge, UK

<sup>1</sup>{qdavid, yanzhanghe}@google.com, <sup>2</sup>ql264@cam.ac.uk

## ABSTRACT

We study the problem of word-level confidence estimation in subword-based end-to-end (E2E) models for automatic speech recognition (ASR). Although prior works have proposed training auxiliary confidence models for ASR systems, they do not extend naturally to systems that operate on word-pieces (WP) as their vocabulary. In particular, ground truth WP correctness labels are needed for training confidence models, but the non-unique tokenization from word to WP causes inaccurate labels to be generated. This paper proposes and studies two confidence models of increasing complexity to solve this problem. The final model uses self-attention to directly learn word-level confidence without needing subword tokenization, and exploits full context features from multiple hypotheses to improve confidence accuracy. Experiments on Voice Search and long-tail test sets show standard metrics (e.g., NCE, AUC, RMSE) improving substantially. The proposed confidence module also enables a model selection approach to combine an on-device E2E model with a hybrid model on the server to address the rare word recognition problem for the E2E model.

**Index Terms**— Automatic speech recognition, confidence, calibration, transformer, attention-based end-to-end models

## 1. INTRODUCTION

Confidence scores are an important feature of automatic speech recognition (ASR) systems that supports many downstream applications to mitigate ASR errors [1–3]. For example, unlabelled utterances with high confidence on the ASR output can be included for semi-supervised learning [4–6]. Words with low word-level confidence can be sent for user correction in spoken dialog systems [5]. An utterance that has low confidence can be further processed by a different recognizer for improvement. System combination also commonly relies on confidence as an indication of uncertainty [7, 8].

In conventional HMM-based hybrid systems, confidence scores are estimated for each output word in the hypotheses. An utterance-level confidence is typically aggregated from the word-level confidence when needed. In such systems, word-level confidence scores can be easily estimated from word posterior probabilities computed from lattices or confusion networks [9, 10]. The estimation can be further improved by model-based approaches to combine word posterior probabilities with optional acoustic, linguistic and duration features using a linear regression model [2, 9], or more recently, using conditional random fields [11], recurrent neural networks [12–14] or graph neural networks [15].

Recently, end-to-end (E2E) ASR models such as the recurrent neural network transducer (RNN-T) [16–18], transformer or con-

former transducer [19–21], attention-based encoder-decoder models [22] *inter alia* have gained popularity and achieved state-of-the-art performance in accuracy and latency [17, 18, 23]. In contrast to conventional hybrid systems, they jointly learn acoustic and language modeling in a single neural network that is E2E trained from data. However, deep neural networks tend to exhibit overconfidence in the prediction [24, 25]. Changes to teacher-forcing maximum likelihood training such as label smoothing [26] and scheduled sampling [27] can make the output probabilities less peaky, but the values still do not correlate with word accuracy well. In our recent work [28], we quantified the impact of these methods and proposed a confidence estimation module (CEM) that directly learns the correctness label for each hypothesized subword using a binary cross-entropy loss with features from the encoder and decoder of the E2E model. Although the CEM is simple and effective, it learns subword confidence scores (word-pieces in [28]) by using a fixed subword tokenization for each word in the reference sequence, while the hypothesis may contain other valid tokenizations (see Table 1). This leads to incorrect ground truth labels for training the CEM.

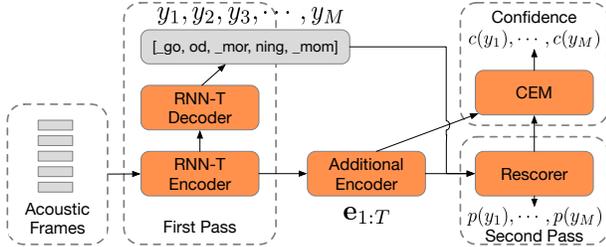
This paper makes the following contributions: 1) propose using self-attention in the CEM to learn the word-level confidence directly for a subword ASR without needing subword tokenization, and 2) leverage cross-attention that attends to both acoustic and linguistic context from multiple hypotheses [29] for additional gains. Experiments show confidence metrics improving substantially from the baseline CEM that learns subword-level confidence [28] to 1) to 2). For application, we test a confidence-based model selection approach: Each utterance is first recognized by an E2E ASR with the proposed word-level CEM on mobile devices for latency and reliability purposes. If the utterance confidence estimated by the CEM is lower than a pre-set threshold, the utterance is sent to the server to be recognized by a conventional hybrid ASR with a large language model (LM) instead for potential quality improvement. Although the on-device E2E model is more accurate overall on a Voice Search test set than the server hybrid model (5.2% vs. 6.4% word error rate (WER)), its lack of LM training data and lexicon causes it to suffer on a rare word test set (17.9% vs. 9.7% WER). The model selection approach achieves the best WER on both sets by applying the same threshold on the CEM output (5.2% on VS and 9.6% on rare word).

