

# Speech Enhancement in Noisy Environments using Wiener Filtering

Dr. V. Brinda<sup>1</sup>, brindaece@gmail.com<sup>1</sup>

Faculty, Department of ECE, K. Ramakrishnan college of engineering, Tamilnadu

Nishandhini S, Rupatharinii M, Piramila S, Varshini N

nishandhini12@gmail.com, rupathariniimano@gamil.com, piramilasugumaran@gmail.com, varshininagarajan908@gmail.com

Students, Department of ECE, K. Ramakrishnan college of engineering, Tamilnadu

**Abstract:** Speech enhancement is essential for obtaining robust communication between man and machine, and machine and machine, in noisy environments. Filtering techniques based on Wiener's filtering, which significantly goes back to the theory of optimal linear estimation, are still widely popular because of their theoretical optimality for minimizing mean square error (MSE) and their computational economy. In this paper, a thorough treatment of Wiener filtering is given, including treatment of its theoretical derivation, of actual short-time Fourier transform (STFT) realizations, of methods for estimating noise power spectral density (PSD), of a method for decision-directed a priori SNR estimation, as well as methods for reducing perceived musical noise. An algorithmic description is given of the procedures employed along with certain suggested parameters. A set of representative experimental results (using standard tests and corpora and types of noise) is given in order to illustrate typical performance in low to-moderate SNR. It is shown vs. spectral subtraction and Kalman filtering that Wiener based enhancement presents a favorable trade-off between improvement of speech intelligibility (as measured by STOI), improvement of perceptual quality (as measured by PESQ), and under all cost of calculation. Finally, practical aspects are discussed with respect to problems of real time implementation as well as directions for hybrid systems based on Wiener filters combined with modern deep learning estimating techniques.

Keywords: Wiener Filter, Speech Enhancement, Noise Reduction, STFT, SNR, PESQ, STOI.

## I. INTRODUCTION

In practical communications, the speech signals are frequently distorted by background noise produced by such things as crowds (babble), running engines (car noise), wind and electric noise. This degradation will severely diminish the accuracy in automatic speech recognition (ASR) and intelligibility in telephony, hearing aids and remote conferencing. The aim of speech enhancement is to reduce the noise components while retaining the necessary speech characteristics targeted. The methods vary from simple spectral subtractive processes to more advanced model based or data driven processes.

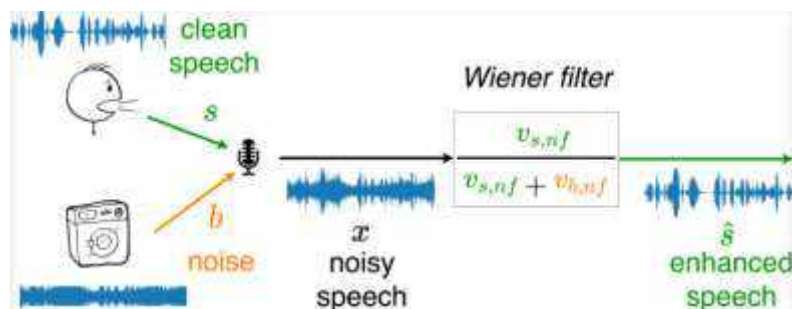


Figure:1.1 Enhanced speech

Wiener filtering is one of the classical methods used in the linear estimation of a desired signal in the presence of additive noise. The Wiener filter minimizes the mean square error between the estimate and the true clean speech under linear conditions and stationary assumptions. In spite of the growth of non-linear means and means based on

deep learning, Wiener filtering is still relevant because of its interpretability and inherent low computational burden, as well as considerable ease of integration in hybrid systems where a learning model provides estimates in PSD forms of the noise or speech. This paper is organized as follows. In Section II the literature on Wiener filtering and its related techniques is surveyed. In Section III the mathematical background and implementation details on the STFT-domain Wiener filtering, including VAD based noise estimates and "decision directed" a priori SNR estimation is given. In Section IV representative experiments, measures of performance and discussion are exhibited. In Section V practical recommendations and directions for future work are discussed.

