# Analysis and synthesis of vocal flutter and vocal jitter

*Jean Schoentgen[1], Philipp Aichinger[2]*

[1]Bio-, Electro- And Mechanical Systems, Faculty of Applied Sciences, Université Libre de Bruxelles, Brussels, Belgium
[2]Medical University of Vienna, Division of Phoniatrics-Logopedics, Department of Otorhinolaryngology, Vienna, Austria

`jschoent@ulb.ac.be, philipp.aichinger@meduniwien.ac.at`

## Abstract

Perturbations of the strict periodicity of the glottal vibrations are relevant features of the voice quality of normophonic and dysphonic speakers. Vocal perturbations in healthy speakers are assigned different names according to the range of the typical perturbation frequencies. The objective of the presentation is to model jitter and flutter, which are in the $> 20Hz$ and $10Hz - 20Hz$ range respectively, via a simulation of the fluctuations of the tension of the thyro-arytenoid muscle and compare simulated perturbations to jitter and flutter observed in vowels sustained by normophonic speakers.

**Index Terms**: vocal frequency jitter, vocal frequency flutter, muscle tension fluctuations, motor units, motor neurones.

## 1. Introduction

Perturbations of the periodicity of the glottal vibrations are relevant features of the voice quality of normophonic and dysphonic speakers [1]. In healthy speakers, one distinguishes vocal perturbations according to the range of the typical perturbation frequencies. One so identifies vocal frequency drift and "physiological" tremor ($< 2Hz$), "neurological" tremor ($2Hz - 10Hz$), vocal frequency flutter ($10Hz - 20Hz$) as well as vocal frequency jitter ($> 20Hz$) [2, 3, 4].

The objective of the presentation is to model jitter and flutter via a simulation of the fluctuations of the tension of the thyro-arytenoid (TA) muscle and compare simulated perturbations to jitter and flutter observed in vowels sustained by normophonic speakers. The purpose of the modelling is to better understand the origin of vocal perturbations and to synthesise different voice qualities. Physiological and neurological tremor are not discussed, because they involve breathing and pulsatile blood flow as well as unsteady firing of the motor neurones, which are not modelled.

The model of the fluctuations of the tension of the TA-muscle is inspired by [5], but drops problematic *a priori* assumptions of the latter. Dropping assumptions enables simulating perturbations with a wider range of parameter values.

Titze has proposed a model that relates observed vocal jitter to the fluctuations of the tension of the TA muscle [5]. The tension of a skeletal muscle is the outcome of the concurrent activity of many motor units of which each comprises a motor neurone that innervates a group of muscle fibres that contract after the arrival of an action potential, i.e. an electrical spike emitted by the neurone. The duration between two adjacent spikes is the inter-spike interval. The simultaneous contraction of several muscle fibres that are under the control of a single neurone is called a muscle twitch. The muscle tension is the outcome of the superposition of many muscle twitches in space because of the concurrent firing of many motor neurones and in time because of the rapid succession of the action potentials of a single motor neurone.

In Titze's model, the typical duration between successive action potentials that initiate muscle twitches is equal to $1/\lambda$, $\lambda$ being the average firing rate of the neurone [5]. Shift values drawn from Gaussian distributions then randomly move the neurone spikes and muscle twitches from their periodic default positions.

Repositioning periodic spikes randomly by normally distributed shifts may be problematic because of the symmetry and infinite support of the Gaussian distribution. Indeed, spikes may be shifted forward and backward in time and spike $n$ may occur before spike $n - 1$ or later than spike $n + 1$ when the variability of the inter-spike intervals is large. In [5], unacceptable shifts have been avoided by keeping the coefficient of variation $\nu$ of the inter-spike intervals small ($\leq 0.15$). This constraint has been based on the histogram of one recording of the spike intervals of one motor unit of a human TA-muscle.

Hereafter, we follow Deger et al.'s [9] advice and model inter-spike intervals by means of a Gamma distribution that generates positive inter-spike intervals for any firing rate and any coefficient of variation, without risking violating causality.

