

Unsupervised Uncertainty Measures of Automatic Speech Recognition for Non-intrusive Speech Intelligibility Prediction

Zehai Tu, Ning Ma, Jon Barker

University of Sheffield, Department of Computer Science, Sheffield, UK

{ztu3, n.ma, j.p.barker}@sheffield.ac.uk

Abstract

Non-intrusive intelligibility prediction is important for its application in realistic scenarios, where a clean reference signal is difficult to access. The construction of many non-intrusive predictors require either ground truth intelligibility labels or clean reference signals for supervised learning. In this work, we leverage an unsupervised uncertainty estimation method for predicting speech intelligibility, which does not require intelligibility labels or reference signals to train the predictor. Our experiments demonstrate that the uncertainty from state-of-the-art end-to-end automatic speech recognition (ASR) models is highly correlated with speech intelligibility. The proposed method is evaluated on two databases and the results show that the unsupervised uncertainty measures of ASR models are more correlated with speech intelligibility from listening results than the predictions made by widely used intrusive methods.

Index Terms: Speech intelligibility prediction, non-intrusive method, unsupervised uncertainty estimation

1. Introduction

Speech intelligibility is usually interpreted as how comprehensible speech is. Accurate intelligibility prediction has always been of great interest for its importance in developing speech enhancement related applications, such as hearing aids. In recent years, non-intrusive intelligibility prediction, which does not require clean signals as references, has drawn increasing attention because of its wider applicability compared to intrusive methods, especially in realistic scenarios. One of the promising candidates for non-intrusive intelligibility prediction is ASR models [1–4], given that they can perform similarly to human speech recognition in certain situations [5–7]. Intelligibility can be characterised by the probability of correct word recognition [8]. Meanwhile, the *uncertainty* of ASR models is associated with the probability of models making correct predictions [9–13].

Motivated by this, this study investigates how to estimate the uncertainty of a strong ASR model and correlate it to speech intelligibility. Specifically, we propose to use an unsupervised ASR uncertainty estimation method, which does not require intelligibility labels or clean references to predict sequence-level speech intelligibility. Our experiments are conducted on both a small vocabulary database with simple noisy scenes and a large vocabulary database with more complex noisy scenes. It is shown that the uncertainty of strong ASR models is highly correlated to speech intelligibility, and the prediction performance can outperform widely used intrusive intelligibility predictors. The experimental results also indicate that the uncertainty of ASR models is better than ASR recognition results at intelligibility prediction.

The next section presents the background for unsupervised ASR uncertainty estimation and recent non-intrusive intelligi-

bility prediction methods. The methodology used to formulate unsupervised ASR uncertainty is explained in Section 3. Section 4 describes the experimental data, setups, and results. Section 5 concludes this work.

2. Background

Uncertainty estimation is crucial for ASR application as it can help improve robustness in critical tasks. Most ASR uncertainty estimation methods construct and optimise an estimator on top of the original ASR model with supervision [11–13]. Recently, a word-level ASR uncertainty estimation method is proposed in [10], and a sequence-level uncertainty estimation method for auto-regressive structured prediction tasks is proposed in [9]. The major advantages of sequence-level uncertainty estimation for intelligibility prediction, which is used in this work, are that it does not require: firstly, human listening results because they are usually noisy and expensive; secondly, token-level labels because the alignment could be difficult and intractable, i.e., listeners may respond little when the speech is not intelligible.

