

INFERENCE AND DENOISE: CAUSAL INFERENCE-BASED NEURAL SPEECH ENHANCEMENT

Tsun-An Hsieh¹, Chao-Han Huck Yang², Pin-Yu Chen³, Sabato Marco Siniscalchi^{2,4}, Yu Tsao¹

¹Research Center for Information Technology Innovation, Academia Sinica, Taiwan

²Georgia Institute of Technology, GA, USA; ³IBM Research, NY, USA

⁴Computer Engineering School, Norwegian University of Science and Technology, Norway

ABSTRACT

This study addresses the speech enhancement (SE) task within the causal inference paradigm by modeling the noise presence as an intervention. Based on the potential outcome framework, the proposed causal inference-based speech enhancement (CISE) separates clean and noisy frames in an intervened noisy speech using a noise detector and assigns both sets of frames to two mask-based enhancement modules (EMs) to perform noise-conditional SE. Specifically, we use the presence of noise as guidance for EM selection during training, and the noise detector selects the enhancement module according to the prediction of the presence of noise for each frame. Moreover, we derived a SE-specific average treatment effect to quantify the causal effect adequately. Experimental evidence demonstrates that CISE outperforms a non-causal mask-based SE approach in the studied settings and has better performance and efficiency than more complex SE models.

Index Terms— Observational Inference, Deep Causal Inference, and Speech Enhancement

1. INTRODUCTION

Recent advances in neural network-based speech enhancement (SE) have demonstrated impressive performance in terms of speech quality and intelligibility scores, such as perceptual evaluation of speech quality (PESQ) [1] and short-time objective intelligibility (STOI) [2] in various speech applications. However, modern SE approaches [3, 4, 5, 6, 7] do not explicitly take the presence of noise into account, and real-world acoustic scenarios often encounter inevitable observational uncertainties. For instance, a meeting could be abruptly disconnected, or a session could be disrupted by temporary noise from the external environment. That is, noise intervention may not affect the entire speech waveform. In such a scenario, conventional neural SE solutions may be unreliable for handling the intermittent/sporadic nature of the noise. By contrast, causal inference (CI) [8] may be a viable paradigm for performing SE. The design of an end-to-end neural SE model within the CI framework is the research question addressed in this study. Causal inference-based machine learning techniques are often featured with the ability to identify unobserved factors or features (also known as confounding variables) with improved model prediction and generalization, i.e., CI-based models are proven to be advantageous of tackling unseen data hence dependable [9, 10]. Furthermore, machine learning models that satisfy the CI training objectives stand to benefit from additional interpretable scores to formally quantify the causal effects, for example, in treatment effect estimation. Previous studies [11, 12] have demonstrated that learning to measure causal variables empowers effective

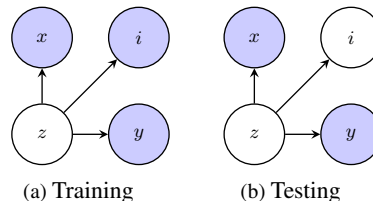


Fig. 1: Causal graphical model (CGM) for the training phase (a) and the testing phase (b). The blue nodes x and y are observable (e.g., noisy speech and speech intelligibility scores). Node z , colored white, is not observable as a parameterized latent variable. Node i , colored blue in (a), is only observable during training, and to be inferred by proposed causal model in (b) the testing time.

model selection. Meanwhile, similar designs are sparse expert models [13, 14]. These approaches divide a large task into small sub-tasks by allocating data categorized in different attributes to several local expert models. Nevertheless, those models do not take into account the assumption of a causal graph; therefore, they can not be evaluated under a formal causal learning settings with treatment effect analysis. Finally, causal effect measurements [15] could incorporate statistical refutation tests to design reliable prediction models.

This study focuses on develop a SE system using the potential outcome framework [16] shown in Fig. 1. Under this paradigm, we model the confounding variable z that causally impacts SE performance y , instead of directly modeling the correlation between the noisy and clean speech. To this end, a two-stage training procedure is adopted to attain satisfactory SE results by leveraging auxiliary intervention labels. To the best of our knowledge, the proposed training process and architecture are the first attempts to introduce CI into an end-to-end neural SE system. In contrast to previous studies that adopt prior causal features [17, 18] to improve system performance, we focus on the inference for vector-to-vector regression network of SE [6, 19] by employing an auxiliary sub-task for state estimation, i.e., noise detection, which could leverage the advantages of causal inference. We design our observational inference enhanced network based on this causal neural architecture. Our contributions is four-fold: (i) a novel neural SE architecture based on causal inference aimed at handling complicated noisy condition is presented; (ii) a novel quantitative measure for the causal effect of the selected intervention is devised; (iii) effectiveness of the proposed solution is demonstrated by showing that CISE outperforms both the non-causal counterpart model and other techniques leveraging complicated EMs in terms of both quality and intelligibility, and (iv) the proposed system merely uses 2.64% of the computational time and 4.96% of GPU memory as compared with the largest CISE variant.

