

MetricGAN+KAN: Kolmogorov-Arnold Networks in Metric-Driven Speech Enhancement Systems

Yemin Mai and Stefan Goetze

Speech and Hearing (SPandH) group, School of Computer Science, The University of Sheffield, Sheffield, United Kingdom
 {ymai5, s.goetze}@sheffield.ac.uk

Abstract—Neural-network-based speech enhancement (SE) approaches have shown to be particularly powerful in combination with perceptually motivated metrics to produce high-quality enhanced speech signals. Among these deep learning (DL)-based SE models, MetricGAN and its extension can generate output signals directly optimising quality metrics. The recently proposed Kolmogorov-Arnold networks (KANs) with learnable activation functions have shown great success replacing multi-layer perceptrons (MLPs). This work, proposes the use of KANs in a MetricGAN framework and analyse their performance replacing different types of network layers. The best-performing proposed MetricGAN+KAN model uses 79.85% fewer parameters and achieves 13.1% higher SE performance (measured by PESQ) on the Voicebank-DEMAND dataset, compared to the MetricGAN+ baseline.

Index Terms—Speech enhancement, quality metrics, Kolmogorov-Arnold network (KAN), Generative adversarial network (GAN), MetricGAN

I. INTRODUCTION

Single-channel speech enhancement (SE) has been a popular research field for some decades [1], focusing on improving the quality [2]–[5] or intelligibility [6], [7] of speech signals in noisy, reverberant environments [8]. Machine learning (ML)-based approaches have led to significant performance gains in recent years, and become the first choice of modelling for SE [9]–[12]. Generative adversarial networks (GANs) [13] and conditional GANs (CGANs) [14] which consist of two sub-models, a generator and a discriminator, have proven to be effective in SE. The MetricGAN [15] approach and its extensions [16]–[24] have achieved the-state-of-the-art results on the Voicebank-DEMAND [25] dataset. However, only limited research exists for optimising the model structure of the MetricGAN framework even though this was already suggested by the authors of [16]. Furthermore, it is time-consuming to train the model with replay buffer, which is necessary for addressing catastrophic forgetting [26].

Recently, KANs [27] have been proposed, integrating learnable activation functions parameterised by B-spline curves into neurons. Authors of KANs have also mentioned that KANs can overcome catastrophic forgetting [27]. Hence, this work analyses the use KANs in the MetricGAN+ framework. The proposed KAN-based SE model is therefore denoted as MetricGAN+KAN, and this work aims at validating some advantages of KANs in a MetricGAN setting. We analyse different model structures, i.e. positions to replace model layer with KAN-based layers and compare performance as well as model parameters to MetricGAN+ on the Voicebank-DEMAND task.

II. REVIEW OF THE METRICGAN+ FRAMEWORK

MetricGAN+ [16] is a spectro-temporal masking-based SE approach. For this, the noisy input signal is first converted to a magnitude spectrogram $X_{f,\tau}$ and a phase spectrogram $\gamma_{f,\tau}$ by the short-time Fourier transform (STFT), where f is the frequency index and τ is the frame index. For the enhancement process, a spectral mask $M_{f,\tau}$ is computed and multiplied with the magnitude spectrogram

of the noisy signal to obtain an estimate of the clean magnitude spectrogram

$$\hat{S}_{f,\tau} = M_{f,\tau} X_{f,\tau}. \quad (1)$$

Then, an estimate of the clean signal is re-synthesised using $\hat{S}_{f,\tau}$ and $\gamma_{f,\tau}$, i.e. by applying the inverse STFT to $\hat{S}_{f,\tau} e^{j\gamma_{f,\tau}}$.

MetricGAN+ [16] consists of two neural networks (NNs), a generator \mathcal{G} aiming to estimate the mask $M_{f,\tau}$ and a discriminator \mathcal{D} assessing the quality of the masking-based SE by metric prediction.

The generator \mathcal{G} takes a noisy spectrogram $X_{f,\tau}$ as the input and outputs the mask $M_{f,\tau}$. Figure 1 visualises the generator \mathcal{G}_0 of the MetricGAN+ baseline [16], which can be split into a part containing recursive layers, i.e. a bidirectional long short-term memory (LSTM) [28] and a part containing non-recursive layers with a leaky rectified linear unit (ReLU) activation function. A learnable sigmoid outputs the mask $M_{f,\tau}$.