Conventional non-intrusive intelligibility predictors, such as SRMR [14] and ModA [15], take advantage of acoustic features related to intelligibility. They heavily rely on prior knowledge on scene acoustics, such as room reverberant characteristics, therefore the application is limited. Another group of methods can be considered as variants of intrusive predictors, like the short-time objective intelligibility (STOI) [16], such as NI-STOI [17], NIC-STOI [18], THMMB-STOI [19]. A clean feature estimation model is usually constructed and used to produce an estimated reference for computing STOI-like scores. Therefore, clean signals are usually required to optimise the estimation model. Meanwhile, transcription or clean signals are sometimes preferred for alignment or voice activity detection. Recently, a number of data-driven methods are proposed, such as [4, 20–22]. These methods train a classification and regression tree or neural networks to predict intelligibility from features of noisy signals, therefore requiring a number of expensive human listening results or scores from intrusive predictors like STOI. Apart from the aforementioned methods, a series of works including FADE [6], and [3, 23] leverage ASR models to predict speech reception thresholds (SRT), i.e., the signal-to-noise ratio (SNR) at which half of words within a group of noisy utterances are correctly recognised, rather than sequence-level intelligibility scores. The most recent work [23] in the series does not require transcripts or reference signals for intelligibility prediction, while it uses identical noises for training and testing, and the recognition results need to be estimated.

3. Method

In this section we describe how two sequence-level ASR uncertainty measures are formulated: *confidence* and *entropy*, with an ensemble method following the derivation in [9]. The ensemble

ble of models can be interpreted from a Bayesian perspective, i.e., regarding model parameters θ as random variables and using a prior $p(\theta)$ to compute the posterior $p(\theta|\mathcal{D})$ with a given dataset \mathcal{D} . As Bayesian inference is usually intractable for models like deep neural networks, it is possible to take advantage of an approximation $q(\theta)$ to $p(\theta|\mathcal{D})$ with a family of models with different parameters [24]. Monte-Carlo Dropout [25] and Deep Ensembles [26] are two major approaches to generate ensembles, and the latter approach is exploited in this work.

3.1. Uncertainty estimation

Given the ASR training dataset containing variable-length sequences of input acoustic features $\{x_1, \dots, x_N\} = \mathbf{x} \in \mathcal{X}$, and the corresponding transcript targets $\{y_1, \dots, y_L\} = \mathbf{y} \in \mathcal{Y}$, an ensemble of M ASR models $\{P(\mathbf{y}|\mathbf{x}; \theta^{(m)})\}$ can be trained to achieve the approximated posterior $q(\theta)$. The sequence-level predictive posterior $P(\mathbf{y}|\mathbf{x}, \theta)$ can be computed as the expectation of the ensemble:

$$\begin{aligned} P(\mathbf{y}|\mathbf{x}, \theta) &= \mathbb{E}_{q(\theta)}[P(\mathbf{y}|\mathbf{x}, \theta)] \\ &\approx \frac{1}{M} \sum_{m=1}^M P(\mathbf{y}|\mathbf{x}, \theta^{(m)}), \end{aligned} \quad (1)$$

where $\theta^{(m)} \sim q(\theta) \approx p(\theta|\mathcal{D})$. The sequence-level entropy $H(\mathbf{y}|\mathbf{x}, \theta)$ can be expressed as:

$$\begin{aligned} H(\mathbf{y}|\mathbf{x}, \theta) &= \mathbb{E}_{P(\mathbf{y}|\mathbf{x}, \theta)}[-\ln P(\mathbf{y}|\mathbf{x}, \theta)] \\ &= - \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y}|\mathbf{x}, \theta) \ln P(\mathbf{y}|\mathbf{x}, \theta). \end{aligned} \quad (2)$$

It is usually not possible to compute the posterior $P(\mathbf{y}|\mathbf{x}, \theta)$ as \mathcal{Y} is an infinite set with variable-length transcript sequences. However, an autoregressive ASR model could factorise the posterior into a product of conditionals:

$$P(\mathbf{y}|\mathbf{x}, \theta) = \prod_{l=1}^L P(y_l|\mathbf{y}_{<l}, \mathbf{x}; \theta), y_l \in \{\omega_1, \dots, \omega_K\}, \quad (3)$$

where ω represents the byte-pair encoding (BPE) token, and K is the size of BPE vocabulary.