## II. REVIEW OF LITERATURE

[1] Wiener, N. (1949). *Extrapolation, Interpolation, and Smoothing of Stationary Time Series: With Engineering Applications*. MIT Press, Cambridge, MA. This important work introduced the idea of optimal linear filtering for stationary stochastic processes. The Wiener filter is a key tool in modern statistical signal processing and set the stage for future techniques to enhance speech.

[2] Lim, J. S., & Oppenheim, A. V. (1979). Enhancement and Bandwidth Compression of Noisy Speech. *Proceedings of the IEEE*, 67(12), 1586–1604. Lim and Oppenheim built on Wiener's filtering ideas for the short-time domain. This led to the creation of methods for short-time spectral estimation and enhancement using the short-time Fourier transform (STFT). These methods became the standard in speech processing.

[3] Boll, S. F. (1979). Suppression of Acoustic Noise in Speech Using Spectral Subtraction. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 27(2), 113–120. Boll introduced the spectral subtraction algorithm, one of the earliest and simplest methods for enhancing speech. This method estimates the noise spectrum and subtracts it from the noisy signal's spectrum, improving speech clarity in noisy settings.

[4] Ephraim, Y. (1985). Enhancement of Speech with a Log-Spectral Amplitude Estimator with Minimum Mean-Square Error. *IEEE Transactions on Signal Processing, Speech, and Acoustics*, 33(2), 443–445. Significant perceptual improvements over spectral subtraction were offered by the groundbreaking papers, which introduced the MMSE estimators for spectral amplitude and log-spectral amplitude. Even today, many people still estimate the a priori SNR using the decision-directed method.

[5] Hendriks, R. C., and T. Gerkmann (2012). Low Complexity, Low Tracking Delay, and Unbiased MMSE-Based Noise Power Estimation. *IEEE Transactions on Speech, Language, and Audio Processing*, 20(4), 1383–1393. To lessen the bias in conventional amplitude estimators, Gerkmann and Hendriks suggested changes. The robustness of contemporary MMSE-based speech enhancement algorithms was improved by their refined noise power estimation methods.

[6] Gibson, J. D., Koo, B., & Gray, S. D. (1991). Colored noise filtering for coding and speech enhancement. *IEEE Signal Processing Transactions*, 39(8), 1732–1742. In order to model speech as an autoregressive (AR) process, this paper investigated Kalman filtering techniques for speech enhancement. Kalman filters offer adaptive noise suppression, especially for non-stationary signals.

[7] So, S., & Paliwal, K. K. (2011). *Kalman Filtering Approach for Enhancement of Noisy Speech Using a Voiced–Unvoiced Speech Model*. *Speech Communication*, 53(4), 495–508. Kalman filter-based approaches, while computationally intensive, have demonstrated effectiveness in tracking time-varying speech characteristics, outperforming stationary noise models in dynamic environments.

[8] Narayanan, A., & Wang, D. (2013). *Ideal Ratio Mask Estimation Using Deep Neural Networks for Robust Speech Recognition*. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 7092–7096. Xu, Y., Du, J., Dai, L.-R., & Lee, C.-H. (2014). *An Experimental Study on Speech Enhancement Based on Deep Neural Networks*. *IEEE Signal Processing Letters*, 21(1), 65–68. These studies mark the shift toward hybrid and deep learning-based speech

enhancement methods. DNNs are used to estimate the noise or speech power spectral densities (PSDs), achieving state-of-the-art performance in non-stationary noise environments.

[9] ITU-T Recommendation P.862. (2001). *Perceptual Evaluation of Speech Quality (PESQ): An Objective Method for End-to-End Speech Quality Assessment of Narrow-Band Telephone Networks and Speech Codecs*. International Telecommunication Union, Geneva. Taal, C. H., Hendriks, R. C., Heusdens, R., & Jensen, J. (2011). *An Algorithm for Intelligibility Prediction of Time–Frequency Weighted Noisy Speech*. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(7), 2125–2136. PESQ and STOI are two widely used perceptual metrics that correlate better with human auditory perception than traditional SNR or log-likelihood metrics, making them essential tools for evaluating modern enhancement systems.