## 2. CONFIDENCE FOR TWO-PASS ASR

The base ASR model uses a state-of-the-art two-pass E2E architecture [17] introduced in [23], where the first pass RNN-T generates four candidates for the second pass transformer decoder [30] to rerank. Such architecture is proven to achieve low latency and high accuracy streaming recognition on mobile devices. We aim to add a light-weight CEM while maintaining the efficiency of the model.

For the second pass, the RNN-T is treated as a black box that

\*Work was done while the author interned at Google.



**Fig. 1:** System diagram for the two-pass ASR with confidence.

generates the sequences of acoustic encodings  $\mathbf{e} \triangleq \mathbf{e}_{1:T}$  and the hypothesized subword sequence  $y_{1:M}$ . Fig. 1 shows the overall architecture. The rescorer scores each subword using

$$p(y_i|\mathbf{e}, y_{1:i-1}) = \text{Softmax}(\text{Linear}(\phi(i|\mathbf{e}, y_{1:i-1}))), \quad (1)$$

where  $\phi$  is the rescorer’s penultimate layer activations. The sequence with the highest second pass log probability  $\sum_{i=1}^M \log(p(y_i|\mathbf{e}, y_{1:i-1}))$  is output as the transcription.  $p(y_i|\mathbf{e}, y_{1:i-1})$  serves as a naive estimate of subword confidence. Until the end of Sec. 2.2, the dependence on  $(\mathbf{e}, y_{1:i-1})$  is active even when unwritten.

A dedicated confidence output  $c$  can be computed as

$$\text{top-}K(i) := K \text{ largest log probabilities at decoder index } i \quad (2)$$

$$b(y_i) = [\text{Emb}(y_i); \phi(i|\mathbf{e}, y_{1:i-1}); \log(p(y_i)); \text{top-}K(i)] \quad (3)$$

$$c(y_i) = \sigma(\text{MLP}(b(y_i))). \quad (4)$$

For CEM proposed in [28], a fully-connected multilayer perceptron (MLP) is used. “Emb” is the input subword and position embedding. The CEM can be trained jointly or separately with the ASR, using a binary cross-entropy loss:  $\mathcal{L} = -\sum_{i=1}^M d(y_i) \log c(y_i) + (1 - d(y_i)) \log(1 - c(y_i))$ , where  $d(y_i) = 1$  if the edit distance between hypothesized and reference subword sequences outputs “correct” for  $y_i$  and  $d(y_i) = 0$  if it outputs “insertion” or “substitution”.

It is worth noting that the features above that are related to the posterior probability are also commonly used in confidence models for conventional ASR [9, 11]. We use them as the baseline for the E2E ASR, but also integrate acoustic and linguistic context with self-attention (Sec. 2.2) and deliberation (Sec. 2.3) in a unified network, which dramatically improves the performance over the baseline.

