



# Adaptive Reverberation Absorption Using Non-Stationary Masking Components Detection for Intelligibility Improvement

Guilherme Zucatelli , *Student Member, IEEE*, and Rosângela Coelho , *Senior Member, IEEE*

**Abstract**—This letter proposes a new time domain absorption approach designed to reduce masking components of speech signals under noisy-reverberant conditions. In this method, the non-stationarity of corrupted signal segments is used to detect masking distortions based on a defined threshold. The non-stationarity is objectively measured and is also adopted to determine the absorption procedure. Additionally, no prior knowledge of speech statistics or room information is required for this technique. Two intelligibility measures (ESII and ASII<sub>ST</sub>) are used for objective evaluation. The results show that the proposed scheme leads to a higher intelligibility improvement when compared to competing methods. A perceptual listening test is further considered and corroborates these results. Furthermore, the updated version of the SRMR quality measure (SRMR<sub>norm</sub>) demonstrates that the proposed technique also attains quality improvement.

**Index Terms**—Reverberation, absorption, non-stationarity, intelligibility.

## I. INTRODUCTION

SPEECH communication commonly takes place in enclosed and urban environments such as concert halls, kitchens and offices. Along with the direct acoustic signal propagation between source and listener locations, the sound reverberates due to reflection in walls and surfaces. While the early reflections (ER) can improve speech intelligibility, late reverberation (LR) may cause quality and intelligibility reduction [1]–[4].

Room impulse response (RIR) typically describes the sound propagation and is generally defined by the reverberation time ( $T_{60}$ ) and the direct-to-reverberant ratio (DRR). Speech signals can also be degraded by background acoustic noises (Babble, Chainsaw and Cafeteria) present in the urban space. Such effects are non-stationary masking components and represent a major drawback to speech intelligibility improvement.

In the literature, speech enhancement solutions were designed to cope with background non-stationary noises [5]–[8] attaining interesting results for quality and intelligibility. However, room

reverberation is not considered by these techniques. Adaptive time-domain pre-processing methods were proposed to improve speech intelligibility by mitigating the effect of masking distortions on non-stationary frames. The Steady State Suppression (SSS) [9] solution considers the importance of transient regions of speech for intelligibility and suppresses steady-state frames to reduce overlap masking effects. A more recent approach, the Adaptive Gain Control (AGC) [10] method uses prior knowledge of speech statistics and the RIR information to adaptively improve or reduce the energy of speech frames. Both methods operate in a single channel prior to speech signal presentation in a room.

This letter proposes a new time-domain method denominated Adaptive Reverberation Absorption with Non-Stationary Detection (ARA<sub>NSD</sub>). Different from SSS and AGC techniques, the main idea of this proposal is to absorb masking components at the listener's position that are detrimental to speech intelligibility and its natural non-stationary property. The approach works similar to a physical element, changing the low absorption characteristic of materials that compose a room. One major advantage of ARA<sub>NSD</sub> is that it adaptively mitigate such distortions, leading to speech intelligibility improvement and restoration of the non-stationarity behavior with no prior knowledge of the RIR or speech statistics. The Index of Non-Stationarity (INS) [11] is selected as an objective measure for the detection of masking components. A non-stationarity threshold is defined for the proposed frame-by-frame absorption procedure.

Extensive experiments are conducted to objectively evaluate the ARA<sub>NSD</sub> for speech intelligibility improvement. The noisy-reverberant scenario is composed of two real reverberant rooms and four background non-stationary acoustic noises with five different SNR values. The ESII [12] and ASII<sub>ST</sub> [13] measures are adopted for the intelligibility prediction. These measures are explicitly designed to deal with the non-stationarity of speech and its distortions. The SRMR<sub>norm</sub> [14] measure is further considered as it is primarily used for signals under reverberation. A subjective listening test is also performed and results show that the proposed method outperforms competing techniques in terms of speech intelligibility.