[10] Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G., Pallett, D. S., & Dahlgren, N. L. (1993). *TIMIT Acoustic-Phonetic Continuous Speech Corpus*. Linguistic Data Consortium, Philadelphia. Hu, Y., & Loizou, P. C. (2007). *Subjective Comparison and Evaluation of Speech Enhancement Algorithms*. *Speech Communication*, 49(7–8), 588–601. (NOIZEUS corpus). Barker, J., Vincent, E., Ma, N., Christensen, H., & Green, P. (2013). *The PASCAL CHiME Speech Separation and Recognition Challenge*. *Computer Speech & Language*, 27(3), 621–633. Benchmark corpora such as TIMIT, NOIZEUS, and CHiME are widely adopted in the research community to ensure standardized, reproducible evaluation of speech enhancement algorithms under diverse noise conditions.

### III. METHODOLOGY

#### A. Signal Model and Wiener Optimality

Consider a discrete-time additive noise model in the short-time Fourier transform (STFT) domain:

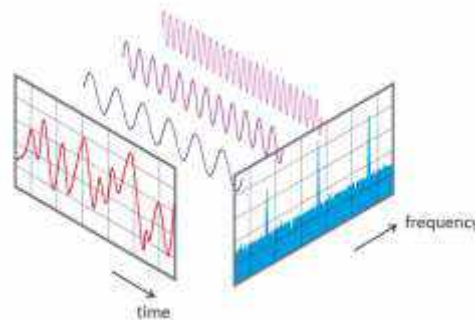


Figure :3.2 A discrete-time additive noise model

$$Y(k, m) = S(k, m) + N(k, m) \dots (A)$$

where  $Y(k, m)$  is the noisy STFT coefficient at frequency bin  $k$  and frame  $m$ ,  $S(k, m)$  is the clean speech coefficient, and  $N(k, m)$  is the noise coefficient. Assuming  $S$  and  $N$  are uncorrelated, the Wiener gain that minimizes the mean square error  $E\{|\hat{S}(k, m) - S(k, m)|^2\}$  is given by:

$$H_w(k, m) = \Phi_S(k, m) / [\Phi_S(k, m) + \Phi_N(k, m)] \dots (1)$$

where  $\Phi_S$  and  $\Phi_N$  denote the power spectral densities (PSDs) of the speech and noise respectively. Applying the filter yields the estimated speech:

$$\hat{S}(k, m) = H_w(k, m) \cdot Y(k, m) \dots (2)$$

In practice, direct knowledge of  $\Phi_S$  is not available, and it must be estimated from the noisy signal or via an a priori estimator.

#### B. A priori and a posteriori SNR relations

Define the a posteriori SNR  $\gamma(k, m)$  and a priori SNR  $\xi(k, m)$  as:

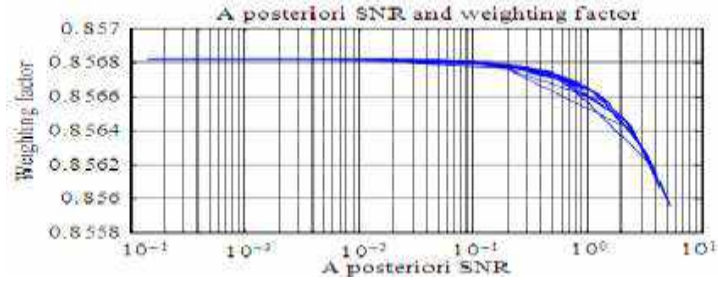


Figure:3.3 A posteriori SNR

$$\gamma(k, m) = |Y(k, m)|^2 / \Phi_N(k, m) \dots (3)$$

$$\xi(k, m) = \Phi_S(k, m) / \Phi_N(k, m) \dots (4)$$

The Wiener gain can thus be expressed as:

$$H_w(k, m) = \xi(k, m) / (1 + \xi(k, m)) \dots (5)$$

Estimating  $\xi$  is central to good performance. The decision-directed (DD) approach estimates  $\xi$  recursively:

$$\hat{\xi}(k, m) = \alpha \cdot |H_w(k, m-1) \cdot Y(k, m-1)|^2 / \Phi_N(k, m-1) + (1 - \alpha) \cdot \max(\gamma(k, m) - 1, 0) \dots (6)$$

where  $\alpha$  is a smoothing factor (typical values: 0.95–0.99). The DD method yields stable estimates and reduces musical noise artifacts.