One observes that small $\nu$ values boost vocal jitter and large values boost vocal flutter. Jitter in human voices must therefore examined together with flutter to discover physiologically acceptable model parameters. But, vocal flutter is only rarely discussed in the literature on voice quality and few data are available [6]. Also, the distinction between flutter and jitter is not explained satisfactorily. Vocal flutter is occasionally referred to as a "rapid tremor", which suggests that the distinction may be auditory [8]. Indeed, "rapid" modulations of the vocal frequency are perceived as rough and "slow" modulations as tremor, with a boundary that is approximately situated at $20Hz$ [7].

Section 2 describes the model of the fluctuations of the tension of the TA muscle; section 3 describes the decomposition of recorded vocal cycle lengths into jitter, flutter and tremor; section 4 reports jitter and flutter in vowels sustained by normophonic speakers; section 5 describes simulations of flutter and jitter, compares these with observed data and discusses muscular causes of vocal frequency perturbations in the context of other causes that have been considered in the literature.

## 2. A model of the neurological sources of vocal aperiodicities

*Muscle twitch model*
The contraction of the muscle twitch model $s(t) = \frac{t}{T_r} e^{1 - \frac{t}{T_r}}$ starts at time $t = 0$. The maximum contraction is reached after a rise time $t = T_r$, subsequently the twitch decays slowly over an interval $> T_r$ [5, 10]. The latency and rise time have been

equal to $2ms$ and $30ms$ respectively [5]. The latency is the short delay between the arrival of a spike and the beginning of the contraction and is taken into account in the spike positions. A muscle twitch $s(t - t_s)$ caused by the arrival of a motor neurone spike at time $t_s$ is positioned by convolving $s(t)$ with the unit spike $\delta(t - t_s)$.

*Spike and twitch time series of a single motor unit*

Deger et al. [9] have shown that Gamma distributions are good approximations of the observed distributions of neural inter-spike intervals $\Delta_{isi}$. The activity of a motor neurone is therefore simulated by drawing $\Delta_{isi}$ values from a Gamma distribution $\mathcal{G}(k, b)$. Parameters $k$ and $b$ depend on the average firing rate $\lambda$ of the neurone and the coefficient of variation $\nu$ of the inter-spike intervals: $k = 1/\nu^2$ and $b = \lambda/\nu^2$. The temporal positions $t_s$ of the motor neurone spikes are obtained iteratively: $t_{s+1} = t_s + \Delta_{isi}$. The physical constraint $\Delta_{isi} > 0$ is satisfied mathematically.

Neurones enter into a refractory period $T_{refr}$ after firing. The physiological constraint $\Delta_{isi} > T_{refr}$ implies $\nu < 1$, which is a necessary condition only. Intervals $\Delta_{isi}$ are therefore monitored and dropped when $< T_{refr}$.

Also, Titze has shifted spike positions $t_s$ by a random constant $T_T$ drawn from a Gaussian distribution $\mathcal{N}(0, 1/\lambda)$ [5]. The purpose has been to prevent motor neurones from firing in quasi-synchrony when the firing rate $\lambda$ is the same for all motor units and the coefficient of variation $\nu$ is small. Spike positions $t_s$ are therefore replaced by updated positions $t_s + T_l + T_T$ taking into account *ad hoc* shift $T_T$ and muscle twitch latency $T_l$ that has been referred to earlier.

The muscle twitch time series $t_w(t)$ owing to one motor unit is then obtained by convolving muscle twitch $s(t)$ with the unit spike sequence: $t_w(t) = \sum_s \delta(t - t_s) \star s(t)$.

*Muscle tension*

Muscle tension $T_m(t)$ is the sum at each time instant $t$ of the muscle twitch time series $t_w(t)$ of $N_{mu}$ motor units. However, the superposition of $N_{mu}$ muscle twitch time series can be replaced by the superposition of $N_{mu}$ spike position time series followed by one convolution when taking into account that a convolution is a linear operation and assuming that twitch shape $s(t)$ is the same for all motor units: $T_m(t) = \sum_{mu,s} \delta(t - t_{mu,s}) \star s(t)$. The gain in computation time is typically of two orders of magnitude.