## II. REVERBERATION AND NON-STATIONARITY

The reverberation effect is usually defined as a linear filtering process such that, given a RIR  $h(n)$ , the reverberated signal can be obtained by convolution. In real environments, acoustic noises are also a common distortion, which means that the resultant noisy-reverberant speech signal  $s(n)$  can be obtained

Manuscript received September 20, 2019; revised October 18, 2019; accepted October 20, 2019. Date of publication October 31, 2019; date of current version January 22, 2020. This work was supported in part by the National Council for Scientific and Technological Development (CNPq) under Grant 307866/2015, in part by the Fundação de Amparo Pesquisa do Estado do Rio de Janeiro (FAPERJ) under Grant 203075/2016, and in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior—Brasil under Grant Code 001. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Maximo Cobos. (Corresponding author: Rosângela Coelho.)

The authors are with the Laboratory of Acoustic Signal Processing (lasp.ime.eb.br), Military Institute of Engineering (IME), Rio de Janeiro, Brazil (e-mail: coelho@ime.eb.br).

Digital Object Identifier 10.1109/LSP.2019.2950618

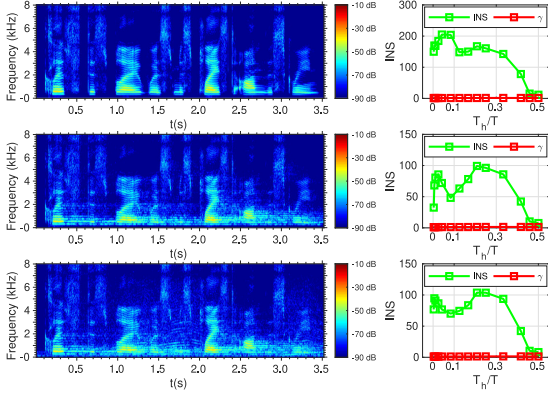


Fig. 1. Spectrogram and related INS for direct signal (top), reverberated signal with  $T_{60} = 4.9$  s and  $SRR = 7.1$  dB (middle), and reverberated signal with Chainsaw noise at  $-3$  dB (bottom).

by  $s(n) = x(n) * h(n) + w(n)$ , where  $x(n)$  is the clean speech signal and  $w(n)$  is the background noise.

The Index of Non-Stationarity (INS) [11] is here defined to objectively examine the non-stationarity of speech signals under noisy-reverberant environments. This measure compares the target signal with stationarity references called surrogates for different time scales  $T_h/T$ , where  $T_h$  is the short-time spectral analysis length and  $T$  is the total signal duration. For each length  $T_h$ , a threshold  $\gamma$  is defined to keep the stationarity assumption considering a 95% confidence degree as

$$INS \begin{cases} \leq \gamma, & \text{signal is stationary} \\ > \gamma, & \text{signal is non-stationary.} \end{cases} \quad (1)$$

Fig. 1 illustrates the spectrograms and INS values obtained for a direct speech signal and its corresponding reverberated version in the Aula Carolina<sup>1</sup> room with  $T_{60} = 4.9$  s in two conditions: without and with a background Chainsaw noise at  $-3$  dB. Note that reverberation and acoustic noise significantly change the temporal and spectral structure of speech signal. These masking effects can engender intelligibility reduction [1]–[3]. Furthermore, the non-stationary behavior of the natural speech signal is considerably attenuated, varying its maximum INS value from 200 to around 100. The background Chainsaw noise increases the INS value in small scales, which means that short-time segments become more distinct from the overall signal. Therefore, the INS is here proposed for detection and reduction of masking components on non-stationary speech regions.

### III. ADAPTIVE REVERBERATION ABSORPTION WITH NON-STATIONARY DETECTION

The ARA<sub>SND</sub> method is presented in this section. The technique is described in two main phases: reverberation detection and acoustic absorption.