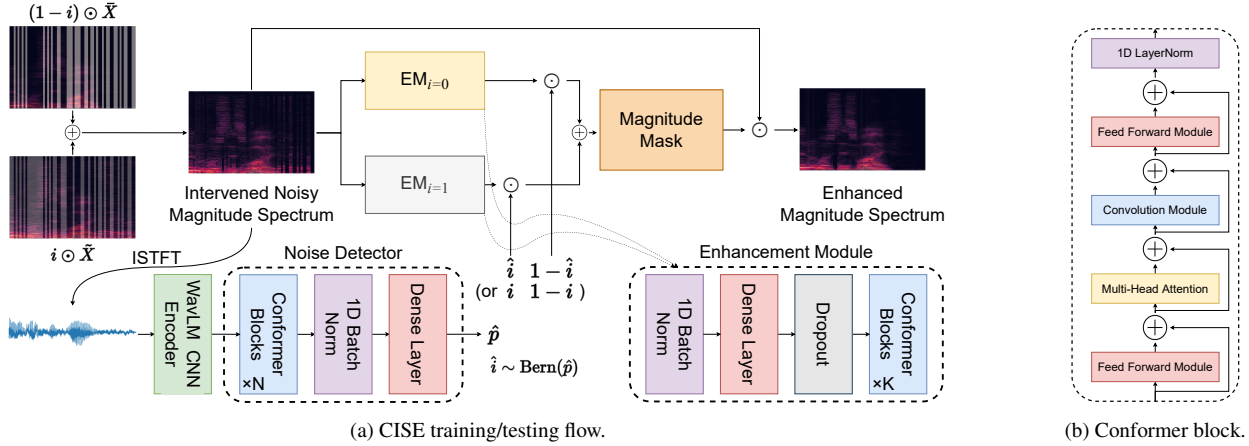


Fig. 2: Proposed causal inference-based SE architecture. Causal inference-based SE (CSIE) is based on two independent neural enhancement module (EM) without parameter sharing.

2. BACKGROUND

2.1. Causal Inference & Representation Learning

Based on Pearl’s causal hierarchy theorem [20], modeling machine learning problems at a higher causal level (e.g., interventional) could provide access to more useful information while extracting relevant features or in learning from proxy variables [21]. Recently, causality learning [8, 22] has been proven successful when combined with representation learning for feature extraction and probabilistic inference in sequence modeling [23]. For instance, causal convolution [17, 18] is a successful approach widely applied in speech synthesis. Moreover, observational inference empowered neural models attaining top performances in clinical learning [24], sequence modelings [25], and robust reinforcement learning [21]. Causal inference [8] is another mainstream approach to causal learning that focuses on learning robust proxy variables and inferencing under unseen dynamics, such as confounding variables.

2.2. Treatment Effect Quantification

To quantify the causal effects of the intervention on the outcomes of interest in a randomized controlled trial (RCT), the average treatment effect (ATE) [26] is a metric often adopted. ATE measures the mean difference between the potential outcomes of the treatment and control groups and is formally defined as

$$\text{ATE} = \mathbb{E}[y|i = 1] - \mathbb{E}[y|i = 0], \quad (1)$$

where i is a binary label that indicates the occurrence of an intervention when its value is equal to one. The two terms on the right-hand side of Eq. (1) denote the expected outcome over the population for the treatment and control groups, respectively. Briefly, a positive/negative ATE implies that the selected intervention has a positive or negative causal effect on the outcome of interest.

3. CAUSAL INFERENCE-BASED SE SYSTEM

3.1. Problem Definition

In general, the causal graphical model of the potential outcome framework [16] can be represented as Fig. 1. In the CGM, x is a noisy observation; y is the outcome for that observation; i denotes

an intervention, and z is a hidden confounding variable that affects the other three variables. Note that because z is learnt implicitly, we cannot know the exact meaning of it. In training, i is an observable variable used as a guidance to learn z efficiently [21]; however, i is unobservable and has to be predicted by the model at a testing time.