*Perturbations of the instantaneous vocal frequency and synthetic vocal cycle lengths*

The instantaneous vocal frequency $f_o(t)$ has been shown to be proportional to the square root of the muscle tension. The relative instantaneous perturbations are therefore the following [5].

$$(f_o(t) - \bar{F}_o)/\bar{F}_o = \sqrt{T_m(t)/\bar{T}_m} - 1. \quad (1)$$

$\bar{F}_o$ designates the unperturbed vocal frequency and $\bar{T}_m$ designates the time-averaged instantaneous muscle tension $T_m$. The presence of ratio $T_m(t)/\bar{T}_m$ in relation (1) enables computing the relative perturbations of the vocal frequency without explicitly modelling the absolute values of the muscle twitches and neural spikes. For speech synthesis, the experimenter may assign a value to the unperturbed vocal frequency $\bar{F}_o$ to obtain instantaneous frequency $f_o(t)$, which also enables calculating the vocal cycle length time series by integrating numerically the equation $d\phi = 2\pi \times f_o(t) \times dt$ and counting the number of time steps that phase $\phi$ takes to increase by $2\pi$.

## 3. Analysis of vocal cycle length time series

A vocal cycle length time series of a normophonic speaker may be decomposed into a jitter, flutter, tremor and residue time series. All perturbations have been removed from the residue up to the ultra-slow modulations owing to $f_o$ drift and physiological tremor.

The decomposition into jitter, flutter and tremor is a straightforward generalisation of the conventional analysis of vocal jitter. Indeed, jitter is obtained by smoothing the cycle length time series by means of a running average, followed by subtracting the smoothed from the original time series and assigning the difference to jitter [11].

The decomposition into jitter, flutter, tremor and a residue involves carrying out the previous step three times. Each step involves smoothing (i.e. low-pass filtering) by a running average, subtracting the smoothed from the un-smoothed time series, storing the difference and replacing the un-smoothed by the smoothed time series before the next step. The successive subtractions guarantee that the decomposition is exact, that is, the sum of the stored differences and the residue is exactly equal to the raw time series. The jitter, flutter, tremor and residue time series are obtained by fixing the cut-off frequency of the running average to $20Hz$, $10Hz$ and $2Hz$ respectively [2]. The residue is the filtered time series after step 3.

The number of samples involved in a running average is obtained as follows. A running average based on a rectangular window is equivalent to a low-pass filter the transfer function of which is a $sinc$ function. The position in $Hz$ where the main lobe of the $sinc$ crosses the frequency axis is equal to the inverse of the length in $sec$ of the rectangular window, which so assigns a value to the cut-off frequency [12]. When $F_c$ designates the cut-off frequency and $F_s$ the sampling frequency of the time series then the length $N_s$ in number of samples of the rectangular window is equal to the ratio $F_s/F_c$ rounded up to the nearest odd integer. In raw vocal cycle length time series, $F_s$ is equal to the average vocal frequency $\bar{F}_o$.

Before decomposition, the cycle length time series can be resampled with a constant sampling step by interpolation. Resampling replaces the (feebly) variable sampling step of the raw cycle length time series by a constant time step and upsampling enables a more fine-grained selection of the length of the running average, which must be an odd number of samples.

## 4. Jitter and flutter in vowels sustained by normophonic speakers

The corpus comprised 36 stable [a] vowel fragments $2sec$ long sustained by 18 male and 18 female normophonic speakers who scored zero on the perceptual G (*grade*) scale. The sampling frequency was $44.1kHz$ [13]. The vowel fragments have been manually segmented and the raw cycle lengths have been obtained by means of PRAAT [11].