### 2.1. Simple Word Confidence From Word-pieces

To decrease the size of the softmax layer and to improve generalization, the subword vocabulary is typically small compared to the word vocabulary. This can be accomplished with graphemes, word-pieces (WP), etc. In this paper, we focus on WP and use it synonymously with “subword”. To compute the WER, the hypothesized WP sequence  $y_{1:M}$  first needs to be converted to its corresponding word sequence  $w_{1:L}$ . This procedure is uniquely determined since each word’s first WP starts with a word boundary indicator (‘-’), e.g. “\_go”, “od”, “\_mor”, “ning”  $\rightarrow$  “good”, “morning”. Similarly for confidence, for a word  $w_j$  consisting of  $Q_j$  WPs, let  $y_{j,q}$  denote the  $q$ -th WP of the  $j$ -th word; a simple way to compute word confidence is  $c_w(w_j) = \text{agg}(c(y_{j,1}), \dots, c(y_{j,Q_j}))$ , where  $\text{agg}$  can be the arithmetic mean, minimum, product, a neural network, etc. In this paper, we experiment with the arithmetic mean aggregator only.

### 2.2. E2E Word Confidence From Word-pieces

The drawback of the approach in Sec. 2.1 lies in the mismatch between WP correctness and word correctness. Even though WP sequences uniquely determine word sequences, the reverse does not

Hyp:	_go	od	_mor	ning	_mom
Ref:	_go	od	_morn	ing	
WP edit:	cor	cor	sub	sub	ins
Word edit:	-	cor	-	cor	ins
$d(w_j)$ :	-	1	-	1	0
$m(y_i)$ :	0	1	0	1	1
$\mathcal{L}(w_j)$ :	-	$\log c_w(w_1)$	-	$\log c_w(w_2)$	$\log(1 - c_w(w_3))$

**Table 1:** Top: example of a reference having non-unique tokenizations that leads to the inaccurate ground truth WP correctness labels. Bottom: example of using an end-of-word mask  $m$  to implement the word-level loss in a CEM with an output at every WP.

hold. Each reference word can be divided into WPs in multiple valid ways. Table 1 shows an example where the word “morning” is correctly transcribed, but results in two substitutions in the WP edit distance output. Searching over all possible reference tokenizations for the one with the fewest WP edits creates an undesirable computational burden during training, and we do not investigate that solution in this paper. Stochastic methods such as BPE-dropout [31] also do not help here, since they assume that any segmentation is equally valid, which does not hold when computing edit distance.

Using word edit distance output as the ground truth training labels would bypass the multiple tokenization problem. However, because ASR / CEM output at the WP level, two design choices need to be made: at which WP to output the word confidence, and how to incorporate information from every WP that makes up the word. We choose to use the confidence output at the final WP of every word as its word confidence, and change the MLP in (4) to a transformer:

$$\mathbf{b} = \{b(y_1), \dots, b(y_{i-1})\} \quad (5)$$

$$c(y_i) = \sigma(\text{Transformer}(\text{CA}(\mathbf{e}), \text{SA}(\mathbf{b}))) \quad (6)$$

$$c_w(w_j) = c(y_{j,Q_j}), \quad (7)$$

where CA and SA denote the cross-attention and self-attention mechanisms [30], respectively, and  $b(y_i)$  is defined in (3). This allows the model to learn how to attend to the features of earlier WPs in the same word in a true E2E fashion. The cross-attention also improves the confidence estimation using the acoustic context.

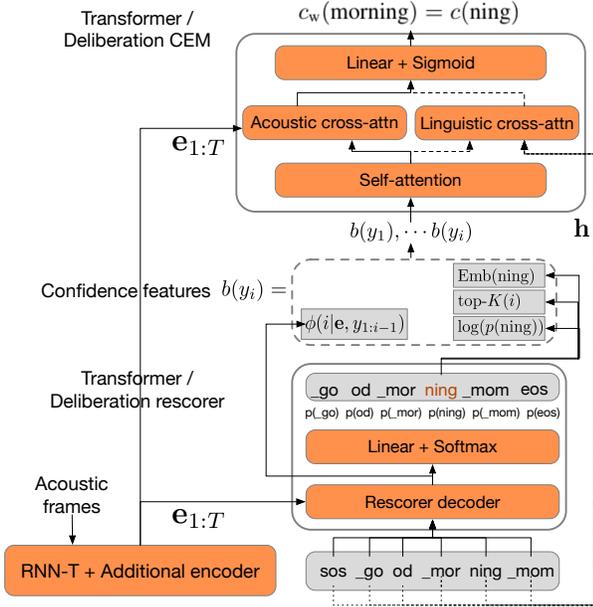
The word-level loss function becomes

$$L = -\sum_{j=1}^L d(w_j) \log c_w(w_j) + (1 - d(w_j)) \log(1 - c_w(w_j))$$

Similar to  $d(y_i)$ , the value of  $d(w_j)$  depends on the word edit distance output. The loss, in a WP-based ASR model, can be easily implemented with an end-of-word masked loss (see Table 1).