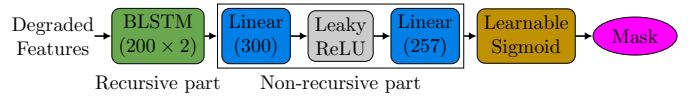


Figure 1: Generator model \mathcal{G}_0 (MetricGAN+ baseline)

Figure 2 visualises the discriminator model structure of MetricGAN+ \mathcal{D}_0 . After a batch normalisation (BN) layer, it can be split into a convolutional part and a non-convolutional part. Subscripts $_0$ in Figures 1 and 2 indicate the baseline [16] model in contrast to model variants introduced later in this work.

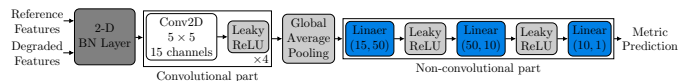


Figure 2: Discriminator model \mathcal{D}_0 (MetricGAN+ baseline)

The discriminator \mathcal{D} predicts a metric score $Q'(\cdot)$ normalised between 0 and 1 (often a normalised version of the perceptual evaluation of speech quality (PESQ) metric) given the noisy (or enhanced) magnitude spectrogram and the corresponding clean spectrogram $S_{f,\tau}$. The discriminator is also a differentiable surrogate function which imitates non-differentiable intrusive measures of audio quality such as PESQ [29], subjective mean opinion score (MOS) [30], or DNSMOS [31] since such perceptually motivated metrics often correlate better with human perception [2], [32] than traditional loss functions such as simple mean squared error (MSE) and hence can address discriminator evaluation mismatch (DEM) [15].

In terms of training, the loss function for the discriminator is given by

$$L_{\mathcal{D}} = \mathbb{E}_{\mathbf{X}, \mathbf{S}} \left[(\mathcal{D}(\mathbf{S}, \mathbf{S}) - Q'(\mathbf{S}, \mathbf{S}))^2 + (\mathcal{D}(\mathcal{G}(\mathbf{X}), \mathbf{S}) - Q'(\mathcal{G}(\mathbf{X}), \mathbf{S}))^2 + (\mathcal{D}(\mathbf{X}, \mathbf{S}) - Q'(\mathbf{X}, \mathbf{S}))^2 \right], \quad (2)$$

where \mathbf{X} is the magnitude spectrogram matrix of the noisy signal containing $X_{f,\tau} \forall f, \tau$ and \mathbf{S} is the respective clean spectrogram matrix. The loss function for the generator is given by

$$L_G = \mathbb{E}_{\mathbf{X}} [(\mathcal{D}(\mathcal{G}(\mathbf{X}), \mathbf{S}) - w)^2], \quad (3)$$

with w being the desired metric score that the discriminator assigned to the enhanced speech, which is set to 1 in [16] maximising the enhancement or varied in [22]. In each epoch, MetricGAN+ is trained using the following procedure:

- 1) Train the generator using back-propagation (BP) [33].
- 2) Store current enhanced signals and the corresponding scores into the so-called replay buffer.
- 3) Train the discriminator using clean signals, current enhanced waves and noisy waves.
- 4) Repeat 3), but use a part of the enhanced waves in the replay buffer, which is controlled by the hyper-parameter `history_portion` [16], [22].

As mentioned above, the replay buffer is used to address catastrophic forgetting in the discriminator. It can greatly improve the performance of the discriminator, and subsequently improve the quality of signals enhanced by the generator. However, training with the replay buffer increases training time. Authors of MetricGAN+ [16] already mention that the structure of the discriminator can be improved which will be analysed in this work by using KANs in the recursive and non-recursive parts of the generator as well as in the convolutional and non-convolutional part of the discriminator.

III. REVIEW OF KOLMOGOROV-ARNOLD NETWORKS (KANs)

KANs [27] are a recently proposed type of NN architecture having gained considerable attention on GitHub¹. The novelty of KANs is that the activation function is placed within the neuron, and is learnable. It is inspired by the Kolmogorov-Arnold representation theorem [34]

$$f(\mathbf{x}) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right), \quad (4)$$

where $f: [0, 1]^n \rightarrow \mathbb{R}$ is smooth, $\phi_{q,p}: [0, 1] \rightarrow \mathbb{R}$, and $\Phi_q: \mathbb{R} \rightarrow \mathbb{R}$. Based on (4), KANs replace the weight matrix in a traditional NN layer with a matrix of functions, denoted by $\Phi = \{\phi_{q,p}\}$ where $p = 1, 2, \dots, n_{\text{in}}$ and $q = 1, 2, \dots, n_{\text{out}}$. $\phi(\cdot)$ is a learnable activation function formulated as the scaled sum of a base activation function $b(x)$ and a learnable curve $g(x)$,