The objective has been to obtain the jitter and flutter in stable vowel fragments and compare their sizes. The sizes of jitter and flutter are summarised by the averages of the absolute values of the corresponding time series obtained via the decomposition described in section 3 [11]. Before decomposition, the raw length time series had been upsampled to $1kHz$ by interpolation. Figure 1 shows an example of the upsampled raw lengths, jitter, flutter, tremor and residue for a $2sec$ fragment of a vowel [a] sustained by a male normophonic speaker.
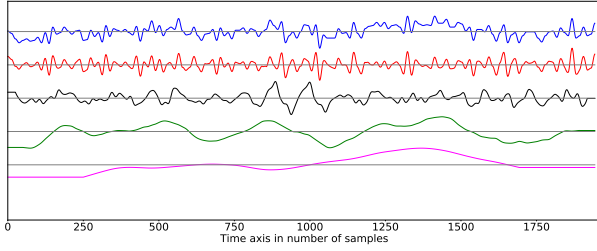
Figure 1: *Upsampled cycle length time series, jitter, flutter, tremor and physiological tremor/drift (from top to bottom). The time series are scaled individually. The average cycle length has been subtracted from the raw cycle length time series before decomposition and display.*

Table 1: *Quartiles of the relative size in % of flutter & jitter and the ratio of flutter to jitter for 36 [a] vowel fragments ($2sec$) sustained by 18 male and 18 female speakers.*

| Quartiles | Flutter (%) | Jitter (%) | Ratio |
|-----------|-------------|------------|-------|
| Min | 0.04 | 0.11 | 0.30 |
| Q1 | 0.10 | 0.13 | 0.64 |
| Med. | 0.15 | 0.17 | 0.77 |
| Q3 | 0.19 | 0.23 | 0.94 |
| Max | 0.34 | 0.45 | 1.33 |

Table 1 reports the quartiles of the size in % of jitter and flutter relative to the average vocal cycle length for the 36 vowel fragments. The rightmost column reports the ratio of the size of flutter to the size of jitter, which has been less than $0.94$ for 75 % of the speakers.

## 5. Simulation of flutter and jitter via a model of muscle tension fluctuations

Perturbations of the vocal cycle lengths have been simulated by means of the model of the muscle tension fluctuations described in section 2. The purpose has been, firstly, to examine the influence of the variability of the neural inter-spike intervals on the size of jitter and flutter to discover a physiologically acceptable range for the coefficient of variation $\nu$ and, secondly, to regress on plausible model parameters the relative sizes of flutter and jitter to examine the influence of the parameters on the perturbations.

The model parameter values have been randomly selected in the range $(M - 0.5M, M + 0.5M)$ relative to a middle value $M$, which has been fixed according to published data [5, 14]. The midrange values are given in Table 2. The sampling frequency of the instantaneous perturbations has been $200kHz$. For each simulation, the model parameters have been the same for all motor units, which is an idealisation that eases the interpretation of statistical analyses. No attempt has been made to fit the simulated to the observed perturbations.

Vocal cycle length time series $2s$ long have been obtained via relation (1) for different unperturbed $\bar{F}_o$ values. The cycle length time series have been upsampled to $1kHz$ via interpolation to facilitate the comparison with the cycle length time series observed for human speakers. Figure 2 shows an example of a synthetic cycle length time series together with the jitter and flutter time series obtained via decomposition

Table 2: *Midrange values $M$ of the model parameters that have been randomly selected in the range $(M - 0.5 \times M, M + 0.5 \times M)$ for simulation.*

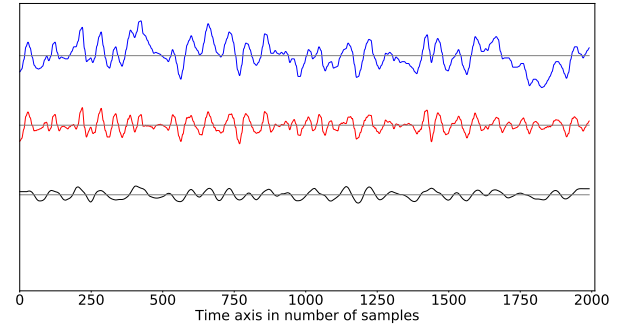| Parameter | Symbol | Midrange $M$ | Reference |
|-----------|--------|--------------|-----------|
| Firing rate | $\lambda$ | $30Hz$ | [14, 5] |
| Coeff. var. isi | $\nu$ | $0.15$ | [5] |
| Coeff. var. isi | $\nu$ | $0.30, 0.45, 0.60$ | |
| Refractory period | $T_{refr}$ | $2.5ms$ | [14] |
| # motor units | $N_{mu}$ | $100$ | [5] |
| Vocal frequency | $\bar{F}_o$ | $150Hz$ | |