Confidence is usually considered as the maximum predicted probability, and the sequence-level confidence \mathcal{C}_S in this work is regarded as a combination of token-level confidence. In order to make fair comparison of sequences with variable lengths, a length normalisation rate is used [27], and \mathcal{C}_S is computed as:

$$\mathcal{C}_S = \exp \left[\frac{1}{L} \ln \sum_{l=1}^L \max \frac{1}{M} \sum_{m=1}^M P(y_l|\mathbf{y}_{<l}, \mathbf{x}; \theta^{(m)}) \right]. \quad (4)$$

Entropy computation is usually challenging as the expectations of \mathbf{y} are practically intractable, i.e., there are K^L possible candidates for a L -length sequence y_L , and a forward-pass inference needs to be conducted for each hypothesis y . Meanwhile, beam-search in ASR inference stage is able to provide high-quality hypotheses and can therefore be considered as an importance-sampling which yields hypotheses from high-probability space. By using B top hypotheses within a beam, the approximated sequence-level entropy \mathcal{H}_S with simple Monte-Carlo estimation can be computed as:

$$\begin{aligned} \mathcal{H}_S &= - \sum_{b=1}^B \frac{\pi_b}{L^{(b)}} \ln P(\mathbf{y}^{(b)}|\mathbf{x}, \theta), \\ \pi_b &= \frac{\exp \frac{1}{T} \ln P(\mathbf{y}^{(b)}|\mathbf{x}, \theta)}{\sum_k^B \exp \frac{1}{T} \ln P(\mathbf{y}^{(k)}|\mathbf{x}, \theta)}, \end{aligned} \quad (5)$$

where a calibration temperature T can be introduced to adjust the distribution of hypotheses, and:

$$\ln P(\mathbf{y}^{(b)}|\mathbf{x}, \theta) = \sum_{l^{(b)}=1}^{L^{(b)}} \ln \frac{1}{M} \sum_{m=1}^M P(y_l^{(b)}|\mathbf{y}_{<l}^{(b)}, \mathbf{x}; \theta^{(m)}). \quad (6)$$

3.2. ASR posterior

The ASR model used in this work is based on the transformer architecture [28], which has shown impressive results recently. The model consists of a convolutional neural network-based front-end, a transformer-based encoder, and a transformer-based decoder. A mechanism combining the Connectionist Temporal Classification (CTC) and attention-based sequence to sequence (seq2seq) is used for the optimisation [29]. When estimating the uncertainty, the predictive posterior for each token is expressed as:

$$\begin{aligned} P(y_l|\mathbf{y}_{<l}, \mathbf{x}; \theta^{(m)}) &= \lambda P_{CTC}(y_l|\mathbf{y}_{<l}, \mathbf{x}; \theta^{(m)}) \\ &\quad + (1 - \lambda) P_{seq2seq}(y_l|\mathbf{y}_{<l}, \mathbf{x}; \theta^{(m)}), \end{aligned} \quad (7)$$

where λ is a weighting coefficient.

4. Experiments and Results

Our experiments are conducted on the Noisy Grid corpus [5] and the round one database of Clarity Prediction Challenge (CPC1) [30] which will be introduced in details in the next subsections. For both experiments, ensembles of six ASR models are employed to estimate uncertainty. As the entropy is supposed to be negatively correlated with intelligibility, negative entropy $-\mathcal{H}_S$ is used for evaluation. In addition, we evaluate word correctness score (WCS), defined as the number of words that are correctly recognised divided by the total number of words in the utterance, from ensembles of ASR models.

ASR models used in this work are all finetuned from the SpeechBrain [31] released LibriSpeech model¹ with only different random seeds. The ASR models take 80-channel log mel-filter bank coefficients of an utterance with 16 kHz sampling rate as input features. The convolutional front-end consists of three 2D convolutional layers, and the encoder and the decoder consists of twelve and six multi-head attention transformer blocks, respectively. The weighting coefficient α is set to 0.4 for uncertainty estimation. The calibration temperature T is kept as 1. The top 10 hypotheses within the beam are used for entropy estimation.