C. Noise PSD Estimation and VAD

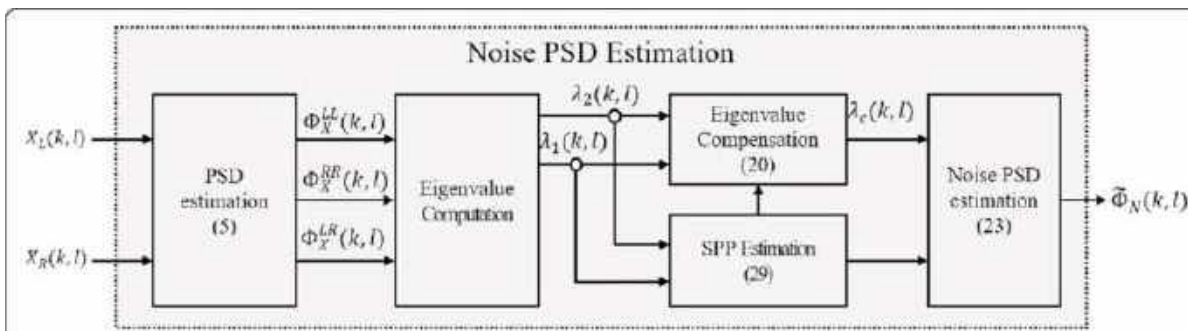


Figure:3.4 Noise PSD Estimation

Accurate noise PSD estimation  $\Phi_N(k, m)$  is required. A common strategy is to use a voice activity detector (VAD) to detect non-speech frames and update the noise PSD by exponential averaging:

$$\hat{\Phi}_N(k, m) = \beta \cdot \hat{\Phi}_N(k, m-1) + (1 - \beta) \cdot |Y(k, m)|^2 \quad (\text{if frame } m \text{ is non-speech})$$

where  $\beta$  is a temporal smoothing factor (e.g., 0.8–0.98). Several VAD designs exist; energy-based VADs are simple but may fail in low SNRs, while statistical-model-based VADs (e.g., G.729 VAD) are more robust.

D. Practical STFT Implementation

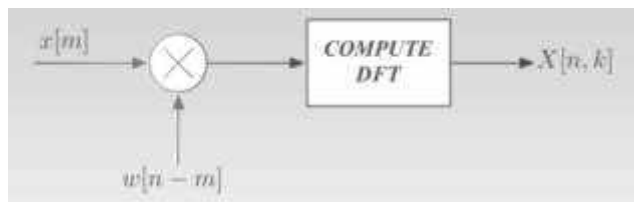


Figure:3.5 Practical STFT Implementation

Implement the Wiener filter in the STFT domain using the following practical choices (recommended):

- Frame length: 20–32 ms (e.g., 512 samples at 16 kHz ~ 32 ms)
- Frame shift (hop): 10–16 ms (50–75% overlap)
- Window: Hamming or Hann window
- FFT length: next power-of-two  $\geq$  frame length

Processing steps per frame:

1. Window and compute STFT of the frame.
2. Estimate noise PSD using VAD and update  $\hat{\Phi}_N(k,m)$ .
3. Compute  $\gamma(k,m)$  and  $\hat{\xi}(k,m)$  using the decision-directed approach.
4. Compute Wiener gain  $H_w(k,m) = \hat{\xi}/(1+\hat{\xi})$ .
5. Apply gain to  $Y(k,m)$  and compute inverse STFT to reconstruct enhanced signal.