### 2.3. Multi-hypotheses Deliberation

Works such as [32] have shown that statistics from multiple beam search decoded hypotheses further improve confidence accuracy. In general, words shared across more hypotheses tend to have higher confidence. Information from multiple hypotheses is even more relevant to the model introduced in Sec. 2.2. In the example from Table 1, having two hypotheses [\_go, od, \_mor, ning, \_mom] and [\_go, od, \_morn, ing, \_mom] attend to each other would inform the model that they concatenate to the same word sequence, and that they should be mapped to similar confidence scores. Additionally, the goal of confidence prediction is to score a known hypothesis. This is different from auto-regressive decoding, where knowing the full hypothesis trivializes the problem. Thus, the CEM can make use of the future context of the hypothesis to score the current word.



**Fig. 2:** Model architecture. For clarity, only the actions of predicting  $c$ (“ning”) are shown. Refer to Table 1 for an example of the masked loss function for the entire sequence. Top- $K$  is defined in (2). All dashed connections and the linguistic cross-attention block are only used in the deliberation CEM but not transformer CEM.

We incorporate two sources of information, acoustic encoding ( $\mathbf{e}$ ) and multiple hypotheses encoding ( $\mathbf{h}$ ), in a learned way through the multi-source attention block in a deliberation model [29]:

$$\mathbf{h} = \left[ \text{BLSTM} \left( y_{1:M_1}^{(1)} \right); \dots; \text{BLSTM} \left( y_{1:M_H}^{(H)} \right) \right] \quad (8)$$

$$c(y_i) = \sigma(\text{Transformer}(\text{CA}(\mathbf{e}) + \text{CA}(\mathbf{h}), \text{SA}(\mathbf{b}))), \quad (9)$$

where  $H$  is the number of hypotheses attended to and  $M_H$  is the number of WPs in the  $H$ -th hypothesis. Fig. 2 shows the model.

### 3. EXPERIMENTAL SETUP

**Architecture and Training Setup:** The RNN-T architecture follows [17] exactly, with 8 LSTM layers in the encoder and 2 LSTM layers in the prediction network. Each LSTM layer is unidirectional, with 2,048 units and a projection layer with 640 units. The transformer rescorer architecture is the same as in [23]. Its encoder consists of the shared encoder from RNN-T and additional 2 LSTM layers, and the decoder consists of 4 self-attention layers, 2 of which contain the cross-attention over the encoder. The deliberation rescorer architecture is described in [29]; for consistency, we replace its LAS decoder with the same 4 layers transformer decoder. We feed up to eight RNN-T beam search results into the linguistic cross-attention mechanism. All internal dimensions in the rescorders are 640. The size of the WP vocabulary is 4,096. The ASR model is frozen during CEM training to not affect the WER.

For the top- $K$  feature in (2), we observe diminishing returns beyond  $K = 4$ , and use this setting for all experiments. Thus, the input features introduced in (3) are 640, 640, 1, 4 dimensional, respectively. “WP MLP” denotes the CEM in (4) with confidence averaged at the word level (Sec. 2.1); we use 3 layers that result in hidden activation with dimensions of 640, 320, 1, respectively. We also replace the second layer with one transformer decoder block and call this setup “WP Xformer”. “E2E Xformer” denotes the CEM in

Sec. 2.2. All confidence models are trained in TensorFlow with the Lingvo [33] toolkit using four hypotheses from the frozen RNN-T, and evaluated on only the top ranked hypothesis. The optimizer is Adam [34] with learning rate 0.0005, and the global batch size is 4,096 across  $8 \times 8$  TPU.