#### A. Reverberation Detection

A reverberation group (RG), denoted as  $s_{RG}(m, n)$ , is here defined as the  $m$ -th segment composed of  $N = 8$  consecutive frames of the corrupted speech. This window duration is selected to enable a long-term temporal observation of the reverberation effect and detect noisy-reverberant masking components.

<sup>1</sup>RIR collected from the AIR database [15].

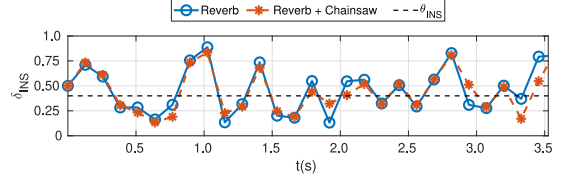


Fig. 2. Non-stationarity variation ( $\delta_{INS}$ ) for the reverberated and noisy-reverberated signals.

Successive RGs are obtained considering 50% overlap between signals.

For each  $s_{RG}(m, n)$ , the INS values are computed considering different scales of  $T_h/T$ . The INS values obtained for all scales are grouped into a vector  $\mathbf{v}_{INS}(m)$  which characterizes the non-stationary behavior of the  $m$ -th RG. Consecutive vectors are then used to compute a normalized variation of the non-stationary property as

$$\delta_{INS}(m) = \frac{\|\mathbf{v}_{INS}(m) - \mathbf{v}_{INS}(m-1)\|}{\|\mathbf{v}_{INS}(m)\| + \|\mathbf{v}_{INS}(m-1)\|}. \quad (2)$$

Fig. 2 shows the  $\delta_{INS}$  values obtained for the reverberated and noisy-reverberated speech signals of Fig. 1. Note that even with masking components, important speech regions, e.g. the ones near 0.2 s, 1.1 s, 1.4 s and 2.7 s (refer to Fig. 1 (top)), are still identified by the highest values of  $\delta_{INS}$  in both conditions. Moreover, masked regions closed to 0.7 s, 1.9 s and 3.0 s attain low  $\delta_{INS}$  values. This demonstrates that the proposed  $\delta_{INS}$  is an interesting detection approach for noise and reverberation masking components. The  $\theta_{INS}$  (black dashed line) in Fig. 2 illustrates a threshold of non-stationarity defined by the median value of  $\delta_{INS}$ . In this example,  $\theta_{INS}$  value is 0.4 indicating the difference of the speech and noisy-reverberant regions.

#### B. Acoustic Absorption

The proposed ARA<sub>NSD</sub> absorption approach is implemented on a frame-by-frame basis and is established depending on the value of  $\theta_{INS}$ . For each frame  $s_{frm}(l, n)$ , a INS vector  $\mathbf{v}_{frm}(l)$  is extracted similarly as in Section III-A. A short-time distance  $d(l) \in [0, 1]$  is then computed as in (2) and determine the  $l$ -th frame absorption.

Sigmoid functions are selected to assign each value of  $d(l)$  to a corresponding absorption  $A(m, l)$  because of their smoothness and monotonic property. The proposed adaptive absorption  $A(m, l)$  is therefore defined in every frame  $l$  by

$$A(m, l) = \begin{cases} F(l) \cdot \frac{L(m) - S}{1 + \exp(-k \cdot (d(l) - d_0))} + S, & \delta_{INS} \leq \theta_{INS}; \\ \frac{L'}{1 + \exp(-k' \cdot (d(l) - d'_0))}, & \delta_{INS} > \theta_{INS}, \end{cases} \quad (3)$$

where  $d_0$  and  $d'_0$  are the inflection points with corresponding growth rate of  $k$  and  $k'$ . The  $S$  stands for a minimum shift in order to avoid total absorption of signal frames. Moreover  $L(m)$  and  $L'$  are the maximum absorption values. As the noisy-reverberant masking effect is non-stationary by nature it is important to determine an adaptive upper bound absorption. Both  $L(m)$  and the  $F(l)$  factor are considered for this task. The first one is updated accounting the overlapped region as  $L(m) = p\delta_{INS} + (1 - p)L(m-1)$ , where  $p$  assigns the importance of the present RG signal. The second term is defined as the factor  $F(l) = d(l)^{1.2-d(l)}$  to guarantee that  $A(m, l) \approx L(m)$  only for  $d(l) \approx 1$ . As  $d(l)$  represents the short-term non-stationarity