Three metrics, including root mean square error (RMSE), normalised cross-correlation coefficient (NCC), and Kendall's Tau coefficient (KT), are used to evaluate the correlation between the intelligibility scores from listening results, which are represented by the WCS, and the ASR WCS, the estimated uncertainty measures \mathcal{C}_S , $-\mathcal{H}_S$. Following the convention of evaluating intelligibility prediction, we report the correlations achieved by applying a logistic mapping function $f(x) = 1/[1 + \exp(ax + b)]$, because RMSE and NCC could be invalid in non-linear cases, and the monotonicity correlation is already of great interest for analysis and optimisation. For the proposed method, the two parameters a and b are optimised

¹huggingface.co/speechbrain/asr-transformer-transformerlm-librispeech

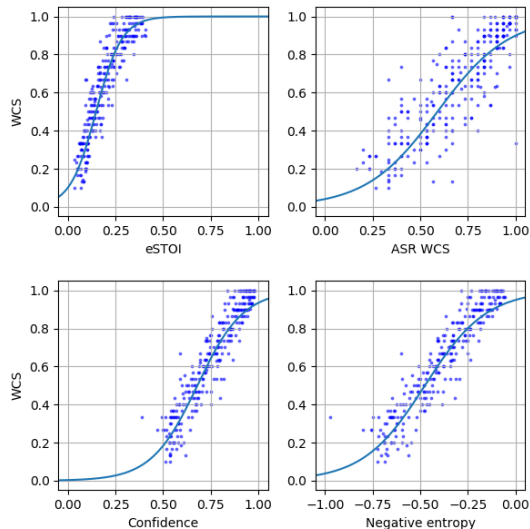


Figure 1: Predicted intelligibility measures on Noisy Grid Corpus test set, including eSTOI, and ASR WCS, confidence C_S , negative entropy $-\mathcal{H}_S$, from the ensemble of ASR models optimised with the NGrid, versus the listening result WCS, in addition with the logistic mapping functions.

in the development set with non-linear least squares², and used in the test set to map the estimated uncertainty to the predicted intelligibility. For the baseline system, the parameters are optimised on the combined training and development sets.

4.1. Noisy Grid corpus

4.1.1. Database

The Noisy Grid corpus is an extension to the original Grid corpus [32] with added speech shaped noise (SSN) at 12 different SNR levels from -14 dB to 40 dB. Each Grid utterance consists of six words following the structure of “command-color-preposition-letter-digit-adverb”, and the words are randomly selected within a limited vocabulary of [4, 4, 4, 25, 10, 4] words. The listeners are asked to identify “color”, “letter”, and “digit” in the listening tests, therefore the WCS for each utterance can only be [0, 1/3, 2/3, 1]. In order to make the distribution of WCS relatively more continuous, the reported WCS is averaged over ten utterances at the same SNR level. The database comprises utterances spoken by 34 speakers, in which the utterances of 22 speakers are used as training set for ASR optimisation, 6 speakers as development set, and 6 speakers as test set. We observed that over 90% utterances, whose SNRs are equal to or higher than 0 dB, have perfect WCS in the listening tests. In order to even the distribution of the database, we report the results of utterances whose SNRs are lower than 0 dB.

4.1.2. Setup

We exploit STOI and extended-STOI (eSTOI) [33]³ as the baseline intelligibility predictors. Both STOI and eSTOI are intrusive measures taking advantage of the correlation between the acoustic features of clean reference signals and corresponding

²docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.curve_fit.html

³<https://github.com/mpariente/pystoi>

Table 1: Correlation evaluation between the listening result WCS and predicted intelligibility measures on Noisy Grid Corpus test set.