$$\phi(x) = w_1 b(x) + w_2 g(x), \quad (5)$$

where w_1 and w_2 are scaling factors which can be learnable. It is noteworthy that in the original paper [27], only one scaling factor is used for $\phi(\cdot)$. Liu *et al.* used sigmoid linear units (SiLUs) [35]

$$b(x) = \frac{x}{1 + e^{-x}} \quad (6)$$

as the base activation functions for KANs. Several implementations have been reviewed in [36]. For the learnable curves, B-splines, which

require basis functions and controlling points were proposed initially with B-spline basis functions of order k defined as [37], [38]

$$B_{i,0}(x) = \begin{cases} 1 & \text{if } x_i \leq x \leq x_{i+1}, \\ 0 & \text{otherwise,} \end{cases} \quad (7)$$

$$B_{i,k}(x) = \frac{x - x_i}{x_{i+k} - x_i} B_{i,k-1}(x) + \frac{x_{i+k+1} - x}{x_{i+k+1} - x_{i+k}} B_{i+1,k-1}(x), \quad (8)$$

with x_i for $i = -k, -k+1, \dots, G+k$ being the predefined boundary, and G the grid size. The spline curve is given by

$$g(x) = \sum_{i=0}^{G+k-1} c_i B_{i,k}(x), \quad (9)$$

where c_i for $i = 0, 1, \dots, G+k-1$ is the trainable controlling point, and $g(x)$ is defined on $[x_0, x_G]$.

KANs have also been integrated into convolutional neural networks (CNNs) and recurrent neural networks (RNNs), called convolutional KANs (CKANs) [39] and recurring KANs (RKANs) [40], respectively. In CKANs, the kernel becomes a group of learnable activation functions. The result of a 2-D convolution (with stride 1) is given by

$$a_{x,y}^{(l)} = \sum_{i=1}^m \sum_{j=1}^n \phi_{i,j}^{(l)} \left(a_{x+i,y+j}^{(l-1)} \right), \quad (10)$$

where $a_{x,y}^{(l)}$ is the feature map at layer l . In RKANs, the prediction at time t is given by

$$\hat{\mathbf{y}}_t = \Phi \mathbf{h}_t, \quad (11)$$

where \mathbf{h}_t is the hidden state at time t .

As mentioned in [27], advantages of KANs are that KANs can overcome catastrophic forgetting, because the update on c_i only changes part of the spline curve, and that KANs have a better scaling law than traditional NNs, *i.e.*, using fewer parameters to achieve similar (or higher) performance comparing to traditional NNs. Bodner *et al.* [39] have shown that the two advantages also work for CKANs.

IV. EXPERIMENTS

A. Dataset

The Voicebank-DEMAND dataset [25] is used for the following experiments, created from the Voicebank dataset [41] which contains 300-hours clean speech audio spoken by approximately 500 healthy speakers from the UK mixed with various types of noises recorded indoor and outdoor from the DEMAND dataset [42] (cafeteria, car interior, a kitchen, meeting, metro station, restaurant, train station and heavy traffic) and two others (babble noise and speech-shaped noise).

The training set of Voicebank-DEMAND consists of 11572 noisy speech signals at 4 signal-to-noise ratios (SNRs) of 0, 5, 10, and 15 dB paired with the respective clean speech reference signals from 28 different speakers (14 male, 14 female), with English or Scottish accents. The testset contains 824 utterances, mixed at SNRs of 2.5, 7.5, 12.5 and 17.5 dB, with five different noises which do not appear in the training set (bus, cafe, office, public square and living room) and contains speech from two (one male, one female) speakers who do not appear in the training set.

B. Implementation

Models are trained using the SpeechBrain framework [43]. For the implementation of KANs, `efficient-kan` [44] is used, and `torch-conv-kan` [36] for CKANs (with parametric ReLU [45] activation at the output). For RKANs, a gated recurrent unit (GRU)

¹<https://github.com/KindXiaoming/pykan>.

[46] version of RKANs, namely GRU-KAN, is implemented. GRU-KAN uses the same formulae as GRU, but uses (11) for computing the prediction.

C. Experiment Setup

In the future work section of the original paper of MetricGAN+ [16], Fu *et al.* have indicated that the discriminator can be improved to achieve better performance, and that the use of replay buffer is time-consuming. KANs proposed by Liu *et al.* [27] can be a possible improvement and mitigate catastrophic forgetting. Thus, it is worth implementing an improved version of MetricGAN+ basing on KANs, namely MetricGAN+KAN, and testing some advantages of KANs.