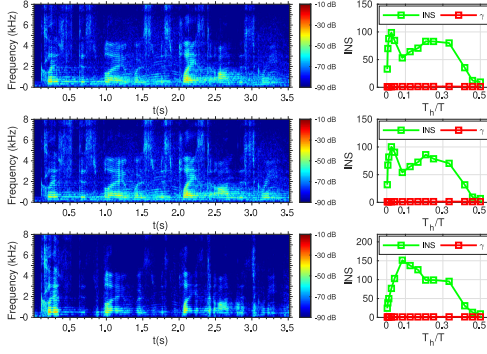


Fig. 3. Spectrogram and related INS for noisy-reverberant signal processed by SSS (top), AGC (middle) and  $ARA_{NSD}$  (bottom).

behavior, the absorption maintain a high value if  $d(l) \approx 1$  for it refers to an important speech region. The processed signal  $s'(n)$  is obtained by overlap add process of absorbed frames  $s'_{frm}(l, n) = A(m, l) \cdot s_{frm}(l, n)$ .

Fig. 3 depicts the spectrograms and INS values of the noisy-reverberant signal (refer to Fig. 1 (bottom)) processed by the baseline SSS, the AGC technique and the proposed  $ARA_{NSD}$ . Note that, the  $ARA_{NSD}$  is able to absorb masking components of the corrupted signal, e.g. near 0.7 s and 1.5 s, which makes the resulting signal more similar to its anechoic version. Moreover, the proposed method restores the natural non-stationarity behavior raising the INS value from 100 up to 150, which is closer to the direct signal.

#### IV. EXPERIMENTS AND DISCUSSION

Several noisy-reverberant conditions are used to evaluate the SSS [9], the AGC [10] and the proposed  $ARA_{NSD}$  technique in terms of intelligibility. A subset of 24 speakers (16 male and 8 female) are randomly selected from the TIMIT speech database [16], which leads to a total of 240 speech signals (ten for each speaker). From these, 100 are arbitrarily chosen for the test signal and the remaining are used on the speech modeling step of AGC. Each speech segment is sampled at 16 kHz and has, on average, 3 seconds. Two real reverberation rooms from the AIR database [15] are considered in the experiments. The Stairway is characterized by a medium reverberation time ( $T_{60} = 1.1$  s) and a small value of DRR = -9.1 dB. The Aula Carolina room presents parameters of  $T_{60} = 4.9$  s and a higher value of DRR = 7.1 dB. Both RIRs are equalized for a total energy of 17.9 dB. The Babble, SSN, Cafeteria and Chainsaw additive background noises are selected, respectively, from the RSG-10 [17], DEMAND [18] and Freesound.org<sup>2</sup> databases. Except for the SSN, all other noises are characterized with non-stationary behavior.

Speech signals are corrupted considering five SNRs values varying from -3 dB up to 1 dB, where the SNRs are measured between the original unprocessed speech and the background noise. The SNR range is adopted to guarantee ESII and  $ASII_{ST}$  scores between 0.45 and 0.75 for the unprocessed (UNP) speech signal in all scenarios. These values are defined as thresholds of poor and good intelligibility [19], [20], respectively. All UNP intelligibility scores are presented in Tables I and II. The smallest value (ESII = 0.48) is achieved for the Stairway room with the highly non-stationary Chainsaw noise at -3 dB. The Aula