	WER	Measure	RMSE ↓	NCC ↑	KT ↑
STOI	-	-	0.154	0.853	0.670
eSTOI	-	-	0.100	0.928	0.762
LS	49.09	C_S	0.172	0.762	0.572
		$-\mathcal{H}_S$	0.166	0.788	0.595
		WCS	0.206	0.607	0.440
CGrid	32.88	C_S	0.224	0.521	0.329
		$-\mathcal{H}_S$	0.235	0.444	0.302
		ASR WCS	0.148	0.825	0.650
DGrid	21.03	C_S	0.098	0.925	0.767
		$-\mathcal{H}_S$	0.099	0.924	0.768
		ASR WCS	0.115	0.901	0.754
NGrid	17.04	C_S	0.093	0.937	0.790
		$-\mathcal{H}_S$	0.094	0.936	0.791
		ASR WCS	0.144	0.844	0.695

degraded processed signals. Because the Grid corpus has a limited vocabulary, the inference of the ASR models is strictly constrained within the Grid dictionary. As the ASR models operate at 16 kHz, the Noisy Grid utterances are downsampled from 25 kHz to 16 kHz.

To investigate the impact of prior knowledge of the ASR models (the data used for ASR optimisation) could have on intelligibility prediction, we employ different ensembles of models including: ASRs finetuned on the training sets of (1) LibriSpeech (LS); (2) clean Grid corpus (CGrid); (3) clean Grid mixed with DEMAND noise [34] at SNRs from -15 dB to 15 dB (DGrid); (4) the original noisy Grid corpus (NGrid). The ensemble of ASR models finetuned on LS are optimised for two epochs, and those finetuned on CGrid, DGrid, NGrid are optimised for 10 epochs.

4.1.3. Results

Table 1 lists the evaluation results on the Noisy Grid test set. Figure 1 shows the eSTOI predicted intelligibility scores, ASR WCS, confidence and negative entropy from the ensemble of ASR models finetuned on NGrid versus the listening result WCS along with their logistic mapping functions. The result shows that the uncertainty estimated by the ensemble of ASR models optimised with NGrid is highly correlated with speech intelligibility and outperforms STOI, eSTOI. In addition, the uncertainty is better at intelligibility prediction than ASR WCS. The confidence is slightly more correlated with intelligibility than entropy in terms of RMSE and NCC, while the entropy performs slightly better in terms of KT.

The word error rates (WER) of Noisy Grid test set for each ensemble of ASR models (which vary by their degree of prior knowledge of the evaluated ASR models) are also shown in Table 1. It shows that a strong prior knowledge of the test data leads to a high correlation between ASR uncertainty and speech intelligibility based on the results of CGrid, DGrid, and NGrid. However, it can be observed that when the ASR models have no knowledge of the noisy signals, the confidence and negative entropy of LS finetuned ensemble could outperform the CGrid finetuned ensemble. It is also worth noting that ASR models optimised on DGrid, i.e., different type of noises from the Noisy Grid test set, could also produce competitive results.

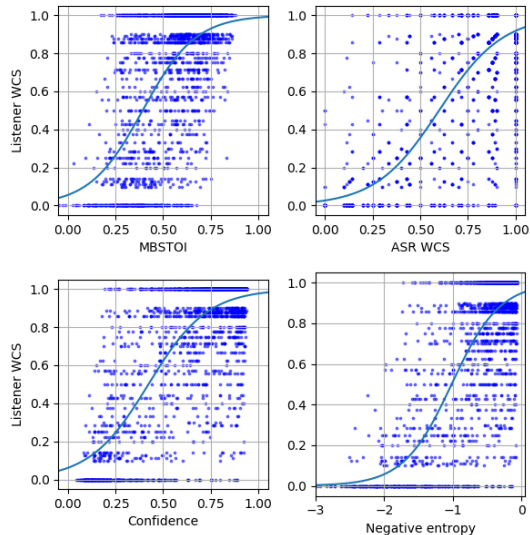


Figure 2: Predicted intelligibility measures on CPC1 closed evaluation set, including the baseline, ASR WCS, and confidence \mathcal{C}_S , negative entropy $-\mathcal{H}_S$, from the ensemble of MSBG+CLS+CPC1 ASR models, versus the listening result WCS, in addition with the logistic mapping functions.