As we focus on forming the SE problem in the interventional causation, we integrate the SE problem with the potential outcome framework, as shown in Fig. 1. Here, we denote variables in the time–frequency domain in uppercase letters. For CISE, we intervene the clean speech by adding noise into it. Therefore, we define the intervened noisy speech $X = (1 - i) \odot \bar{X} + i \odot \tilde{X}$, where \bar{X} and \tilde{X} are the clean and noisy speech, and i is a 0/1-mask that indicates in what frames the intervention occurs. In other words, the noise signal does not prevail over all time stamps. Fig. 2a (a) illustrates how i intervenes the clean speech. The outcome y can be any measure of speech quality, intelligibility, or distance. As for CISE training, we set y as the l_1 distortion distance. Therefore, CISE learns to model z guided by a predefined intervention, which is further used to predict an intervention in the testing phase.

3.2. Enhancement Modules and Noise Detector

The convolution-augmented Transformer (Conformer) [27] has been proven effective in various speech applications [27, 28], including SE [29, 30]; therefore, we adopt it as a main component of our SE system. A Conformer consists of a series of half-step feed forward modules, a multi-head attention, and a convolution module as shown in Figure 2b. We also employed half-step feed forward layers [31] and relative sinusoidal positional encoding for improved performance. In Fig. 2a, two enhancement modules manage the processing of speech frames belonging to the treatment and control groups, respectively. Similar to [28], each enhancement module consists of a batch normalization layer, linear transformation, followed by a few Conformer blocks.

In testing, we need to identify when an intervention occurs. Since the presence of noise is regarded as an intervention in our study, the identification approach is implemented as a noise detector. Although Mel-frequency cepstral coefficients (MFCCs) could be employed as input features for noise detection, we observed that MFCCs led to a severe overfitting of the training data, yielding a 36% difference between the training and testing accuracy - domain mismatch may have caused that. To circumvent overfitting, we

Table 1: Speech qualities and intelligibility of CISE on the curated VoiceBank–DEMAND dataset. DA is the noise detection accuracy.

Model	DA	PESQ	CSIG	CBAK	COVL	STOI
Oracle	1.00	3.21	4.74	4.36	4.04	0.98
CISE	0.92	3.15	4.70	4.27	3.98	0.97
CISE-C [34]	0.92	2.82	4.18	3.75	3.52	0.95
CISE-D [35]	0.92	2.78	4.18	3.71	3.49	0.96
CISE-M [36]	0.92	2.60	3.98	3.64	3.30	0.95
MFCC	0.56	2.56	4.15	3.61	3.38	0.95
Random \hat{i}	0.50	2.50	4.02	3.50	3.28	0.94
Vanilla	–	2.44	3.68	3.20	3.05	0.93
$1 - i$	0.00	2.23	3.57	3.05	2.89	0.92

used WavLM [32] CNN as an encoder to extract more generalizable embedding. Since noise presence prediction can be thought of as a sound event detection task, in which Conformers attained top accuracies [28, 33], we used a Conformer-based noise detector. In the bottom left in Fig. 2a, the noise detector is fed with embedding extracted by the WavLM module. The temporal dependencies among embedding are then modeled by the Conformer blocks and mapped into a sequence of two-dimensional vectors, representing the probabilities of speech fragments being noisy or clean. Finally, the predicted intervention \hat{i} is sampled from a Bernoulli distribution using the predicted probability \hat{p} .

3.3. CISE Training and Testing

As shown in Fig. 2a, we take an intervened noisy speech magnitude spectrum, a staggered combination of the clean and the noisy speech frames in time, as the input. The intervention i , which guides the training procedure, is generated from a Bernoulli distribution (see Section 4.1 for details); in testing, the noisy detector, shown in Fig. 2a(a), generates \hat{i} in order to select different enhancement modules. Next, two Conformer-based enhancement modules, namely $EM_{i=0}$ and $EM_{i=1}$, which estimate magnitude masks belonging to the treatment ($i = 1$) and control ($i = 0$) groups, respectively, remix the predicted magnitude masks based on the intervention labels, and the remixed magnitude mask is multiplied by the intervened noisy magnitude spectrum in an element-wise manner. The intervention labels (or predictions during testing) are then used for mixing the outputs of different enhancement modules. Finally, the SE process is accomplished by multiplying the intervened noisy magnitude spectrum by the magnitude mask.

As CISE simultaneously learns to identify noise occurrences and to perform enhancement, we characterize CISE training as a multi-task learning process, and thus we formulate the loss function as

$$\mathcal{L}_{total}(\bar{X}, \hat{X}, i, \hat{i}) := \mathcal{L}_{l_1}(\bar{X}, \hat{X}) + \mathcal{L}_{CE}(i, \hat{i}), \quad (2)$$

where \hat{X} denotes enhanced speech. We select the l_1 distance and cross-entropy (CE) for the regression and classification tasks, respectively. For magnitude mask estimation, we simply minimize the l_1 distance between the enhanced and target spectra. Meanwhile, CISE training also minimizes the CE between the distribution of the predicted interventions \hat{i} and that of the corresponding ground truth.