TABLE I  
ESII INTELLIGIBILITY MEASURE [%] FOR UNP SPEECH SIGNALS

Noises	SNR (dB)	Stairway ( $T_{60} = 1.1$ s)					Aula Carolina ( $T_{60} = 4.9$ s)				
		-3	-2	-1	0	1	-3	-2	-1	0	1
Babble		0.53	0.53	0.54	0.55	0.56	0.64	0.65	0.66	0.67	0.67
Cafeteria		0.54	0.55	0.56	0.56	0.57	0.65	0.66	0.67	0.67	0.68
Chainsaw		0.48	0.49	0.50	0.51	0.52	0.57	0.58	0.59	0.61	0.62
SSN		0.52	0.52	0.53	0.54	0.55	0.62	0.63	0.64	0.65	0.66

TABLE II  
 $ASII_{ST}$  INTELLIGIBILITY MEASURE [%] FOR UNP SPEECH SIGNALS

Noises	SNR (dB)	Stairway ( $T_{60} = 1.1$ s)					Aula Carolina ( $T_{60} = 4.9$ s)				
		-3	-2	-1	0	1	-3	-2	-1	0	1
Babble		0.58	0.59	0.60	0.61	0.61	0.68	0.68	0.69	0.70	0.71
Cafeteria		0.60	0.61	0.61	0.62	0.62	0.69	0.70	0.70	0.71	0.71
Chainsaw		0.55	0.56	0.57	0.58	0.58	0.62	0.63	0.64	0.65	0.66
SSN		0.58	0.58	0.59	0.60	0.61	0.66	0.67	0.68	0.69	0.70

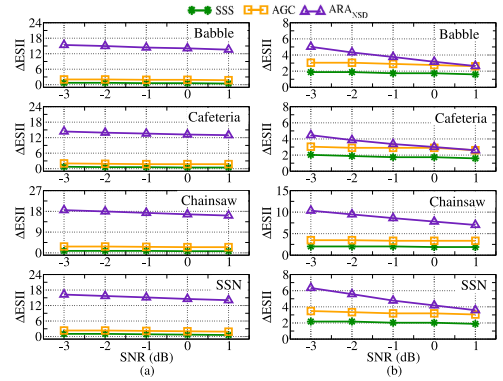


Fig. 4.  $\Delta$ ESII intelligibility improvement [ $\times 10^{-2}$ ] for (a) the Stairway ( $T_{60} = 1.1$  s) and (b) the Aula Carolina ( $T_{60} = 4.9$  s) rooms.

Carolina room with Cafeteria noise at 1 dB presents the highest score of  $ASII_{ST} = 0.71$ . The  $ARA_{NSD}$  operates with 32 ms frames and  $\theta_{INS} = 0.4$ . The maximum value for relevant speech regions  $L'$  is set to 1.2 and the RG importance to  $p = 0.7$  in all scenarios. The sigmoid parameters are fixed to  $k = 17$  for  $d = -0.2$  and  $k' = 13$  for  $d' = 0.5$ . The minimum shift  $S$  is set to 0.05.

##### A. Objective Evaluation of Intelligibility

The ESII [12] and  $ASII_{ST}$  [13] measures are adopted to evaluate the intelligibility improvement under non-stationary noisy-reverberant conditions. The direct path speech signal  $s_{dir}(n)$  is chosen as the reference signal. Since late reverberation and additive background noise reduce intelligibility and are uncorrelated to  $s_{dir}(n)$ , the jointly distortion is obtained by the subtraction  $s(n) - s_{dir}(n)$ . These objective measures are normalized by the intelligibility achieved for the clean unprocessed signal corrupted by SSN noise at 20 dB, considered here as a good intelligibility reference.

The ESII intelligibility improvement ( $\Delta$ ESII) is presented on Fig. 4 for the Stairway and Aula Carolina rooms. In the first case, the  $ARA_{NSD}$  outperforms the competing methods accomplishing more than seven times the AGC value in most of the cases. For the Cafeteria scenario at -1 dB the proposed technique achieves an improvement of 13.6, which corresponds to an assessment eight times the value of 1.7 for the AGC method. As the Stairway room presents a DRR of

<sup>2</sup>Available at www.freesound.org.