4.2. CPC1

4.2.1. Database

For the purpose of advancing hearing aid intelligibility prediction, the CPC1 database provides a large number of binaural signals and their corresponding responses made by hearing impaired listeners. Each signal corresponds to a noisy scene, which is simulated by mixing a target utterance and a segment of noise in a room, and enhanced by a machine learning hearing-aid system based on the listener’s hearing loss measure. The complete database consists of 6 speakers, 10 hearing aid systems and 27 listeners. Two separate but related tracks are included in CPC1: (1) *closed-set*, in which the evaluation hearing-aid systems and listeners are the same as those in the training data; (2) *open-set*, in which the hearing-aid systems or listeners in the evaluation set are different from those in the training data. Readers are referred to [30] for full details. In both tracks, the training/development scenes are split between 70 % and 30 %, and the results on the extra evaluation set are reported.

4.2.2. Setup

Since in CPC1 the listeners are hearing impaired and the signals are binaural, the CPC1 baseline system employs a combination of Cambridge MSBG hearing loss simulator [35–38] and MB-STOI [39]. The MSBG simulator applies simulated degradation to an input signal according to the hearing loss measures of a listener, and MBSTOI is a refined version of binaural STOI.

To estimate the uncertainty of a binaural signal from an ensemble of ASR models, the signal is resampled to 16 kHz after processing with the MSBG model. Uncertainty of the left and right channel of each binaural signal is estimated independently, and a better ear principle is applied for the binaural uncertainty, i.e., the higher value of \mathcal{C}_S or $-\mathcal{H}_S$ is regarded as the binaural uncertainty. The same better ear rule is also applied to the left and right ASR WCS.

Table 2: Correlation evaluation between the listening result WCS and predicted intelligibility measures on CPC1 evaluation set.

	WER	Measure	RMSE ↓	NCC ↑	KT ↑
<i>Closed-set</i>					
CPC1 Baseline	-	-	0.285	0.621	0.398
Proposed without MSBG	25.17	\mathcal{C}_S	0.241	0.751	0.472
		$-\mathcal{H}_S$	0.239	0.754	0.477
		ASR WCS	0.249	0.730	0.525
Proposed with MSBG	30.33	\mathcal{C}_S	0.234	0.767	0.497
		$-\mathcal{H}_S$	0.233	0.768	0.499
		ASR WCS	0.249	0.731	0.526
<i>Open-set</i>					
CPC1 Baseline	-	-	0.365	0.529	0.391
Proposed with MSBG	30.93	\mathcal{C}_S	0.248	0.729	0.512
		$-\mathcal{H}_S$	0.246	0.734	0.512
		ASR WCS	0.253	0.717	0.530

The pretrained ASR models are first finetuned on the LibriSpeech (LS) for two epochs. Furthermore, they are optimised with LS *train-clean-100* added with noises from the training set in the first round Clarity Enhancement Challenge [40] for 10 epochs. Finally, the models are optimised with the CPC1 training set for another 10 epochs. Therefore, the ASR models possess knowledge of clean, noisy, and processed speech signals. For the *closed-set* experiments, we trained two ensembles of ASR models with and without the MSBG hearing loss model processed signals.

4.2.3. Results

For the CPC1 *closed-set*, the uncertainty estimated from the ensemble of ASR models are more strongly correlated with speech intelligibility than the baseline, and negative entropy gains a slight advantage over confidence. In terms of RMSE and NCC, the uncertainty also outperforms ASR WCS. On the contrary, ASR WCS performs better with regard to KT as WCS are discrete, i.e., *tied* pairs are more likely to appear. In addition, the results show that the MSBG model could provide a slight advantage for intelligibility prediction.

The results on the *open-set* are consistent with those on the *closed-set*. It is also worth noting that, the baseline has a large performance drop as the evaluation signals are very different from the ones in the training set. However, the ASR models are quite robust to this mismatch as the WERs are similar, and achieve similar performances.