3.4. ATE for Speech Enhancement

For SE tasks, ATE can be defined for each chosen evaluation metric. Specifically, for a given metric \mathcal{M} (e.g., PESQ or STOI), ATE is

Table 2: Computational overheads. Each entry shows the time and memory usage of processing 1 second speech signal. The left and right of the slash shows the usage for batch size of 1 and 16. Both forward and backward pass are considered.

	CISE	CISE-C	CISE-D	CISE-M
CPU Time (s)	2.20/2.39	2.50/49.79	1.79/53.19	2.56/50.46
GPU Time (s)	0.03/0.12	0.142/4.54	0.08/3.90	0.10/4.02
GPU Mem. (GB)	2.19/2.91	4.47/58.62	2.55/23.11	2.76/29.51

defined as

$$\text{ATE}_{\mathcal{M}} = \mathbb{E}[\mathcal{M}(\bar{x}, \mathcal{E}(x)) | i = 1] - \mathbb{E}[\mathcal{M}(\bar{x}, \mathcal{E}(\bar{x})) | i = 0]. \quad (3)$$

In Equation (3), \bar{x} denotes unobservable clean speech, x is the observable noisy version of \bar{x} , and $\mathcal{E}(x)$ is the CISE-enhanced speech signal. For example, a positive ATE_{PESQ} implies that, on average, the addition of noise has a positive causal effect on the enhanced speech in terms of quality. The ATE is independent of optimization; therefore, CISE does not intentionally increase the ATE by distorting clean speech.

4. EXPERIMENTAL SETUP & RESULTS

4.1. Data Curation

To evaluate the proposed CISE approach, we use the Voice Bank–DEMAND dataset [37], which we curate to make it suitable for causal inference. In the original Voice Bank–DEMAND, clean speech from 30 speakers are recorded in a studio room at sample rate of 48 KHz. Among those speakers, speech material from 28 speakers is used for training, and the rests are used for testing. The training set includes 10 types of noises added to the clean speech at 4 signal-to-noise-ratio (SNR) levels, ranging from 0 dB to 15 dB. For the test set, 5 unseen noises are added to the clean speech at SNR from 2.5 dB to 17.5 dB.

For CISE, we need to know where the intervention takes place. Therefore, an additional information is needed, namely a label indicating the presence of noise in a given speech frame. To this end, we remix noisy data through combining clean and noisy speech $X = (1 - i) \odot \bar{X} + i \odot \tilde{X}$ where $i \sim \text{Bern}(0, 1)$, $i \in \{0, 1\}$ indicates of the appearance of noise, and i is drawn from a Bernoulli distribution with $p(i = 1) = 0.5$. X is the intervened noisy speech randomly mixed by the clean speech \bar{X} and the original noisy speech \tilde{X} . With the intervened noisy speech X and the corresponding i , a causal inference-based SE system can be implemented as shown in Fig. 2a. Technically, we sample i with the length of X , and then repeat each time stamp to the frequency dimension of X .

4.2. Experimental Setup

To match the down sampling rate and the features size of the WavLM CNN encoder, we set the size of Fourier transform to 1023, the length of the analysis window to 1023 sample points (approximately 0.064 seconds), and the step of sliding window to 320. The dropout rate is set to 25%, and we use two layers of Conformer blocks to encode temporal dependencies. For intervention prediction, we use the WavLM CNN encoder to extract a 512-dimensional general purpose representation of the speech, which is then used for noise detection. The noise detector comprises four Conformer blocks stacked together, along with a batch normalization layer over the channels

Table 3: Causal inference explanation for speech intelligibility indexes. ATEs with controlled noise detection accuracy p . The left columns denote quality/intelligibility scores; on the right are the ATEs of the corresponding metrics.