Figure 2: *Simulated raw cycle lengths, jitter and flutter (from top to bottom). Jitter and flutter have been obtained by decomposing raw cycle lengths obtained by simulation. The average cycle length has been subtracted from the raw cycle length time series before decomposition and display. The model parameters have been: $\bar{F}_o = 137Hz$, $\lambda = 16Hz$, $\nu = 0.215$, $N_{mu} = 141$. The length of the time series is $2s$.*

(section 3).

A total of $4 \times 1000$ simulations have been carried out with the midrange coefficient of variation $\nu$ of the neural inter-spike intervals fixed to $0.15, 0.30, 0.45, 0.60$ respectively. Table 3 reports quartiles of the flutter and jitter in % as well as the flutter to jitter ratio for each midrange coefficient of variation $\nu$.

One observes that, firstly, jitter and flutter are of the expected order of magnitude when considered independently; secondly, jitter as well as flutter increase with coefficient $\nu$, but flutter increases faster than jitter; thirdly, a comparison of the flutter to jitter ratios in Tables 1 and 3 suggests that physiologically acceptable $\nu$ values are $< 0.3$, which corroborates assumptions made in [5].

Table 4 reports for midrange $\nu = 0.15$ the linear regression of the sizes of flutter and jitter as well as of the flutter to jitter ratio on the model parameters listed in Table 2. The parameter and perturbation values have been z-normalised so that the regression weights may be compared.

Predictibly, parameters the increase of which one would expect to smooth muscle tension (i.e. an increase of the firing rate $\lambda$ and number of motor units $N_{mu}$) decrease flutter and jitter, whereas parameters the increase of which one would expect to coarsen muscle tension (i.e. an increase of coefficient of variation $\nu$ of neural inter-spike intervals) increase flutter and jitter. Finally, Table 4 reports a mild dependence of relative jitter and flutter on unperturbed vocal frequency $\bar{F}_o$, which has been inserted into relative perturbations (1) to synthesise vocal cycle length time series.

Table 3: *Quartiles of the flutter, jitter (in %) and flutter to jitter ratio of* 1000 *simulations with the midrange coefficient of variation values fixed to* 0.15, 0.30, 0.45, 0.60 *respectively.*

|  | **Quartiles** | $\nu = 0.15$ | 0.30 | 0.45 | 0.60 |
|---|---|---|---|---|---|
| **Flutter** | Q1 | 0.07 | 0.13 | 0.20 | 0.26 |
|  | Med. | 0.10 | 0.18 | 0.27 | 0.35 |
|  | Q3 | 0.30 | 0.25 | 0.36 | 0.45 |
| **Jitter** | Q1 | 0.12 | 0.16 | 0.19 | 0.22 |
|  | Med. | 0.19 | 0.22 | 0.26 | 0.29 |
|  | Q3 | 0.30 | 0.31 | 0.33 | 0.35 |
| **Ratio** | Q1 | 0.40 | 0.72 | 0.94 | 1.10 |
|  | Med. | 0.53 | 0.89 | 1.10 | 1.20 |
|  | Q3 | 0.70 | 1.03 | 1.20 | 1.40 |

Table 4: *Regression weights of z-scaled model parameters listed in Table* 2 *with regard to z-scaled flutter, jitter and flutter to jitter ratio for* 1000 *simulations with midrange* $\nu = 0.15$. *Weights in italic are statistically not significant. The bottom row reports the "adjusted" coefficient of determination* $R^2$.