Therefore, the first aim of the experiment is to compare MetricGAN+KAN with MetricGAN+ in terms of SE performance. Several modifications on the discriminator and the generator are experimented. The secondary aim is to validate two advantages of KANs. One is that KANs have a better scaling law. The other is that KANs can mitigate catastrophic forgetting.

In the presented experiment, there are six generators (\mathcal{G}_0 to \mathcal{G}_5) and five discriminators (\mathcal{D}_0 to \mathcal{D}_4). Generator \mathcal{G}_0 (see Figure 1) and discriminator \mathcal{D}_0 (see Figure 2) are used in MetricGAN+, which is the baseline, and other discriminators and generators are modified basing on them respectively (see Tables Ia and Ib). Tables Id and Ie show the number of parameters of each model. The naming of MetricGAN+KAN follows the form $\text{mgk-g}__-\text{d}__ < (\text{NHP}) >$, where mgk is the abbreviation of MetricGAN+KAN, $\text{g}__$ specifies the generator, $\text{d}__$ specifies the discriminator, and (NHP) means training without replay buffer, *i.e.*, history_portion is set to 0. All the models are trained for 400 epochs. For KAN hyper-parameters, the range of splines is $[-1, 1]$, the grid size is 5, and the spline order is set to 3. Other hyper-parameters used in MetricGAN+KAN are the same as MetricGAN+.

V. RESULTS

Results in Table I show that KANs have a better scaling law. Discriminators \mathcal{D}_1 and \mathcal{D}_2 which integrated with KANs improved the performance of SE marginally with slightly fewer parameters, comparing to MetricGAN+. In terms of discriminators integrated with CKANs, discriminators \mathcal{D}_3 and \mathcal{D}_4 showed great improvement with much more parameters, and discriminator \mathcal{D}_5 showed slight improvement with much fewer parameters. For generators, generators \mathcal{G}_1 and \mathcal{G}_2 which integrated with KANs slightly degraded the performance even with much more parameters. Generators \mathcal{G}_3 uses smaller size of hidden states and smaller number of layers with significantly fewer parameters, and still reached similar performance. Generator \mathcal{G}_5 showed that GRU worked better than LSTM, and generator \mathcal{G}_4 showed that KANs further improved the performance of generator \mathcal{G}_5 with much fewer parameters. Generator \mathcal{D}_6 showed that GRU-KAN also worked well comparing to MetricGAN+. In summary, the integration of KANs has marginal improvement, while that of CKANs has significant improvement. Both KANs and CKANs can have similar performance as traditional NNs using fewer parameters, and CKANs can even outperform CNNs.

Among the tested variants of MetricGAN+KAN, $\text{mgk-g}_4\text{-d}_4$ is chosen as the best model. It uses 79.85% fewer parameters (85.35% fewer for the generator, 468.9% more for the discriminator) and achieves 13.15% higher PESQ scores comparing to MetricGAN+.

Figure 3 shows the comparison among the clean, the noisy, MetricGAN+ enhanced and $\text{mgk-g}_4\text{-d}_4$ enhanced magnitude spectrograms.

Results also show that MetricGAN+KAN still suffers from catastrophic forgetting, and the use of replay buffer is still necessary.