Accuracy p	PESQ	CSIG	CBAK	COVL	STOI	SSNR	Average Treatment Effect					
							PESQ	CSIG	CBAK	COVL	STOI	SSNR
0.0	2.23	3.57	3.05	2.89	0.92	8.37	-1.1872	-0.8659	-1.2546	-1.1394	-0.0477	-9.819
0.1	2.27	3.65	3.12	2.96	0.92	9.04	-1.1275	-0.7827	-1.1659	-1.0606	-0.0423	-9.0432
0.3	2.37	3.83	3.29	3.11	0.93	10.58	-0.9617	-0.5979	-0.9384	-0.8662	-0.0314	-7.0541
0.5	2.51	4.03	3.50	3.29	0.94	12.50	-0.7247	-0.3834	-0.6063	-0.6081	-0.0193	-3.9817
0.7	2.69	4.26	3.77	3.51	0.95	14.95	-0.3772	-0.1548	-0.1055	-0.2598	-0.0063	1.2719
0.9	2.99	4.56	4.13	3.82	0.97	17.99	0.2283	0.0556	0.4185	0.2610	0.0099	10.0948
1.0	3.21	4.74	4.36	4.04	0.98	19.83	0.9245	0.1013	0.4860	0.5576	0.0204	16.4458

of features vectors and a fully-connected layer to reduce the dimensionality of the hidden states from 512 to 2, representing $p(i = 0)$ and $p(i = 1)$, respectively. For convergence stability, in the training stage, we use the intervention labels for enhancement module switching; however, we only use the predicted interventions during testing. Finally, a standard Adam [38] with learning rate 10^{-4} , $\alpha = 0.9$, and $\beta = 0.999$ is adopted for the optimization.

4.3. Speech Quality and Intelligibility Results

In this section, we compare the proposed CISE system with its variants using the state-of-the-art SE models as EMs and analyze the importance of the noise detector. We report several metrics often used to assess the speech quality and intelligibility in Table 1. PESQ estimates the perceptual speech quality by assigning a score ranging from -0.5 to 4.5. The STOI is an intrusive measure of the intelligibility of degraded speech signals. CSIG, CBAK, and COVL are composite measures [39] of speech quality. CSIG focuses on the quality of foreground speech; conversely, CBAK estimates the extent of the intrusion of background noise, with a higher score indicating less intrusion; COVL evaluates overall quality combined with previous scores. SSNR represents for segmental signal-to-noise ratio.

In Table 1, from top to bottom, we report the speech quality and intelligibility scores of different SE systems. Oracle denotes the ideally achievable result with CISE when the intervention labels are known at the testing time. From the second to the fifth rows are CISE and its variants, where CISE, CISE-C, CISE-D, and CISE-M use Conformer-based (proposed), CMGAN [34], DEMUCS [35], and MANNER [36] as EMs, respectively. Comparing these CISE variants, we observed that despite sharing the same detection accuracy, the results of using different EMs for CI training is uneven. CISE has overall highest scores; CISE-C and CISE-D have a similar behavior; CISE-M performs even worse than Random \hat{i} (described in the following context) for CSIG and COVL. In addition, CISE with Conformer building blocks attains better results using less processing time, and with a lower memory consumption, as shown in Table 2. Accordingly, we conclude that the proposed CISE with Conformer building blocks best matches the potential outcome framework among the studied settings. Starting from the sixth row, MFCC refers to a noise detector that uses MFCC features as inputs with identical classification structure as CISE presented in Section 3.2; Random \hat{i} refers to a CISE architecture with a malfunctioning noise detector that samples \hat{i} at random following a Bernoulli distribution as described in Section 4.1. As expected, Random \hat{i} achieves slightly worse results than MFCC since the detection accuracy of Random \hat{i} is slightly higher than the detection accuracy of MFCC. In the eighth row, Vanilla leverages the same EM used by the CISE without the module switching mechanism; that is, only a single enhancement

module is present. The last row, $1 - i$, offers the worst case of the CISE system, where clean and noisy frames are swapped (i.e., the detection accuracy is 0), and thereby assigned to an incorrect enhancement module in Fig. 2a. Through controlling the EM structure and comparing the CISE with the last four rows, we see that, CISE outperforms its non-causal version and other cases by taking advantage of being informed by the hidden factor behind observable variables, and thus CISE yields the best overall results. Notably, although the same Conformer-based SE architecture is used, Random \hat{i} can still outperform Vanilla even with a randomized intervention.

4.4. ATE with Speech Evaluation Metrics

In Table 3, we control the detection accuracy p (shown in the left column) moving from 0.0 to 1.0. The speech quality and intelligibility scores discussed in Section 4.3 are reported in the second to the seventh columns; in the seventh column, the SSNR is presented. A CISE with higher detection accuracy leads to better results, demonstrating the importance of accurate intervention prediction. The remaining columns in Table 3 report the ATE values for each of the metrics given in the second to the sixth row. The ATE represents the causal effect of the selected intervention (i.e., the presence of noise) on an outcome (i.e., metric score). As each ATE monotonically increases as p increases, the detection accuracy and ATE values are positively correlated. In addition, a p greater than 0.9 is required to obtain positive ATEs in each metric, showing a high demand for an accurate noise detector for CISE inference.