| **Parameters** | **Flutter(%)** | **Jitter(%)** | **Ratio** |
|---|---|---|---|
| $\bar{F}_o$ | +0.05 | +0.16 | −0.28 |
| $\lambda$ | −0.62 | −0.87 | +0.63 |
| $\nu$ | +0.52 | +0.07 | +0.55 |
| $N_{mu}$ | −0.28 | −0.25 | *−0.02* |
| $T_{refr}$ | *+0.02* | *+0.01* | *+0.02* |
| $R^2$ | +0.76 | +0.87 | +0.74 |

The weights of the flutter to jitter ratio reflect the individual dependencies of jitter and flutter. The dependence of the ratio on coefficient of variation $\nu$ is the outcome of a boost with increasing $\nu$ that is stronger for flutter than for jitter. The dependence of the ratio on firing rate $\lambda$ is the outcome of a decline with increasing $\lambda$ that is stronger for jitter than for flutter. The lack of dependence of the ratio on the number of motor units $N_{mu}$ is the consequence of a decline with increasing $N_{mu}$ that is similar for jitter and flutter.

*Fluctuations of the tension of the TA-muscle in the context of other causes of vocal jitter*

Hereafter, we would like to discuss the relevance of the model of the vocal perturbations described in sections 2 and 4 in the context of other causes that have been listed in the literature as possible explanations of vocal perturbations in healthy and dysphonic speakers [1].

Model [5] has been written about in the literature occasionally only. One reason is that it does not explain changes in vocal jitter that cannot be directly related to perturbations due to the activity of the TA muscle. For instance, vocal jitter may increase in healthy speakers after vocal loading in dry air, after injecting atropine in the folds or owing to dehydration, menstruation or light laryngitis, which are conditions that are expected to increase the viscosity of the cover of the vocal folds [15]. Similarly, one observes that vocal jitter decreases in high-pitched or falsetto voices, which involve the activity of muscles other than the TA muscle.

We discuss here a simple model, which suggests that the size of existing perturbations may be modulated fold-internally

and that, therefore, an observed increase or decrease of vocal jitter does not necessarily imply that an autonomous cause of perturbation has been activated or deactivated.

The model is kinematic and assumes that the body (muscle) and cover of a vocal fold vibrate sinusoidally at the same average frequency, but at different amplitudes $A_m$ and $A_c$ [16]. The instantaneous phases $\phi_m$ and $\phi_c$ of the body and cover are assumed to be the same up to a perturbation $\theta_{pert}$ of the phase of the fold body. The instantaneous phase of the fold cover is assumed to be unjittered. When the folds do not touch, the movement of the fold edge $x_{edge}$ is the sum of the sinusoidal motions of the body and cover, disregarding a constant abduction.

$$x_{edge} = A_c \times \sin(\phi_c) + A_m \times \sin(\phi_m) \quad (2)$$

$$\theta_{pert} = \phi_c - \phi_m \quad (3)$$

Assuming for simplicity's sake that $\theta_{pert}$ is small, that is, $\cos\theta_{pert} \approx 1, \sin\theta_{pert} \approx \theta_{pert}$ and $\tan^{-1}\theta_{pert} \approx \theta_{pert}$ and applying an elementary trigonometric relation [17], one obtains the following approximate expression.

$$x_{edge} \approx (A_c + A_m) \times \sin(\phi_c - \frac{1}{1 + \frac{A_c}{A_m}} \times \theta_{pert}) \quad (4)$$

The argument of the sinusoid in (4) shows that the perturbations of the frequency of vibration of the vocal folds are modulated by the amplitude ratio $A_c/A_m$ of the cover and muscle. This would suggest that a relative decrease of the amplitude of vibration $A_m$ of the muscle in falsetto voices would decrease observed perturbations, whereas a relative decrease of the amplitude of vibration $A_c$ of the cover of the fold owing to an increased viscosity, for instance, would increase observed perturbations. We do not claim that model (4) is physiologically accurate. But, it qualitatively predicts observed phenomena, even though it is a very simple model of possible fold-internal modulations. Fold-internal modulation of the perturbations of the vocal frequency may therefore warrant further investigation.