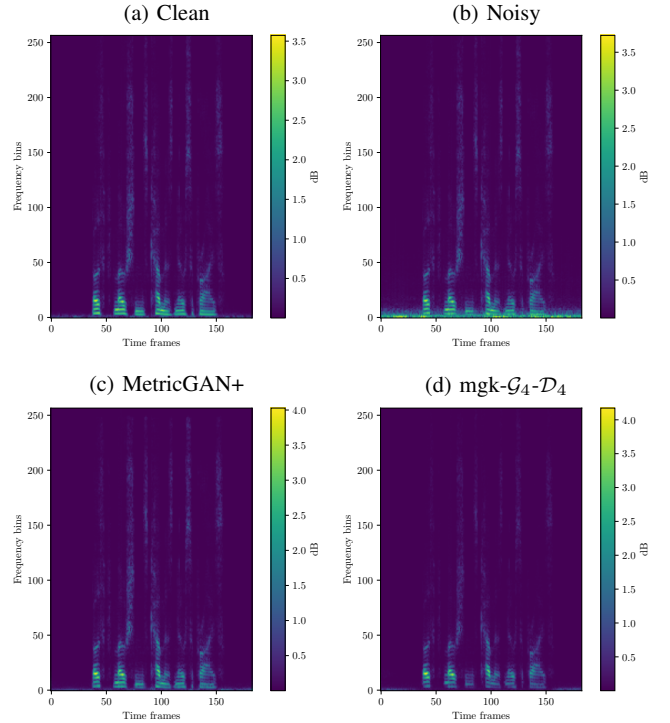


Figure 3: Comparison of magnitude spectrograms

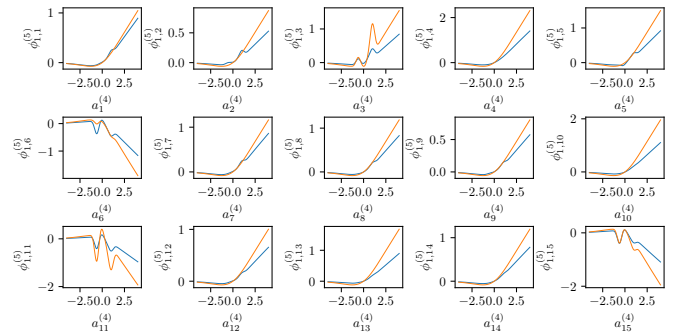


Figure 4: Changes of Eq. (5) in $\text{mgk-g}_0\text{-d}_4$, taken from epoch 100 (in blue) and epoch 400 (in red). $\phi_{i,j}^{(\ell)}$ represents the learnable activation function of input dimension i , output dimension j at layer ℓ , with the corresponding input $a_j^{(\ell-1)}$.

Models trained without replay buffer are degraded significantly. A possible reason is that the locality of splines only exists in Eq. (9). In Eq. (5), the spline curve is scaled, which may reduce the effectiveness of the locality. In other words, in some neurons, the base activation function has much influence than the spline curve (see Fig. 4). However, this reason needs further experiments to show whether it holds or not.

VI. FUTURE WORK

In this work, MetricGAN+KAN is tested on Voicebank-DEMAND. The model needs further experiments on a different dataset. Besides, the reason for the existence of catastrophic forgetting in MetricGAN+KAN (or KANs) needs further investigations. In terms of the further improvement of MetricGAN+KAN, Close *et al.* [22] have shown that introducing a de-generator to MetricGAN+ can improve

(a) Results of only changing generators. Recursive part specifies the type of RNN and non-recursive part specifies the type of feed-forward layers. The number in the parenthesis specifies the hidden size (or the number of neurons), and \times specifies the number of layers (default value is 1).

	Recursive part	Non-recursive part	PESQ	CSIG	CBAK	COVL
Noisy			1.97	3.35	2.44	2.63
MetricGAN+	BLSTM (200 \times 2)	Linear (300, 257)	2.89	3.78	2.92	3.31
mgk- \mathcal{G}_1 - \mathcal{D}_0	BLSTM (200 \times 2)	KAN (80), Linear (257)	2.85	3.73	2.90	3.26
mgk- \mathcal{G}_1 - \mathcal{D}_0 (NHP)	BLSTM (200 \times 2)	KAN (80), Linear (257)	2.68	3.93	2.73	3.30
mgk- \mathcal{G}_2 - \mathcal{D}_0	BLSTM (200 \times 2)	KAN (257)	2.82	3.80	2.93	3.28
mgk- \mathcal{G}_3 - \mathcal{D}_0	BLSTM (40)	KAN (257)	2.85	3.69	2.83	3.23
mgk- \mathcal{G}_4 - \mathcal{D}_0	BGRU (40)	KAN (257)	2.94	3.82	2.88	3.35
mgk- \mathcal{G}_5 - \mathcal{D}_0	BGRU (100)	Linear (300, 257)	2.88	3.94	2.80	3.38
mgk- \mathcal{D}_6 - \mathcal{D}_0	BGRU-KAN (40 \times 2)	KAN (257)	2.93	3.84	2.93	3.36

(b) Results of only changing discriminators. Convolutional part specifies the type of CNN and non-convolutional part specifies the type of feed-forward layers. The number in the parenthesis specifies the number of output channels. The shape of all the convolutional kernels is 5×5 .