5. CONCLUSION

This study provided an effective and efficient solution to incorporate the neural SE training and inference with the potential outcome framework and consequently leveraged the power of causal inference. The experimental evidence showed that the proposed Conformer-based CISE system is capable of conducting outstanding performances regarding quality, intelligibility, and computational overheads with respect to computational time and memory usage. Furthermore, CISE outperforms its non-causal counterpart and all the other variants using powerful state-of-the-art SE as EMs by a large margin, showing the imperative necessity in searching the suitable model for the EMs. Aside from the model performance, we also defined an ATE specialized for SE to quantify the causal effect on a metric given an intervention, and by investigating the change in ATE along with the manipulated detection accuracy, we argued that the accurate intervention prediction is crucial for inference since it causally impacted the performance of a SE system. Our implementation will be available at: <https://github.com/alexiehta/Causal-SE>.

6. REFERENCES

- [1] A.W. Rix, J.G. Beerends, M.P. Hollier, and A.P. Hekstra, “Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs,” in *Proc. ICASSP*, 2001.
- [2] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, “A short-time objective intelligibility measure for time-frequency weighted noisy speech,” in *Proc. ICASSP*. IEEE, 2010.
- [3] X. Lu et al., “Speech enhancement based on deep denoising autoencoder,” *Proc. Interspeech*, 2013.
- [4] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, “A regression approach to speech enhancement based on deep neural networks,” *IEEE/ACM TASLP*, vol. 23, no. 1, pp. 7–19, 2014.
- [5] J. Wang and S. Li, “Self-attention mechanism based system for dcase 2018 challenge task 1 and task 4,” in *DCASE*, 2018.
- [6] J. Qi, J. Du, S. M. Siniscalchi, X. Ma, and C. Lee, “On mean absolute error for deep neural network based vector-to-vector regression,” *IEEE SPL*, 2020.
- [7] W. Hartmann et al., “A direct masking approach to robust asr,” *IEEE TASLP*, vol. 21, no. 10, pp. 1993–2005, 2013.
- [8] J. Pearl, “Causal inference,” *Causality: objectives and assessment*, pp. 39–58, 2010.
- [9] C. Louizos, U. Shalit, J. M. Mooij, D. Sontag, R. Zemel, and M. Welling, “Causal effect inference with deep latent-variable models,” *Proc. NeurIPS*, vol. 30, 2017.
- [10] F. Johansson, U. Shalit, and D. Sontag, “Learning representations for counterfactual inference,” in *Proc. ICML*, 2016.
- [11] S. Athey, “Machine learning and causal inference for policy evaluation,” in *Proc. KDD*, 2015, pp. 5–6.
- [12] S. Tang and J. Wiens, “Model selection for offline reinforcement learning: Practical considerations for healthcare settings,” in *Proc. Mach Learn Res.*, 2021.
- [13] Aswin Sivaraman and Minje Kim, “Sparse mixture of local experts for efficient speech enhancement,” in *Proc. Interspeech*, 2020.
- [14] Kenichi Kumatani, Robert Gmyr, Felipe Cruz Salinas, Linqun Liu, Wei Zuo, Devang Patel, Eric Sun, and Yu Shi, “Building a great multi-lingual teacher with sparsely-gated mixture of experts for speech recognition,” *arXiv preprint arXiv:2112.05820*, 2021.
- [15] J. Pearl, “The seven tools of causal inference, with reflections on machine learning,” *Communications of the ACM*, vol. 62, no. 3, pp. 54–60, 2019.
- [16] Donald B Rubin, “Causal inference using potential outcomes: Design, modeling, decisions,” *Journal of the American Statistical Association*, vol. 100, no. 469, pp. 322–331, 2005.
- [17] P. Agrawal, J. Carreira, and J. Malik, “Learning to see by moving,” in *Proc. ICCV*, 2015, pp. 37–45.
- [18] A. van den Oord et al., “Wavenet: A generative model for raw audio,” in *Prof. ISCA SSW*, 2016, pp. 125–125.
- [19] S. M. Siniscalchi, “Vector-to-vector regression via distributional loss for speech enhancement,” *IEEE SPL*, 2021.
- [20] E. Bareinboim, J. D. Correa, D. Ibeling, and T. Icard, “On pearl’s hierarchy and the foundations of causal inference,” in *Probabilistic and Causal Inference: The Works of Judea Pearl*, pp. 507–556. 2022.
- [21] C.-H. H. Yang, I.-T. Hung, Y. Ouyang, and P.-Y. Chen, “Training a resilient q-network against observational interference,” in *Proc. AAAI*, 2022.
- [22] J. M Robins, A. Rotnitzky, and L. P. Zhao, “Analysis of semiparametric regression models for repeated outcomes in the presence of missing data,” *Journal of the American statistical Association*, vol. 90, no. 429, pp. 106–121, 1995.
- [23] V. Melnychuk, D. Frauen, and S. Feuerriegel, “Causal transformer for estimating counterfactual outcomes,” in *Proc. ICML*, 2022.
- [24] T. W Killian, M Ghassemi, and S. Joshi, “Counterfactually guided policy transfer in clinical settings,” in *Conference on Health, Inference, and Learning*. PMLR, 2022, pp. 5–31.
- [25] X. Wang et al., “Inferbert: a transformer-based causal inference framework for enhancing pharmacovigilance,” *Frontiers in Artificial Intelligence*, vol. 4, 2021.
- [26] P. W Holland, “Statistics and causal inference,” *Journal of the American statistical Association*, vol. 81, no. 396, 1986.
- [27] A. Gulati et al., “Conformer: Convolution-augmented transformer for speech recognition,” *arXiv preprint arXiv:2005.08100*, 2020.
- [28] K. Miyazaki, T. Komatsu, T. Hayashi, S. Watanabe, T. Toda, and K. Takeda, “Conformer-based sound event detection with semi-supervised learning and data augmentation,” *DCASE Workshop*, vol. 1, pp. 4, 2020.
- [29] S. Kataria, J. Villalba, and N. Dehak, “Perceptual loss based speech denoising with an ensemble of audio pattern recognition and self-supervised models,” in *Proc. ICASSP*, 2021.
- [30] E. Kim and H. Seo, “SE-Conformer: Time-Domain Speech Enhancement Using Conformer,” in *Proc. Interspeech 2021*, 2021, pp. 2736–2740.
- [31] Y. Lu et al., “Understanding and improving transformer from a multi-particle dynamic system point of view,” *arXiv preprint arXiv:1906.02762*, 2019.
- [32] S. Chen et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing,” *arXiv preprint arXiv:2110.13900*, 2021.
- [33] Q. Wang et al., “A four-stage data augmentation approach to resnet-conformer based acoustic modeling for sound event localization and detection,” *arXiv:2101.02919*, 2021.
- [34] Ruizhe Cao, Sherif Abdulatif, and Bin Yang, “CMGAN: Conformer-based Metric GAN for Speech Enhancement,” in *Proc. Interspeech 2022*, 2022, pp. 936–940.
- [35] A. Défossez, “Hybrid spectrogram and waveform source separation,” in *Proc. ISMIR*, 2021.
- [36] H.-J. Park et al., “Manner: Multi-view attention network for noise erasure,” in *Proc. ICASSP*, 2022, pp. 7842–7846.
- [37] C. Valentini-Botinhao, X. Wang, S. Takaki, and J. Yamagishi, “Speech Enhancement for a Noise-Robust Text-to-Speech Synthesis System Using Deep Recurrent Neural Networks,” in *Proc. Interspeech*, 2016.
- [38] D. P Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [39] Y. Hu and P. Loizou, “Evaluation of objective quality measures for speech enhancement,” *IEEE TASLP*, vol. 16, 2008.