**Training Set:** The models are trained on the multi-domain training set used in [17], which spans domains of search, farfield, telephony and YouTube. All datasets are anonymized and hand-transcribed; the transcription for YouTube utterances is done in a semi-supervised fashion [35]. Multi-condition training (MTR) [36] and random data downsampling to 8kHz [37] are used to further increase data diversity.

**Test Sets:** The main test set includes  $\sim 14\text{K}$  Voice Search (VS) utterances extracted from Google traffic, which is anonymized and hand-transcribed. To test the generalizability of the CEM, we use an in-house named entity tagger to identify a list of proper nouns that are common in the LM training data for conventional ASR but rare in the audio-text paired multi-domain training data for E2E ASR models. We select 10,000 sentences from the LM test data for the maps domain, each of which contains at least one of these rare proper nouns, then synthesize audio for these sentences with a TTS system (as in [38]) to create the Long-tail Maps test set.

**Evaluation Metrics:** We evaluate on standard metrics found in prior works in confidence [13]. To measure per-word confidence accuracy, we use normalized cross-entropy (NCE) [39]. To measure whether the confidence score is highly correlated with word correctness, we use area under the receiver operating characteristic curve (AUC-ROC) and the precision recall curve (AUC-PR). Because the ASR model achieves high word correct ratio ( $\text{WCR} = \frac{\#\text{correctly hypothesized words}}{\#\text{all hypothesized words}}$ ) on all datasets, we compute AUC-PR with the PR curve for the less frequent incorrect class (wrongly hypothesized words). Higher is better for these metrics.

Utterance-level confidence is computed by averaging word-level confidence:  $\frac{1}{L} \sum_{j=1}^L c_w(w_j)$ . To measure its accuracy, we use the root mean squared error (RMSE) between the utterance confidence and either the ground truth WCR or  $(1 - \text{WER})$ . Although WER is the gold standard, because this version of the CEM does not explicitly predict deletions, WCR RMSE is a better indicator of the CEM’s quality. Lower is better for these metrics.

## 4. RESULTS

### 4.1. Simple vs E2E Word Confidence from Word-pieces

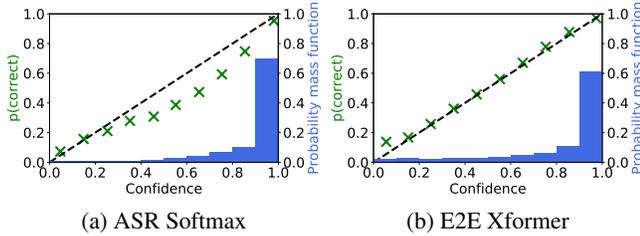
This section compares the performance of the naive confidence from the ASR softmax, WP CEM (Sec. 2.1), and E2E CEM (Sec. 2.2). Table 2 shows the results of the representative models on VS and Maps. Except AUC-PR on Maps, all other metrics improve going from softmax to WP to E2E, showing the effectiveness of our proposed technique. Deep neural networks usually exhibits overconfidence on long-tail data, and Fig. 3 shows the improved calibration curve [24] on Maps arising from E2E word confidence training.

### 4.2. Effects of Input Features

To determine the effects of input embedding in (3), we compare using the ASR input embedding against training a new embedding layer for confidence (“Conf emb”). To quantify the usefulness of the ASR softmax posteriors, we compare using the full set of features in (3), removing the top- $K$  ( $K = 4$ ) feature, and further removing the log posterior  $\log(p(y_i))$ . Table 3 reports the metrics from the input feature ablation study. Having a dedicated confidence embedding layer improves all metrics, at the cost of  $4096 \times 640$  extra parameters. The model degrades when it does not use the log probability