	Convolutional part	Non-convolutional part	PESQ	CSIG	CBAK	COVL
Noisy			1.97	3.35	2.44	2.63
MetricGAN+	Conv2d (15 \times 4)	Linear (50, 10, 1)	2.89	3.78	2.92	3.31
mgk- \mathcal{G}_0 - \mathcal{D}_1	Conv2d (15 \times 4)	Linear (50), KAN (1)	2.94	4.00	2.91	3.45
mgk- \mathcal{G}_0 - \mathcal{D}_2	Conv2d (15 \times 4)	KAN (1)	2.93	3.97	3.01	3.44
mgk- \mathcal{G}_0 - \mathcal{D}_3	CKAN2d (15 \times 2)	KAN (1)	3.02	4.03	3.02	3.50
mgk- \mathcal{G}_0 - \mathcal{D}_4	CKAN2d (15 \times 3)	KAN (1)	3.30	4.02	3.04	3.63
mgk- \mathcal{G}_0 - \mathcal{D}_4 (NHP)	CKAN2d (15 \times 3)	KAN (1)	2.72	3.96	2.75	3.32
mgk- \mathcal{G}_0 - \mathcal{D}_5	CKAN2d (20)	KAN (1)	2.96	4.15	3.19	3.55

(c) Results of unhandled noisy speeches, the baseline, and some combinations of discriminators and generators

	PESQ	CSIG	CBAK	COVL
Noisy	1.97	3.35	2.44	2.63
MetricGAN+	2.89	3.78	2.92	3.31
mgk- \mathcal{G}_4 - \mathcal{D}_3	3.00	3.98	2.95	3.46
mgk- \mathcal{G}_4 - \mathcal{D}_3 (NHP)	2.66	3.85	2.88	3.24
mgk- \mathcal{G}_4 - \mathcal{D}_4	3.27	3.97	2.97	3.59
mgk- \mathcal{G}_5 - \mathcal{D}_3	3.07	4.10	3.00	3.57
mgk- \mathcal{G}_5 - \mathcal{D}_4	3.08	4.08	2.99	3.56
mgk- \mathcal{D}_6 - \mathcal{D}_3	2.99	4.03	3.04	3.49
mgk- \mathcal{D}_6 - \mathcal{D}_4	3.12	3.90	2.95	3.48

(d) Generator parameter count

\mathcal{G}_0	1 895 514
\mathcal{G}_1	2 038 674
\mathcal{G}_2	2 725 857
\mathcal{G}_3	301 537
\mathcal{G}_4	277 617
\mathcal{G}_5	353 314
\mathcal{D}_6	361 267

(e) Discriminator parameter count

\mathcal{D}_0	19 010
\mathcal{D}_1	18 989
\mathcal{D}_2	17 839
\mathcal{D}_3	57 531
\mathcal{D}_4	108 157
\mathcal{D}_5	9 205

Table I: Experimental results in terms of PESQ and composite metrics on test set, and parameter count of generators and discriminators

SE performance and become more robust to unseen noises, which may also work in MetricGAN+KAN. Additionally, Drokin’s study [36] have presented several implementations of the learnable curve, which can be a possible direction for experiments.

VII. CONCLUSION

This work analysed the use of KANs for a MetricGAN+ SE system. The integration of KANs, CKANs and RKANs can improve the SE performance of MetricGAN+. The best model, mgk- \mathcal{G}_4 - \mathcal{D}_4 , achieves 13.15% higher PESQ scores with 79.85% fewer parameters comparing to MetricGAN+, which shows a better scaling law. Experimental results also show that the use of KANs in MetricGAN+KAN cannot mitigate catastrophic forgetting, and the replay buffer is still necessary.