Appendix

A. Causal Hierarchy and Speech Processing

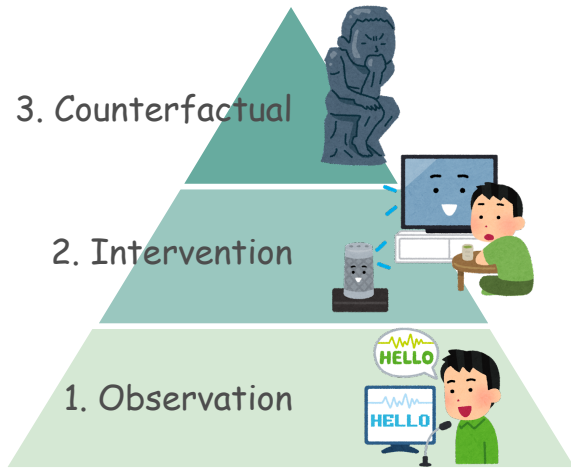


Fig. 3: Automatic speech recognition (ASR) examples at each level of Pearl’s causal hierarchy [20]. At the observation level, given an input “Hello,” an ASR algorithm recognizes speech by matching the input with the most likely word in its dictionary. Based on the first level, an ASR algorithm at the intervention level assumes that all observable variables are noisy and searches for a robust proxy variable representing words in the input speech. Given certain types of interventions, an ASR algorithm at the counterfactual level needs to imagine the following: “*what would the prediction be if different noise signals are involved?*” Nevertheless, examples at this level require further investigation.