Confidence Models	Voice Search					Long-tail Maps				
	NCE	AUC ROC	AUC PR	WCR RMSE	(1-WER) RMSE	NCE	AUC ROC	AUC PR	WCR RMSE	(1-WER) RMSE
ASR Softmax	0.241	0.873	0.280	0.140	0.244	0.286	0.882	0.635	0.221	0.334
WP MLP [28]	0.269	0.885	0.329	0.138	0.233	0.360	0.887	0.684	0.198	0.297
WP Xformer	0.280	0.885	0.347	0.137	0.231	0.365	0.889	0.690	0.194	0.292
E2E Xformer	0.367	0.928	0.466	0.130	0.221	0.389	0.901	0.682	0.186	<b>0.281</b>
+Delib 1-Hyp	0.361	0.923	0.474	0.128	0.206	0.405	0.908	0.691	<b>0.184</b>	0.299
+Delib 8-Hyp	<b>0.425</b>	<b>0.941</b>	<b>0.508</b>	<b>0.127</b>	<b>0.204</b>	<b>0.416</b>	<b>0.911</b>	<b>0.700</b>	<b>0.184</b>	0.283

**Table 2:** Confidence metrics comparing the WP, E2E, and deliberation CEM. The models in the last 4 rows are proposed in this paper.



**Fig. 3:** Calibration curves for ASR Softmax and E2E Xformer confidence models for word confidence on Long-tail Maps. The black and green curves show the ideal and actual calibration curves, respectively. The blue bar plot shows the probability mass in each bin.

features. Given the modest increase in dimensionality from these five additional features, it is worth including them in the CEM.

	NCE	AUC ROC	AUC PR	WCR RMSE	(1-WER) RMSE
E2E Xformer	0.367	0.928	0.466	0.130	0.221
+Conf emb	0.374	0.930	0.477	0.129	0.219
-top- $K(i)$	0.358	0.924	0.453	0.131	0.222
- $\log(p(y_i))$	0.338	0.924	0.435	0.136	0.223

**Table 3:** Input feature ablation studies on Voice Search.

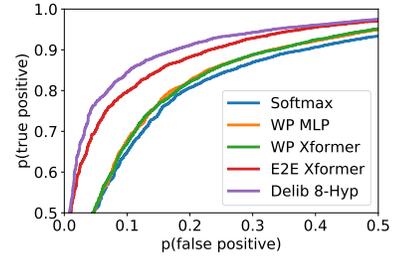
### 4.3. Multi-hypotheses Deliberation Results

This section examines adding multi-hypotheses deliberation to further improve the E2E Xformer CEM. The Delib 1-Hyp model uses a BLSTM to encode the current hypothesis and changes every transformer layer in the ASR and CEM to use multi-source attention that attends to both the acoustic and the BLSTM hypothesis encodings. This allows the CEM to see future context in the hypothesis. The Delib 8-Hyp model uses the RNN-T to generate eight hypotheses to be encoded and attended to. Encodings across different hypotheses are concatenated without any position information. This allows the CEM to use consensus among multiple hypotheses.

Table 2 shows that using full context on the single hypothesis (Delib 1-Hyp) slightly improves some of the metrics. However, consensus from multiple hypotheses (Delib 8-Hyp) greatly improves most metrics over the basic E2E Xformer model. Fig. 4 demonstrates that in a ROC curve, where E2E CEM is clearly better than WP CEM, and adding deliberation improves further.

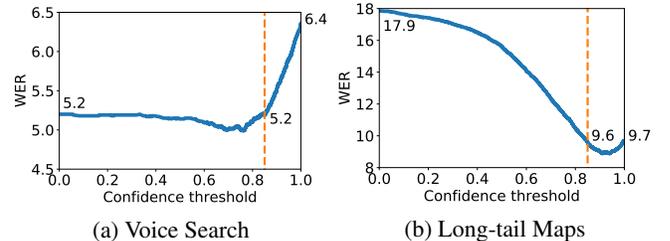
### 4.4. Long-tail Utterances Filtering for Model Selection

For application, we experiment with a confidence-based model selection approach to combine 1) an on-device two-pass E2E ASR with deliberation rescoring and an E2E Delib 8-Hyp confidence



**Fig. 4:** ROC curve on Voice Search for different models.

module (5.2% VS WER), and 2) a conventional hybrid ASR on the server [40] with a large LM (6.4% VS WER). Despite the lower WER on VS, the E2E model’s lack of large-scale LM training data and lexicon causes it to suffer on long-tail utterances (17.9% on Maps compared to 9.7% with server). Thus, the objective is to send all utterances with the confidence score below a pre-set threshold to the server, while keeping the majority of utterances on-device to gain quality, latency, and reliability. Fig. 5 shows the WER of the overall system with different confidence thresholds. When the threshold is set to 0.85, 87% of the VS utterances are processed on-device, and the overall WER is equal to on-device only (5.2%). With the same threshold, only 53% of the Maps utterances are processed on-device, and the overall WER of 9.6% is better than either individual system.