REFERENCES

- [1] P. Vary and R. Martin, *Digital Speech Transmission: Enhancement, Coding and Error Concealment*. John Wiley & Sons, Ltd, 2006.
- [2] S. Goetze, A. Warzybok, I. Kodrasi, J. Jungmann, B. Cauchi, J. RENNIES, E. Habets, A. Mertins, T. Gerkmann, S. Doclo, and B. Kollmeier, “A Study on Speech Quality and Speech Intelligibility Measures for Quality Assessment of Single-Channel Dereverberation Algorithms,” in *Proc. IWAENC’14*, Sep. 2014.
- [3] S.-W. Fu, C.-F. Liao, and Y. Tsao, “Learning with learned loss function: Speech enhancement with quality-net to improve perceptual evaluation of speech quality,” *IEEE Signal Processing Letters*, vol. 27, 2020.
- [4] B. Cauchi, K. Siedenburg, J. F. Santos, T. H. Falk, S. Doclo, and S. Goetze, “Non-Intrusive Speech Quality Prediction Using Modulation Energies and LSTM-Network,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, pp. 1151–1163, July 2019.
- [5] G. Close, W. Ravenscroft, T. Hain, and S. Goetze, “Perceive and predict: self-supervised speech representation based loss functions for speech enhancement,” in *Proc. ICASSP 2023*, 2023.
- [6] C. Völker, A. Warzybok, and S. Ernst, “Comparing Binaural Pre-processing Strategies III: Speech Intelligibility of Normal-Hearing and Hearing-Impaired Listeners,” *Trends in Hearing*, vol. 19, 2015.
- [7] R. Mogridge, G. Close, R. Sutherland, T. Hain, J. Barker, S. Goetze, and A. Ragni, “Non-Intrusive Speech Intelligibility Prediction for Hearing-Impaired Users using Intermediate ASR Features and Human Memory Models,” in *Proc. ICASSP’24*, (Seoul, South Korea), Apr. 2024.
- [8] F. Xiong, B. Meyer, N. Moritz, R. Rehr, J. Anemüller, T. Gerkmann, S. Doclo, and S. Goetze, “Front-end technologies for robust ASR in reverberant environments - spectral enhancement-based dereverberation and auditory modulation filterbank features,” *EURASIP Journal on Advances in Signal Processing*, vol. 2015, no. 1, 2015.
- [9] J. Barker, R. Marxer, E. Vincent, and S. Watanabe, “The third CHiME speech separation and recognition challenge: dataset, task and baselines,” in *Proc. ASRU*, 2015.
- [10] W. Ravenscroft, S. Goetze, and T. Hain, “Att-TasNet: Attending to Encodings in Time-Domain Audio Speech Separation of Noisy, Reverberant Speech Mixtures,” *Frontiers in Signal Processing*, vol. 2, 2022.
- [11] M. Tammen and S. Doclo, “Deep multi-frame MVDR filtering for single-microphone speech enhancement,” in *Proc. ICASSP*, 2021.
- [12] N. Moritz, K. Adiloğlu, J. Anemüller, S. Goetze, and B. Kollmeier,