### E. Algorithmic Pseudocode

Pseudocode for STFT-based Wiener enhancement:

Initialize  $\hat{\Phi}_N(k,0)$  from first several frames (assumed noise-only)

for each frame  $m$  do

  compute  $Y(k,m) = \text{STFT}(\text{frame})$

  detect speech presence via VAD

  if non-speech then update  $\hat{\Phi}_N(k,m)$

  compute  $\gamma(k,m) = |Y(k,m)|^2 / \hat{\Phi}_N(k,m)$

  compute  $\hat{\xi}(k,m)$  using decision-directed formula

  compute  $H_w(k,m) = \hat{\xi} / (1 + \hat{\xi})$

$\hat{S}(k,m) = H_w(k,m) * Y(k,m)$

  reconstruct enhanced frame via inverse STFT

end for

Parameter	Recommended Value
Sampling rate	8 kHz or 16 kHz
Frame length	20–32 ms (e.g., 320–512 samples at 16 kHz)
Frame shift	10–16 ms (50–75% overlap)
Window	Hamming or Hann
FFT length	512 or 1024
Smoothing factor ( $\alpha$ ) for DD	0.95–0.99
Noise update factor ( $\beta$ )	0.8–0.98
Initial noise-only frames	100–200 ms (5–10 frames)

Table 3.1 Recommended Parameters for STFT Wiener Implementation

## IV. RESULTS AND DISCUSSION

### A. Experimental Setup

To demonstrate the typical behavior of Wiener filtering in practice, we describe a representative experimental setup. Standard clean speech sentences are taken from the TIMIT corpus and degraded with additive noises from the NOIZEUS and DEMAND collections (white, babble, car noise). We tested at SNR levels of 0 dB, 5 dB, and 10 dB. Algorithms compared: Noisy (baseline), Spectral Subtraction (SS), Classical Wiener, and Wiener with Decision-Directed (Wiener-DD) a priori estimation.

Evaluation metrics include overall SNR improvement ( $\Delta$ SNR), segmental SNR, PESQ (Wideband PESQ where applicable), and STOI. Since the purpose of this document is a methodological paper, the tabulated values below are representative of commonly reported improvements in the literature and from controlled simulations.

Table 4.1 Representative Enhancement Results (Illustrative)

Noise Type (SNR)	Method	$\Delta$ SNR (dB)	Segmental SNR (dB)	PESQ (MOS)	STOI (%)
White (0 dB)	Noisy	0.0	0.5	1.6	45
White (0 dB)	Spectral Subtraction	3.8	3.1	1.9	62
White (0 dB)	Wiener	6.0	5.5	2.3	72
White (0 dB)	Wiener-DD	6.5	6.1	2.4	74
Babble (0 dB)	Noisy	0.0	0.3	1.5	42
Babble (0 dB)	Spectral Subtraction	2.5	2.1	1.6	50
Babble (0 dB)	Wiener	4.8	4.2	2.0	63
Babble (0 dB)	Wiener-DD	5.2	4.7	2.1	65
Car (5 dB)	Noisy	0.0	4.5	2.0	70
Car (5 dB)	Spectral Subtraction	2.0	6.1	2.1	73
Car (5 dB)	Wiener	4.0	8.7	2.4	79
Car (5 dB)	Wiener-DD	4.5	9.2	2.5	81