## 6. Conclusion

The simulation of vocal jitter and flutter owing to muscle tension fluctuations generates perturbations the size of which is of the expected order of magnitude $(0\%-1\%)$. Also, a comparison of the sizes of flutter and jitter in vowels sustained by healthy speakers confirms the appropriateness of the assumption made in model [5] that the variability of the inter-spike intervals of a motor neurone must be small. The observation that the neural firing rate and variability of the neural inter-spike intervals influence not only the size but also the frequency distribution of the vocal perturbations suggests that the distinction between jitter and flutter may not be perceptual only.

## 7. Acknowledgments

# 8. References

[1] J. Kreiman and D. Sidtis, *Foundations of Voice Studies: An Interdisciplinary Approach to Voice Production and Perception*. Wiley-Blackwell, 2011, page 55.

[2] E. H. Buder and E. A. Strand, "Quantitative and graphic acoustic analysis of phonatory modulations: the modulogram," *J. Speech, Language, Hearing Res.*, vol. 46, pp. 475–490, 2003.

[3] P. R. Cook, Ed., *Music, Cognition and Computerized Sound:An Introduction to Psychoacoustics*. Cambridge, Massachusetts: The MIT Press, 1999, page 199.

[4] I. R. Titze, *Principles of Voice Production*. Englewood Cliffs, NJ: Prentice Hall, 1994, page 332.

[5] I. Titze, "A model for neurologic sources of aperiodicity in vocal fold vibration," *J. Speech, Hearing Res.*, vol. 34, pp. 460–472, 1991.

[6] G. H. Alzamendi and G. Schlotthauer, *Describing voice period variability by means of time series structural analysis*, 2017, proceedings 10th International Workshop: Models and Analysis of Vocal Emissions for Biomedical Applications, Firence, Italy, pages 11-14.

[7] H. Fastl and E. Zwicker, *Psychoacoustics, Facts and Models*. Springer, 2007, e-book.

[8] A. E. Aronson, S. R. Silber, L. O. Ramig, and W. S. Winholtz, "Rapid voice tremor, or 'flutter,' in amyotrophic lateral sclerosis," *Annals of Otology, Rhinology, and Laryngology*, vol. 101, pp. 511– 518, 1992.

[9] M. Deger, M. Helias, C. Boucsein, and S. Rotter, "Statistical properties of superimposed stationary spike trains," *J. Comput. Neurosci.*, vol. 32, pp. 443–463, 2012.

[10] I. P. Herman, *Physics of the Human Body*. Berlin, Heidelberg: Springer, 2007, page 281.

[11] P. Boersma and D. Weeninck, *Praat: doing phonetics by computer [Computer program]*, 2014, [Version 5.4.04, retrieved 2014 from http://www.praat.org]. [Online]. Available: http://www.praat.org

[12] S. W. Smith, *The Scientist and Engineer's Guide to Digital Signal Processing*. California Technical Pub, 1997, e-book, page 277.

[13] U.MA., "Database for perceptual analysis of voice quality," 2018, [accessed 23-January-2018]. [Online]. Available: http://www.atic.uma.es/index-atic.html

[14] R. Roark, C.L., J. Li, S. Schaefer, A. Adam, and C. D. Luca, "Multiple motor unit recordings of laryngeal muscles: The technique of vector laryngeal electromyography," *The Laryngoscope*, vol. 112, pp. 2196–2202, 2002.

[15] E. Abberton and A. Fourcin, *Instrumental Clinical Phonetics: Electrolaryngography*, M. J. Ball and C. Code, Eds. Wiley, 2006, page 119.

[16] B. H. Story and I. R. Titze, "Voice simulation with a body-cover model of the vocal folds," *J. Acoustic. Soc. Am.*, vol. 97, pp. 1249–1260, 1995.

[17] H. S. Black, *Modulation Theory*. Van Nostrand, 1953, page 221.