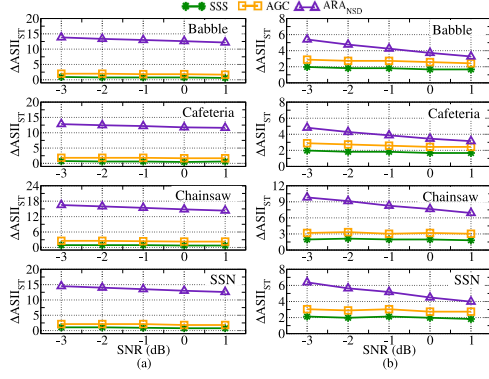


Fig. 5.  $\Delta ASII_{ST}$  intelligibility improvement [ $\times 10^{-2}$ ] for (a) the Stairway ( $T_{60} = 1.1$  s) and (b) the Aula Carolina ( $T_{60} = 4.9$  s) rooms.

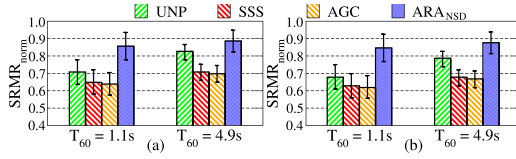


Fig. 6.  $SRMR_{norm}$  for (a) Noise-free reverberation and (b) reverberation with SSN at 0 dB.

−9.1 dB, the reverberation energy in this room is considerably higher than the energy related to the direct signal. This means that the masking components are highlighted in this scenario. As the proposed  $ARA_{NSD}$  is an absorption approach designed to detect such effects, it is able to effectively reduce the temporal coloration. The SSS technique presents the smallest overall intelligibility improvement. Considering the Aula Carolina room, the proposed method also achieves the highest improvement for most of the cases. This is observed for all noisy-reverberant conditions contemplating SNRs below or equal to 0 dB. The best  $\Delta ESII$  results are obtained by  $ARA_{NSD}$  considering the most challenge condition of Chainsaw acoustic noise. The  $ARA_{NSD}$  technique presents similar improvement as AGC for both Babble and Cafeteria noises at 1 dB.

Fig. 5 depicts the  $\Delta ASII_{ST}$  values for both reverberation rooms. For the Stairway room, the proposed method effectively attenuates masking components and attains the highest intelligibility improvement results for all conditions with  $\Delta ASII_{ST}$  values above 10. The  $ARA_{NSD}$  accomplished the highest overall  $\Delta ASII_{ST}$  of 16.5 for the highly non-stationary Chainsaw noise at −3 dB. Baseline technique SSS is outperformed by the  $ARA_{NSD}$  and AGC algorithms in all scenarios. The  $ARA_{NSD}$  also presents the best  $\Delta ASII_{ST}$  intelligibility results for the Aula Carolina room. Once again, the highly non-stationary Chainsaw noise leads to the most challenge condition. In this case, the  $ARA_{NSD}$  is still able to achieve an average improvement of 8.3, compared to 3.6 and 1.2 for the AGC and SSS techniques, respectively.

The  $SRMR$  quality metric [21] estimates the human perceived reverberation effect on speech signals. Its updated version, the  $SRMR_{norm}$  [14], is also selected for objective evaluation. The goal is to distinguish among the three approaches the ones that can better mitigate temporal coloration on speech signals. The direct signal is used as a reference for normalization, such that the  $SRMR_{norm}$  presents values ranging between [0, 1], where 1 determines a reverberation free signal.