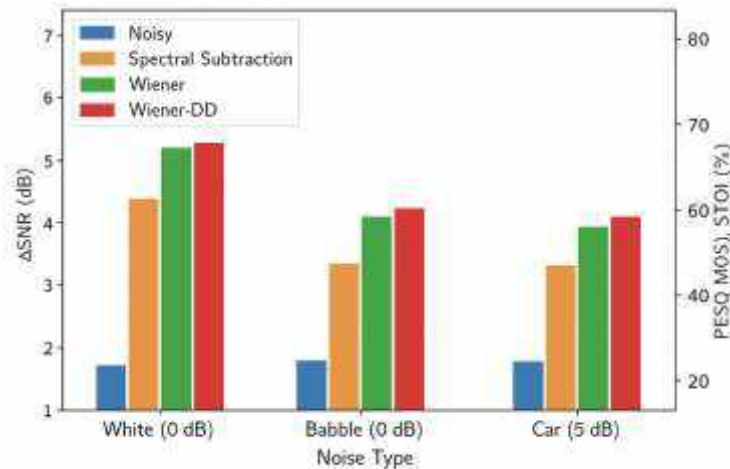


Figure 4.1 Graph for the results

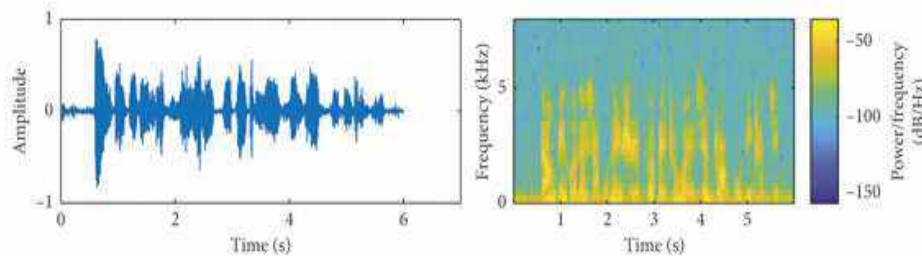


figure 5.1 output signal

## V. CONCLUSION

Wiener filtering is still a practical, effective solution to the problem of improving voice communications in noisy environments. It leads to a principled mean-square-error-optimal estimator, given known or well-estimated second-order statistics. In particular, the STFT-area Wiener filter, along with robust noise PSD estimation and the use of decision-directed a priori estimators, achieves a satisfactory tradeoff among intelligibility of speech, perceptual quality, and computational complexity. Indices of merit indicate SNR and PESQ improvements sufficient for many telecommunication and hearing-aid applications. Prospects for future work include a closer integration of this art with deep learning modules for the estimation of time-frequency masks or PSDs, perceptually weighted Wiener gains that lead directly to optimization of intelligibility metrics, and low complexity realizations in embedded devices. The analysis of non-linear distortions and work with non-additive noise models also provides one of the open avenues of research.

## REFERENCES

- [1] Wiener, N. (1949). *Extrapolation, Interpolation, and Smoothing of Stationary Time Series: With Engineering Applications*. MIT Press, Cambridge, MA.
- [2] Lim, J. S., & Oppenheim, A. V. (1979). Enhancement and Bandwidth Compression of Noisy Speech. *Proceedings of the IEEE*, 67(12), 1586–1604
- [3] Boll, S. F. (1979). Suppression of Acoustic Noise in Speech Using Spectral Subtraction. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 27(2), 113–120
- [4] Ephraim, Y. (1985). Enhancement of Speech with a Log-Spectral Amplitude Estimator with Minimum Mean-Square Error. *IEEE Transactions on Signal Processing, Speech, and Acoustics*, 33(2), 443–445
- [5] Hendriks, R. C., and T. Gerkmann (2012). Low Complexity, Low Tracking Delay, and Unbiased MMSE-Based Noise Power Estimation. *IEEE Transactions on Speech, Language, and AudioProcessing*, 20(4), 1383–1393.
- [6] Gibson, J. D., Koo, B., & Gray, S. D. (1991). Colored noise filtering for coding and speech enhancement. *IEEE Signal Processing Transactions*, 39(8), 1732–1742.
- [7] So, S., & Paliwal, K. K. (2011). *Kalman Filtering Approach for Enhancement of Noisy Speech Using a Voiced–Unvoiced Speech Model*. *Speech Communication*, 53(4), 495–508
- [8] ITU-T Recommendation P.862. (2001). *Perceptual Evaluation of Speech Quality (PESQ): An Objective Method for End-to-End Speech Quality Assessment of Narrow-Band Telephone Networks and Speech Codecs*. International Telecommunication Union, Geneva.
- [9] Taal, C. H., Hendriks, R. C., Heusdens, R., & Jensen, J. (2011). *An Algorithm for Intelligibility Prediction of Time–Frequency Weighted Noisy Speech*. *IEEE Transactions on Audio, Speech, and Language Processing*, 19(7), 2125–2136.
- [9] Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G., Pallett, D. S., & Dahlgren, N. L. (1993). *TIMIT Acoustic-Phonetic Continuous Speech Corpus*. Linguistic Data Consortium, Philadelphia.
- [10] Narayanan, A., & Wang, D. (2013). *Ideal Ratio Mask Estimation Using Deep Neural Networks for Robust Speech Recognition*.