5. Conclusions

In this paper, we have shown that the sequence-level uncertainty of DNN-based ASR models is strongly correlated with speech intelligibility. Therefore, the estimated confidence and entropy from an ensemble of ASR models can be used as effective non-intrusive intelligibility predictors. In addition, the uncertainty estimation is unsupervised requiring no explicit references, i.e., no listening WCS nor reference clean signals are needed for training the predictor. The experimental results on two databases show that the proposed method can outperform STOI and its variants, and is better than ASR WCS at intelligibility prediction.

6. References

- [1] I. Holube and B. Kollmeier, "Speech intelligibility prediction in hearing-impaired listeners based on a psychoacoustically motivated perception model," *The Journal of the Acoustical Society of America*, vol. 100, no. 3, pp. 1703–1716, 1996.
- [2] T. Jürgens and T. Brand, "Microscopic prediction of speech recognition for listeners with normal hearing in noise using an auditory model," *The Journal of the Acoustical Society of America*, vol. 126, no. 5, pp. 2635–2648, 2009.
- [3] C. Spille *et al.*, "Predicting speech intelligibility with deep neural networks," *Computer Speech & Language*, vol. 48, pp. 51–66, 2018.
- [4] M. Karbasi *et al.*, "Non-intrusive speech intelligibility prediction using automatic speech recognition derived measures," *arXiv preprint arXiv:2010.08574*, 2020.
- [5] J. Barker and M. Cooke, "Modelling speaker intelligibility in noise," *Speech Communication*, vol. 49, no. 5, pp. 402–417, 2007.
- [6] M. R. Schädler *et al.*, "Matrix sentence intelligibility prediction using an automatic speech recognition system," *International Journal of Audiology*, vol. 54, no. sup2, pp. 100–107, 2015.
- [7] L. Fontan *et al.*, "Automatic speech recognition predicts speech intelligibility and comprehension for listeners with simulated age-related hearing loss," *Journal of Speech, Language, and Hearing Research*, vol. 60, no. 9, pp. 2394–2405, 2017.
- [8] J. B. Allen, "How do humans process and recognize speech?" in *Modern methods of speech processing*. Springer, 1995, pp. 251–275.
- [9] A. Malinin and M. Gales, "Uncertainty estimation in autoregressive structured prediction," in *ICLR*, 2020.
- [10] D. Oneață *et al.*, "An evaluation of word-level confidence estimation for end-to-end automatic speech recognition," in *SLT*. IEEE, 2021, pp. 258–265.
- [11] K. Kalgaonkar *et al.*, "Estimating confidence scores on asr results using recurrent neural networks," in *ICASSP*. IEEE, 2015, pp. 4999–5003.
- [12] A. Ragni *et al.*, "Confidence estimation and deletion prediction using bidirectional recurrent neural networks," in *SLT*. IEEE, 2018, pp. 204–211.
- [13] P. Swarup *et al.*, "Improving asr confidence scores for alexa using acoustic and hypothesis embeddings." in *INTERSPEECH*, 2019, pp. 2175–2179.
- [14] T. H. Falk *et al.*, "A non-intrusive quality and intelligibility measure of reverberant and dereverberated speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 7, pp. 1766–1774, 2010.
- [15] F. Chen *et al.*, "Predicting the intelligibility of reverberant speech for cochlear implant listeners with a non-intrusive intelligibility measure," *Biomedical signal processing and control*, vol. 8, no. 3, pp. 311–314, 2013.
- [16] C. H. Taal *et al.*, "An algorithm for intelligibility prediction of time–frequency weighted noisy speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 7, pp. 2125–2136, 2011.
- [17] A. H. Andersen *et al.*, "A non-intrusive short-time objective intelligibility measure," in *ICASSP*. IEEE, 2017, pp. 5085–5089.
