



Published in final edited form as:

*IEEE Trans Biomed Eng.* 2010 March ; 57(3): 626–633. doi:10.1109/TBME.2009.2033037.

## Automatic Detection of Swallowing Events by Acoustical Means for Applications of Monitoring of Ingestive Behavior

**Edward S. Sazonov [Member, IEEE],**

The Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699, USA (phone: 315-268-3914)

**Oleksandr Makeyev [Student Member, IEEE],**

The Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699, USA

**Stephanie Schuckers [Member, IEEE],**

The Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699, USA

**Paulo Lopez-Meyer,**

The Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699, USA

**Edward L. Melanson, and**

The Division of Endocrinology, Metabolism, and Diabetes, University of Colorado Denver, Aurora, CO 80045, USA

**Michael R. Neuman [Senior Member, IEEE]**

The Department of Biomedical Engineering, Michigan Technological University, Houghton, MI 49931, USA

Edward S. Sazonov: esazonov@ieee.org; Oleksandr Makeyev: mckehev@cias.clarkson.edu; Stephanie Schuckers: sschucke@clarkson.edu; Paulo Lopez-Meyer: lopezmp@clarkson.edu; Edward L. Melanson: ed.melanson@ucdenver.edu; Michael R. Neuman: mneuman@mtu.edu

### Abstract

Our understanding of etiology of obesity and overweight is incomplete due to lack of objective and accurate methods for Monitoring of Ingestive Behavior (MIB) in the free living population. Our research has shown that frequency of swallowing may serve as a predictor for detecting food intake, differentiating liquids and solids, and estimating ingested mass. This paper proposes and compares two methods of acoustical swallowing detection from sounds contaminated by motion artifacts, speech and external noise. Methods based on mel-scale Fourier spectrum, wavelet packets, and support vector machines are studied considering the effects of epoch size, level of decomposition and lagging on classification accuracy. The methodology was tested on a large dataset (64.5 hours with a total of 9,966 swallows) collected from 20 human subjects with various degrees of adiposity. Average weighted epoch recognition accuracy for intra-visit individual models was 96.8% which resulted in 84.7% average weighted accuracy in detection of swallowing events. These results suggest high efficiency of the proposed methodology in separation of swallowing sounds from artifacts that originate from respiration, intrinsic speech, head movements, food ingestion, and ambient noise. The recognition accuracy was not related to body mass index, suggesting that the methodology is suitable for obese individuals.

### Index Terms

biomedical signal processing; obesity; pattern recognition; swallowing; wearable devices

## I Introduction

The world is still losing in the battle with the obesity epidemic. According to WHO, in 2005 there were approximately 1.6 billion overweight and at least 400 million obese adults worldwide [1]. Current trends are unsettling: 2015 projections predict 2.3 billion overweight and 700 million obese adults worldwide. Obesity is one of the risk factors for development of chronic diseases and presents a serious health problem. A recent study [2] suggested that effects of obesity on global health may be comparable to those of cancer. Though the etiology of obesity is a topic of ongoing scientific debate, regulation of food intake may be an important factor for maintaining a healthy weight [3] in the environment that provides abundance of inexpensive, highly palatable and energy dense foods, while requiring only minimal levels of physical activity [4].

While various methods have been developed for accurate and objective characterization of physical activity [5], at the present time, there is no accurate, inexpensive, non-intrusive way for objective Monitoring of Ingestive Behavior (MIB) in free living conditions. The most precise method of measuring energy intake is the Doubly-Labeled Water (DLW) technique which provides accurate estimates of caloric energy intake over long periods of time (10–14 days), if subjects remain weight stable. However, the DLW technique cannot identify daily intake patterns. Dietary self-report methods like food frequency questionnaires [6], self-reported diet diaries [7], and multimedia diaries [8] have been shown to be inaccurate and underreport daily intake.