Fig. 6 illustrates the average  $SRMR_{norm}$  values for the Stairway and Aula Carolina rooms under a noise-free reverberation

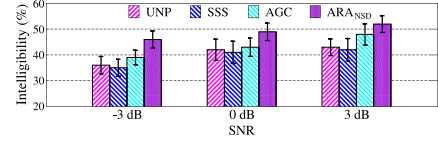


Fig. 7. Perceptual intelligibility evaluation for ISM room ( $T_{60} = 1.0$  s) and SSN additive acoustic noise.

condition (a) and a noisy-reverberation scenario with SSN background noise at 0 dB (b). Note that the  $ARA_{NSD}$  attains the best  $SRMR_{norm}$  values for all situations with a mean of 0.85 and 0.89 for the Stairway and Aula Carolina rooms, respectively. This implies that the proposed method achieved an average quality increment of 0.16 and 0.08 for these rooms when compared with the UNP case. The SSS and AGC techniques present similar behavior, attaining the worst average  $SRMR_{norm}$  values in these scenarios. These results reinforce the capacity of the proposed method to absorb masking components providing intelligibility and quality improvement.

### B. Subjective Intelligibility Evaluation

A listening test [22] with ten native male Brazilian volunteers was conducted considering a closed scenario of phonetic balanced words.<sup>3</sup> Their ages ranged from 22 to 41 years with an average of 32. A simulated room with  $7.0 \times 5.2 \times 3.0$  m<sup>3</sup> and  $T_{60} = 1.0$  s was generated by the image source method (ISM) [23]. The SSN acoustic noise was adopted with SNRs of −3 dB, 0 dB and 3 dB. Ten words were applied for each of 12 test conditions, i.e., three SNR levels for three methods plus the unprocessed case. Participants were introduced to the task in a training session with 8 words. The material was diotically presented using a pair of Roland RH-200S headphones. Listeners heard each word once in an arbitrary presentation order and were asked to indicate the word in a sheet list.

The average intelligibility scores and standard deviations values for each method are presented in Fig. 7. The  $ARA_{NSD}$  improves the intelligibility under all conditions over competing approaches. The proposed method obtained intelligibility improvement scores of 9, 7 and 9 compared to 3, 1 and 5 for the AGC technique for SNR values of −3 dB, 0 dB and 3 dB, respectively. In accordance with findings of [10], [13], SSS attains scores less than or equal to the UNP case.

### V. CONCLUSION

This letter proposed a new time domain absorption approach designed to reduce masking components of speech signals under noisy-reverberant conditions. In this method, the non-stationarity of segments of the corrupted signal is used to detect masking distortions based on a defined threshold. The non-stationarity degree was objectively measured with the INS and was also adopted to determine the absorption procedure. Two reverberant rooms and four acoustic noises were used to compose the noisy-reverberant scenarios. Two intelligibility measures were used for objective evaluation. The results showed that the proposed scheme leads to a higher intelligibility improvement when compared to competing methods. A perceptual listening test corroborated these results. An objective quality measure demonstrated that the proposed technique attains quality improvement.

<sup>3</sup>The complete test database is available at [lasp.ime.eb.br](http://lasp.ime.eb.br).