Fig. 3 depicts Pearl’s causal hierarchy in the context of automatic speech recognition (ASR). At the *observation level*, a non-causal ASR system is trained by maximizing the likelihood between the input speech and the corresponding labels (i.e., word index). At the *intervention level*, the ASR system learns robust proxy variables from the intervened data, resulting in a better recognition performance when noise is involved. At the *counterfactual level*, one considers “*What would the outcome be if I took another action?*” However, ASR applications at this level are yet difficult to deploy. Inferring this logic-aware information and adopting different effective strategies could be learned naturally through human perception. As a summary, how to design an end-to-end neural speech model equipped with the **power of “inference”** is still an open topic and deserved more investigation. For an SE system, it aims to improve speech quality and intelligibility of a speech signal. One major goal is to reduce noise contaminated speech. Given a clean speech signal \hat{x}_t and a noise signal n_t at time index t , a DL-based SE model tries to learn a mapping from noisy speech signal $x_t = \hat{x}_t + n_t$ to \hat{x}_t so as to acquire clean speech. However, this formation could overly simplify the process of SE since, in many cases, noises are not prevailing in the whole utterance and likely unseen at test phase.

In the regime of causal learning, algorithms aim to model the hidden variable z , a confounder that obfuscates level 1 learning methods during testing, to achieve more accurate outcome prediction. In the scenario of SE, we define x as noisy speech, y as the likelihood of the enhanced and clean speech, and z can be viewed as some concepts of clean speech that SE model wants to learn. By

Algorithm 1 Causal Inference-based SE system

- 1. Inputs:** enhancement module, \mathcal{G} ; noise detector, \mathcal{F} ; speech data, \mathcal{D} ; #Iteration, K ; clean speech \bar{X} ; original noisy speech \tilde{X} ; intervened noisy speech X ; enhanced speech \hat{X} ; ISTFT inverse short-time Fourier transform.
- 2. Randomly Initialize Weights:** $\theta_{i=0}^{(0)}$, $\theta_{i=1}^{(0)}$, and $\theta_{cls}^{(0)}$.
- 3. Training Iteration:**
for $\bar{X}, \tilde{X} \sim \mathcal{D}$; k^{th} iteration **do**
 $i \sim \text{Bern}(p = 0.5)$
 $X \leftarrow (1 - i) \odot \bar{X} + i \odot \tilde{X}$
if training **then**
 $\hat{X} \leftarrow (1 - i) \odot \mathcal{G}(X; \theta_{i=0}^{(n)}) + i \odot \mathcal{G}(X; \theta_{i=1}^{(n)})$
 $\hat{i} \leftarrow \mathcal{F}(\text{ISTFT}(X); \theta_{cls}^{(n)})$
 $\theta_{i=0}^{(k+1)} \leftarrow \theta_{i=0}^{(k)} - \gamma \nabla_{\theta_{i=0}^{(k)}} \mathcal{L}_{l1}(\bar{X}, \hat{X})$
 $\theta_{i=1}^{(k+1)} \leftarrow \theta_{i=1}^{(k)} - \gamma \nabla_{\theta_{i=1}^{(k)}} \mathcal{L}_{l1}(\bar{X}, \hat{X})$
 $\theta_{cls}^{(k+1)} \leftarrow \theta_{cls}^{(k)} - \gamma \nabla_{\theta_{cls}^{(k)}} \mathcal{L}_{CE}(i, \hat{i})$
else
 $\hat{i} \leftarrow \mathcal{F}(\text{ISTFT}(x); \theta_{cls}^{(N)})$
 $\hat{X} \leftarrow (1 - \hat{i}) \odot \mathcal{G}(X; \theta_{i=0}^{(N)}) + \hat{i} \odot \mathcal{G}(X; \theta_{i=1}^{(N)})$
end if
end for

giving i , CISE can learn confounding variable z more efficiently and robustly.

B. Algorithm of Causal Inference-based SE

As presented in Algorithm 1, the intervention label i is used to train the enhancement modules and noise detector. During testing, the predicted intervention \hat{i} is used to mix the estimated mask of each enhancement module. Therefore, the effectiveness significantly depends on the noise detector performance, and the modeling of the confounding variables.

Acknowledgement

The authors want to thank insightful comments from Prof. Minje Kim, Indiana University Bloomington, on the preliminary draft.