- “Multi-channel speech enhancement and amplitude modulation analysis for noise robust automatic speech recognition,” *Computer Speech & Language*, vol. 46, 2017.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in Neural Information Processing Systems* (Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, eds.), vol. 27, Curran Associates, Inc., 2014.
- [14] M. Mirza and S. Osindero, “Conditional generative adversarial nets,” *CoRR*, vol. abs/1411.1784, 2014.
- [15] S.-W. Fu, C.-F. Liao, Y. Tsao, and S.-D. Lin, “MetricGAN: Generative Adversarial Networks based Black-box Metric Scores Optimization for Speech Enhancement,” in *Proc. 36th Int. Conf. on Machine Learning*, Jun 2019.
- [16] S. Fu, C. Yu, T. Hsieh, P. Plantinga, M. Ravanelli, X. Lu, and Y. Tsao, “MetricGAN+: An Improved Version of MetricGAN for Speech Enhancement,” *CoRR*, vol. abs/2104.03538, 2021.
- [17] J. Kong, J. Kim, and J. Bae, “Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis,” *Advances in neural information processing systems*, vol. 33, pp. 17022–17033, 2020.
- [18] C. Donahue, B. Li, and R. Prabhavalkar, “Exploring speech enhancement with generative adversarial networks for robust speech recognition,” in *ICASSP*, pp. 5024–5028, 2018.
- [19] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, “An experimental study on speech enhancement based on deep neural networks,” *IEEE Signal processing letters*, vol. 21, no. 1, pp. 65–68, 2013.
- [20] Y. Wang, A. Narayanan, and D. Wang, “On training targets for supervised speech separation,” *IEEE/ACM transactions on audio, speech, and language processing*, vol. 22, no. 12, pp. 1849–1858, 2014.
- [21] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, “Speech enhancement based on deep denoising autoencoder,” in *Interspeech*, vol. 2013, 2013.
- [22] G. Close, T. Hain, and S. Goetze, “MetricGAN+/-: Increasing Robustness of Noise Reduction on Unseen Data,” in *EUSIPCO 2022*, 2022.
- [23] G. Close, T. Hain, and S. Goetze, “PAMGAN+/-: Improving Phase-Aware Speech Enhancement Performance via Expanded Discriminator Training,” in *AES 154th Conv.*, May 2023.
- [24] G. Close, W. Ravenscroft, T. Hain, and S. Goetze, “CMGAN+/-: The University of Sheffield CHiME-7 UDASE Challenge Speech Enhancement System,” in *Proc. 7th Int. Workshop on Speech Processing in Everyday Environments (CHiME 2023)*, Aug. 2023.
- [25] C. Valentini-Botinhao, X. Wang, S. Takaki, and J. Yamagishi, “Investigating RNN-based speech enhancement methods for noise-robust Text-to-Speech,” in *SSW*, pp. 146–152, 2016.
- [26] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, “An empirical investigation of catastrophic forgetting in gradient-based neural networks,” *arXiv preprint arXiv:1312.6211*, 2013.
- [27] Z. Liu, Y. Wang, S. Vaidya, F. Ruehle, J. Halverson, M. Soljačić, T. Y. Hou, and M. Tegmark, “Kan: Kolmogorov-Arnold networks,” *arXiv preprint arXiv:2404.19756*, 2024.
- [28] M. Schuster and K. Paliwal, “Bidirectional recurrent neural networks,” *IEEE Transactions on Signal Processing*, vol. 45, 12 1997.
- [29] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, “Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs,” in *2001 IEEE International Conference on Acoustics, Speech, and Signal Processing*, vol. 2, pp. 749–752, 2001.
- [30] R. C. Streijl, S. Winkler, and D. S. Hands, “Mean opinion score (MOS) revisited: methods and applications, limitations and alternatives,” *Multimedia Systems*, vol. 22, no. 2, pp. 213–227, 2016.
- [31] C. K. A. Reddy, V. Gopal, and R. Cutler, “Dnsmos: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6493–6497, 2021.
- [32] A. Avila, B. Cauchi, S. Goetze, S. Doclo, and T. Falk, “Performance Comparison of Intrusive and Non-Intrusive Instrumental Quality Measures for Microphone-Array Processed Speech,” in *Proc. International Workshop on Acoustic Signal Enhancement IWAENC 2016*, 2016.
- [33] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*, vol. 323, no. 6088, 1986.
- [34] A. N. Kolmogorov, “On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition,” in *Doklady Akademii Nauk*, vol. 114, pp. 953–956, Russian Academy of Sciences, 1957.
- [35] S. Elfving, E. Uchibe, and K. Doya, “Sigmoid-weighted linear units for neural network function approximation in reinforcement learning,” *CoRR*, vol. abs/1702.03118, 2017.
- [36] I. Drokin, “Kolmogorov-Arnold convolutions: Design principles and empirical studies,” *arXiv preprint arXiv:2407.01092*, 2024.
- [37] W. J. Gordon and R. F. Riesenfeld, “B-spline curves and surfaces,” in *Computer Aided Geometric Design*, pp. 95–126, Academic Press, 1974.
- [38] C. Boor, “Subroutine package for calculating with b-splines,” 1971.
- [39] A. D. Bodner, A. S. Tepsich, J. N. Spolski, and S. Pourteau, “Convolutional Kolmogorov-Arnold networks,” 2024.
- [40] R. Genet and H. Inzirillo, “Tkan: Temporal Kolmogorov-Arnold networks,” *arXiv preprint arXiv:2405.07344*, 2024.
- [41] C. Veaux, J. Yamagishi, and S. King, “The voice bank corpus: Design, collection and data analysis of a large regional accent speech database,” in *Int. conf. oriental COCOSDA, jointly with 2013 Conf. on Asian spoken language research and evaluation (O-COCOSDA/CASLRE)*, 2013.
- [42] J. Thiemann, N. Ito, and E. Vincent, “The Diverse Environments Multi-channel Acoustic Noise Database (DEMAND): A database of multichannel environmental noise recordings,” *Proc. of Meetings on Acoustics*, vol. 19, no. 1, 2013.
- [43] M. Ravanelli, T. Parcollet, P. Plantinga, A. Rouhe, S. Cornell, L. Lugosch, C. Subakan, N. Dawalatabad, A. Heba, J. Zhong, J.-C. Chou, S.-L. Yeh, S.-W. Fu, C.-F. Liao, E. Rastorgueva, F. Grondin, W. Aris, H. Na, Y. Gao, R. D. Mori, and Y. Bengio, “SpeechBrain: A general-purpose speech toolkit,” 2021. *arXiv:2106.04624*.
- [44] Blealtan and A. Dash, “efficient-kan.” <https://github.com/Blealtan/efficient-kan>, 2024.
- [45] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” *CoRR*, vol. abs/1502.01852, 2015.
- [46] K. Cho, B. van Merriënboer, Ç. Gülçehre, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” *CoRR*, vol. abs/1406.1078, 2014.