Our recent research [9] has shown that frequency of swallowing can serve as a predictor for accurate detection of food intake, differentiation between liquid and solid foods and estimation of ingested mass, with high frequency of swallowing being indicative of ingestion. Thus, an affordable wearable MIB device can be created for objective characterization of food intake. Such a device would use the proposed acoustical method to detect both spontaneous and food intake swallows as they happen throughout the day without any conscience input from the user. A second-stage algorithm [9] would use the recognized swallows to detect and characterize food intake based on the frequency of swallowing. Potentially, such a device can reduce intake underreporting because: 1) monitoring is objective and does not rely on self-report; 2) continuous capture of spontaneous swallows indicates whether the sensor system is being worn or not, thus preventing or detecting intentional misreport.

While the weight gain is ultimately determined by energy balance (energy intake minus energy expenditure) and the proposed MIB device by itself cannot capture the energy content of a meal, such a device can provide valuable information about ingestion that is not available at this time. The device can also help diagnose and treat dangerous behaviors leading to weight gain, such as unconscious snacking [10], night eating [11], and evening [12] or weekend overeating [13]. The device may also find applications in diagnostics and treatment of disorders not directly related to obesity such as inadvertent weight loss (cachexia), anorexia and bulimia as well as dysphagia.

The proposed method for acoustical detection of swallowing events is the first and fundamental step in implementation of the wearable MIB device. The swallowing detection does not need to differentiate between spontaneous and food intake swallows as the food intake detection relies only on frequency of swallowing events. Algorithms presented in [9] can be applied as the second step of processing to detect and characterize food intake from the sequence of swallows.

This paper demonstrates high accuracy of swallowing event detection by acoustical means on the largest dataset to date by the methodologies based on mel-scale Fourier Spectrum (msFS) and Wavelet Packet Decomposition (WPD) for time-frequency representation, and Support Vector Machines (SVM) for automatic recognition of characteristic sound of swallowing. It also contains assessment of the size of a near-optimal time decomposition window and effects of the decomposition level and epoch lagging on accuracy of swallowing detection suggesting that epoch duration used in earlier publications may not be optimal. Furthermore, assessment of recognition accuracy as a function of subject's Body Mass Index (BMI) shows that the proposed acoustical method is suitable for obese individuals. Finally, it is demonstrated that proposed methods have substantial tolerance to the sound artifacts resulting from food intake, intrinsic speech and background noise and thus may be suitable for free living conditions.

The description is organized as follows: Section 2 presents the background on assessment of swallowing sound signals and currently used automatic swallowing detection methods. Section 3 provides a brief description of the data collection process. Section 4 presents a detailed description of the proposed methodology. Experimental results are presented in Section 5 followed by the Discussion and Conclusions.

## II Acoustical Detection of Swallowing Events

At the present time videofluoroscopy and EMG are considered the gold standard in studies of deglutition. Videofluoroscopy depends on bulky and potentially unsafe equipment while EMG is too invasive due to frequently used subcutaneous placement of electrodes in the masseter, suprahyoid and infrahyoid muscles [14] to avoid interference from the muscles of the neck. Other reported sensors include a variety of strain devices [14–16]. However, most of the reported results indicate that detection of swallowing by a laryngeal strain sensor is not appropriate for obese subjects since under chin adipose deposits inhibit reliable detection of swallows. Use of accelerometer placed over the suprasternal notch of trachea as suggested by [17–19] may also be not appropriate for obese individuals for the same reasons. Detection of the characteristic swallowing sound created by the specific motion of laryngopharynx can be performed by a microphone which is significantly less invasive and more effective for obese individuals than the methods listed above.

Several methods have been proposed for assessment of swallowing sounds using signal processing and pattern recognition techniques. Papers [17–20] presented methodologies for automatic decomposition of the tracheal sound signal into swallowing and respiratory segments in applications to dysphagia. The signal decomposition techniques utilized such features as autoregressive coefficients, root mean square values of the signal in time domain, average power of the signal within several frequency bands, waveform fractal dimension and Discrete Wavelet Transform on time windows (epochs) ranging in duration from 12.5 to 100 ms. Reported recognition rates were in the range from 78.54% [17] to 93% [19] although the sound recordings did not include any speech or noise.