**Fig. 5:** Overall WER at different operating points for model selection. When confidence threshold is 0, all utterances are processed on-device. When it is 1, all utterances are processed on the server.

## 5. CONCLUSION

We propose an extension to a light-weight confidence estimation module for E2E ASR models to directly estimate word-level confidence with self-attention and deliberation, by learning from the full acoustic and linguistic context of subword sequence and multiple hypotheses. Experimental results show the proposed approach is critical to improving confidence metrics substantially when applied to a state-of-the-art two-pass E2E system. It also enables a confidence-based model selection approach to address the rare word recognition problem for the E2E system.

## 6. REFERENCES

- [1] F. Wessel, R. Schluter, K. Macherey, & H. Ney, "Confidence measures for large vocabulary continuous speech recognition," *IEEE Trans. on Speech and Audio Processing*, 2001.
- [2] H. Jiang, "Confidence measures for speech recognition: A survey," *Speech communication*, 2005.
- [3] D. Yu, J. Li, & L. Deng, "Calibration of confidence measures in speech recognition," *IEEE Trans. on Audio, Speech, and Language Processing*, 2011.
- [4] H.Y. Chan & P.C. Woodland, "Improving broadcast news transcription by lightly supervised discriminative training," in *ICASSP*, 2004.
- [5] G. Tur, D. Hakkani-Tür, & R.E. Schapire, "Combining active and semi-supervised learning for spoken language understanding," *Speech Communication*, 2005.
- [6] D. Park, Y. Zhang, Y. Jia, W. Han, C.C. Chiu, *et al.*, "Improved noisy student training for automatic speech recognition," in *Interspeech*, 2020.
- [7] G. Evermann & P.C. Woodland, "Posterior probability decoding, confidence estimation and system combination," in *NIST Speech Transcription Workshop*, 2000.
- [8] J.G. Fiscus, "A post-processing system to yield reduced word error rates: Recognizer output voting error reduction (ROVER)," in *ASRU*, 1997.
- [9] G. Evermann & P.C. Woodland, "Large vocabulary decoding and confidence estimation using word posterior probabilities," in *ICASSP*, 2000.
- [10] L. Mangu, E. Brill, & A. Stolcke, "Finding consensus in speech recognition: Word error minimization and other applications of confusion networks," *Computer Speech and Language*, 2000.
- [11] M.S. Seigel & P.C. Woodland, "Combining information sources for confidence estimation with CRF models," in *Interspeech*, 2011.
- [12] K. Kalgaonkar, C. Liu, Y. Gong, & K. Yao, "Estimating confidence scores on ASR results using recurrent neural networks," in *ICASSP*, 2015.
- [13] A. Ragni, Q. Li, M.J.F. Gales, & Y. Wang, "Confidence estimation and deletion prediction using bidirectional recurrent neural networks," in *SLT*, 2018.
- [14] P. Swarup, R. Maas, S. Garimella, S.H. Mallidi, & B. Hoffmeister, "Improving ASR confidence scores for alexa using acoustic and hypothesis embeddings," *Interspeech*, 2019.
- [15] Q. Li, P. Ness, A. Ragni, & M.J.F. Gales, "Bi-directional lattice recurrent neural networks for confidence estimation," in *ICASSP*, 2019.
- [16] Y. He, T. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, *et al.*, "Streaming end-to-end speech recognition for mobile devices," in *ICASSP*, 2019.
- [17] T. Sainath, Y. He, B. Li, A. Narayanan, R. Pang, *et al.*, "A streaming on-device end-to-end model surpassing server-side conventional model quality and latency," in *ICASSP*, 2020.