- [18] C. Sørensen *et al.*, "Non-intrusive intelligibility prediction using a codebook-based approach," in *EUSIPCO*. IEEE, 2017, pp. 216–220.
- [19] M. Karbasi *et al.*, "Twin-hmm-based non-intrusive speech intelligibility prediction," in *ICASSP*. IEEE, 2016, pp. 624–628.
- [20] A. H. Andersen *et al.*, "Nonintrusive speech intelligibility prediction using convolutional neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1925–1939, 2018.
- [21] D. Sharma *et al.*, "A data-driven non-intrusive measure of speech quality and intelligibility," *Speech Communication*, vol. 80, pp. 84–94, 2016.
- [22] R. E. Zezario *et al.*, "Stoi-net: A deep learning based non-intrusive speech intelligibility assessment model," in *APSIPA ASC*. IEEE, 2020, pp. 482–486.
- [23] A. M. C. Martinez *et al.*, "Prediction of speech intelligibility with dnn-based performance measures," *Computer Speech & Language*, p. 101329, 2021.
- [24] L. Hoffmann and C. Elster, "Deep ensembles from a bayesian perspective," *arXiv preprint arXiv:2105.13283*, 2021.
- [25] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *ICML*. PMLR, 2016, pp. 1050–1059.
- [26] B. Lakshminarayanan *et al.*, "Simple and scalable predictive uncertainty estimation using deep ensembles," *NeurIPS*, vol. 30, 2017.
- [27] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.
- [28] A. Vaswani *et al.*, "Attention is all you need," in *NeurIPS*, 2017, pp. 5998–6008.
- [29] S. Kim *et al.*, "Joint ctc-attention based end-to-end speech recognition using multi-task learning," in *ICASSP*. IEEE, 2017, pp. 4835–4839.
- [30] J. Barker *et al.*, "The 1st clarity prediction challenge: A machine learning challenge for hearing aid intelligibility prediction," in *INTERSPEECH*, 2022.
- [31] M. Ravanelli *et al.*, "Speechbrain: A general-purpose speech toolkit," *arXiv preprint arXiv:2106.04624*, 2021.
- [32] M. Cooke *et al.*, "An audio-visual corpus for speech perception and automatic speech recognition," *The Journal of the Acoustical Society of America*, vol. 120, no. 5, pp. 2421–2424, 2006.
- [33] J. Jensen and C. H. Taal, "An algorithm for predicting the intelligibility of speech masked by modulated noise maskers," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 11, pp. 2009–2022, 2016.
- [34] J. Thiemann *et al.*, "The diverse environments multi-channel acoustic noise database (demand): A database of multichannel environmental noise recordings," in *Proceedings of Meetings on Acoustics ICA2013*, vol. 19, no. 1. Acoustical Society of America, 2013, p. 035081.
- [35] T. Baer and B. C. J. Moore, "Effects of spectral smearing on the intelligibility of sentences in noise," *JASA*, vol. 94, no. 3, pp. 1229–1241, 1993.
- [36] —, "Effects of spectral smearing on the intelligibility of sentences in the presence of interfering speech," *JASA*, vol. 95, no. 4, pp. 2277–2280, 1994.
- [37] B. C. J. Moore and B. R. Glasberg, "Simulation of the effects of loudness recruitment and threshold elevation on the intelligibility of speech in quiet and in a background of speech," *JASA*, vol. 94, no. 4, pp. 2050–2062, 1993.
- [38] M. A. Stone and B. C. Moore, "Tolerable hearing aid delays. i. estimation of limits imposed by the auditory path alone using simulated hearing losses," *Ear and Hearing*, vol. 20, no. 3, pp. 182–192, 1999.
- [39] A. H. Andersen *et al.*, "Refinement and validation of the binaural short time objective intelligibility measure for spatially diverse conditions," *Speech Communication*, vol. 102, pp. 1–13, 2018.
- [40] S. Graetzer *et al.*, "Clarity-2021 challenges: Machine learning challenges for advancing hearing aid processing," in *INTERSPEECH*, 2021, pp. 686–690.