## REFERENCES

- [1] R. Bolt and A. MacDonald, "Theory of speech masking by reverberation," *J. Acoust. Soc. Amer.*, vol. 21, no. 6, pp. 577–580, 1949.
- [2] A. Nabelek, "Communication in noisy and reverberant environments," *Acoustical Factors Affecting Hearing Aid Performance*. Needham Heights, MA, USA: Allyn & Bacon, 1993, pp. 15–28.
- [3] P. Assmann and Q. Summerfield, "The perception of speech under adverse conditions," in *Speech Processing in the Auditory System*. Berlin, Germany: Springer, 2004, pp. 231–308.
- [4] J. S. Bradley, H. Sato, and M. Picard, "On the importance of early reflections for speech in rooms," *J. Acoust. Soc. Amer.*, vol. 113, no. 6, pp. 3233–3244, 2003.
- [5] T. Gerkmann and R. C. Hendriks, "Unbiased MMSE-based noise power estimation with low complexity and low tracking delay," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 4, pp. 1383–1393, May 2012.
- [6] L. Zão, R. Coelho, and P. Flandrin, "Speech enhancement with EMD and hurst-based mode selection," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 5, pp. 899–911, May 2014.
- [7] R. Coelho and L. Zão, "Empirical mode decomposition theory applied to speech enhancement," in *Signals and Images: Advances and Results in Speech, Estimation, Compression, Recognition, Filtering and Processing*, R. Coelho, V. Nascimento, R. Queiroz, J. Romano, and C. Cavalcante, Eds. Boca Raton, FL, USA: CRC Press, 2015.
- [8] R. Tavares and R. Coelho, "Speech enhancement with nonstationary acoustic noise detection in time domain," *IEEE Signal Process. Lett.*, vol. 23, no. 1, pp. 6–10, Jan. 2016.
- [9] T. Arai, N. Hodoshima, and K. Yasu, "Using steady-state suppression to improve speech intelligibility in reverberant environments for elderly listeners," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 7, pp. 1775–1780, Sep. 2010.
- [10] P. N. Petkov and Y. Stylianou, "Adaptive gain control for enhanced speech intelligibility under reverberation," *IEEE Signal Process. Lett.*, vol. 23, no. 10, pp. 1434–1438, Oct. 2016.
- [11] P. Borgnat, P. Flandrin, P. Honeine, C. Richard, and J. Xiao, "Testing stationarity with surrogates: A time-frequency approach," *IEEE Trans. Signal Process.*, vol. 58, no. 7, pp. 3459–3470, Jul. 2010.
- [12] K. S. Rhebergen and N. J. Versfeld, "A speech intelligibility index-based approach to predict the speech reception threshold for sentences in fluctuating noise for normal-hearing listeners," *J. Acoust. Soc. Amer.*, vol. 117, no. 4, pp. 2181–2192, 2005.
- [13] R. C. Hendriks, J. B. Crespo, J. Jensen, and C. H. Taal, "Optimal near-end speech intelligibility improvement incorporating additive noise and late reverberation under an approximation of the short-time SII," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 5, pp. 851–862, May 2015.
- [14] J. Santos, M. Senoussaoui, and T. Falk, "An improved non-intrusive intelligibility metric for noisy and reverberant speech," in *Proc. IEEE 14th Int. Workshop Acoust. Signal Enhancement*, 2014, pp. 55–59.
- [15] M. Jeub, M. Schafer, and P. Vary, "A binaural room impulse response database for the evaluation of dereverberation algorithms," in *Proc. IEEE 16th Int. Conf. Digit. Signal Process.*, 2009, pp. 1–5.
- [16] J. S. Garofolo, L. F. Lamel, W. M. Fisher, J. G. Fiscus, and D. S. Pallett, "DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1," Philadelphia, PA, USA: NASA STI/Recon, Tech. Rep. N, vol. 93, 1993.
- [17] H. J. Steeneken and F. W. Geurtsen, "Description of the RSG-10 noise database," TNO Inst. Perception, Soesterberg, The Netherlands, Tech. Rep. IZF 3, 1988.
- [18] J. Thiemann, N. Ito, and E. Vincent, "Demand: A collection of multi-channel recordings of acoustic noise in diverse environments," in *Proc. Meetings Acoust.*, 2013.
- [19] *American National Standard: Methods for Calculation of the Speech Intelligibility Index*. New York, NY, USA: Amer. Nat. Standards Inst., 1997.
- [20] B. Sauert and P. Vary, "Near end listening enhancement: Speech intelligibility improvement in noisy environments," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. Proc.*, vol. 1, 2006, pp. I-493–I-496.
- [21] T. H. Falk, C. Zheng, and W.-Y. Chan, "A non-intrusive quality and intelligibility measure of reverberant and dereverberated speech," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 7, pp. 1766–1774, Sep. 2010.
- [22] S. Ghimire, "Speech intelligibility measurement on the basis of ITU-T Recommendation P.863," 2012.
- [23] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *J. Acoust. Soc. Amer.*, vol. 65, no. 4, pp. 943–950, 1979.