- [18] J. Li, R. Zhao, Z. Meng, Y. Liu, W. Wei, *et al.*, "Developing RNN-T models surpassing high-performance hybrid models with customization capability," in *Interspeech*, 2020.
- [19] Q. Zhang, H. Lu, H. Sak, A. Tripathi, E. McDermott, *et al.*, "Transformer transducer: A streamable speech recognition model with transformer encoders and RNN-T loss," in *ICASSP*, 2020.
- [20] C.F. Yeh, J. Mahadeokar, K. Kalgaonkar, Y. Wang, D. Le, *et al.*, "Transformer-transducer: End-to-end speech recognition with self-attention," *arXiv:1910.12977*, 2019.
- [21] A. Gulati, J. Qin, C.C. Chiu, N. Parmar, Y. Zhang, *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," in *Interspeech*, 2020.
- [22] J.K. Chorowski, D. Bahdanau, D. Serdyuk, *et al.*, "Attention-based models for speech recognition," in *NeurIPS*, 2015.
- [23] W. Li, J. Qin, C.C. Chiu, R. Pang, & Y. He, "Parallel rescaling with transformer for streaming on-device speech recognition," in *Interspeech*, 2020.
- [24] C. Guo, G. Pleiss, Y. Sun, & K.Q. Weinberger, "On calibration of modern neural networks," in *ICML*, 2017.
- [25] Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, *et al.*, "Can you trust your model's uncertainty? Evaluating predictive uncertainty under dataset shift," in *NeurIPS*, 2019.
- [26] J. Chorowski & N. Jaitly, "Towards better decoding and language model integration in sequence to sequence models," in *Interspeech*, 2017.
- [27] S. Bengio, O. Vinyals, N. Jaitly, & N. Shazeer, "Scheduled sampling for sequence prediction with recurrent neural networks," in *NIPS*, 2015.
- [28] Q. Li, D. Qiu, Y. Zhang, B. Li, Y. He, *et al.*, "Confidence estimation for attention-based sequence-to-sequence models for speech recognition," *arXiv:2010.11428*, 2020.
- [29] K. Hu, T.N. Sainath, R. Pang, & R. Prabhavalkar, "Deliberation model based two-pass end-to-end speech recognition," in *ICASSP*, 2020.
- [30] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, *et al.*, "Attention is all you need," in *NeurIPS*, 2017.
- [31] I. Provilkov, D. Emelianenko, & E. Voita, "BPE-Dropout: Simple and effective subword regularization," in *ACL*, 2020.
- [32] N. Itoh, T.N. Sainath, D.N. Jiang, J. Zhou, & B. Ramabhadran, "N-best entropy based data selection for acoustic modeling," in *ICASSP*, 2012.
- [33] J. Shen, P. Nguyen, Y. Wu, Z. Chen, M.X. Chen, *et al.*, "Lingvo: A modular and scalable framework for sequence-to-sequence modeling," *arXiv:1902.08295*, 2019.
- [34] D.P. Kingma & J. Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.
- [35] H. Liao, E. McDermott, & A. Senior, "Large scale deep neural network acoustic modeling with semi-supervised training data for youtube video transcription," in *ASRU*, 2013.
- [36] C. Kim, A. Misra, K.K. Chin, T. Hughes, A. Narayanan, *et al.*, "Generation of large-scale simulated utterances in virtual rooms to train deep-neural networks for far-field speech recognition in Google Home," in *Interspeech*, 2017.
- [37] J. Li, D. Yu, J.T. Huang, & Y. Gong, "Improving wide-band speech recognition using mixed-bandwidth training data in CD-DNN-HMM," in *SLT*, 2012.
- [38] X. Gonzalvo, S. Tazari, C.a. Chan, M. Becker, A. Gutkin, *et al.*, "Recent advances in Google real-time HMM-driven unit selection synthesizer," in *Interspeech*, 2016.
- [39] M. Siu, H. Gish, & F. Richardson, "Improved estimation, evaluation and applications of confidence measures for speech recognition," in *Eurospeech*, 1997.
- [40] G. Pundak & T.N. Sainath, "Lower frame rate neural network acoustic models," in *Interspeech*, 2016.