Rejection of artifacts arising from ingestion, intrinsic speech and external noise is an issue that needs serious consideration. In the MIB applications, artifacts such as breathing, talking, throat cleaning, head movements, etc. may be confused with swallowing thus decreasing the accuracy of the recognition [21]. The feasibility of sound artifact rejection was tested in [22] where swallowing sound recognition was performed using the Limited Receptive Area neural classifier in combination with short-time Fourier transform and continuous wavelet transform. The methods in [22] achieved 100% accuracy in classification of swallowing sounds on a limited dataset containing swallowing sounds, motion artifacts, talking and music, although practical applications to large datasets were limited by high computational burden of the method.

A recently reported method of automated swallowing detection that was tested in the presence of artifacts originating in talking, head movements, food ingestion, and respiration was presented in [23]. The data was collected from six healthy subjects using a sensor collar containing surface electromyography electrodes and a stethoscope electret microphone. A total of 7.93 hours of data with 1,265 swallows was acquired. Feature similarity search combined with an agreement of the detectors fusion method was used for classification. Four-fold cross validation was used with three folds used for training and one for validation. The average recognition rate of 70% was obtained for labeling epochs of 250ms as belonging to swallows or non-swallows.

In summary, acoustical detection of swallowing events, as presented herein, may present a non-invasive and convenient method suitable for use by obese individuals. However, the field of swallowing sound detection is relatively unexplored with a significant need to focus on realistic conditions with presence of various sound artifacts. Another key consideration is the choice of the epoch duration and lagging for signal analysis. Epoch sizes used in [18,20,23] are substantially shorter (12.5–250 ms) than the average duration of a swallow (0.86 s) [24] and thus may represent only a partial segment of a swallowing sound or require a large number of time lags. The goal of the methodology proposed here is to consider acoustical swallow recognition as a method which may be appropriate for obese individuals; compare two popular signal time-frequency decompositions; investigate selection of key parameters of time-frequency transforms such as epoch duration and level of decomposition; and to test the proposed methods on a challenging dataset that resembles free living conditions and includes artifacts of various origins.

### III Data collection

The data used in this paper were collected in human study reported in [25] where the details of the protocol, hardware, sensors and reliability of the manual scoring procedure are reported, but with no attempt to automatically recognize swallowing events. The following is a summary of the human study. The subject population included 20 volunteers, of which 7 had BMI greater than 30 (obese). Each subject participated in four visits, each of which consisted of a 20-minute resting period, followed by a meal, followed by another 20-minute resting period. Out of 80 collected visits, 10 were discarded due to data collection errors [25]. Selection and sequence of foods were fixed for each meal and represented different physical properties of the food such as crispiness, softness/hardness and tackiness, all of which may impact both the artifacts arising from chewing sounds and the swallowing sound itself. To evaluate the impact of a mealtime conversation on the accuracy of swallowing detection, the subjects were involved in a dialogue with a member of the research team during the second and fourth visits and ate in silence during the first and third visits. Additionally, background noise (city noise, restaurant noise and music) were played during the second and fourth visits to simulate realistic environments where people may be eating. The subjects were monitored by a multi-modal sensor system which included an IASUS NT (IASUS Concepts Ltd) throat microphone located over laryngopharynx. The microphone provided a dynamic range of  $46 \pm 3$  dB with a frequency range of 20 Hz to 8000 Hz. Amplified signals were recorded through a line-in input of a standard sound card at a sampling rate of 44100 Hz. The recordings were manually scored to mark the boundaries of each swallow. The evaluation of inter-rater reliability reported in [25] showed high reliability of manual scores (0.98 average intra-class correlation) for manual scoring of swallows.

### IV. Methods

The proposed methods are based on popular time-frequency decompositions: mel-scale Fourier Spectrum (msFS) and Wavelet Packet Decomposition (WPD) with classification performed

by Support Vector Machines (SVM). Time-frequency decomposition and feature extraction based on WPD and msFS is widely used for processing of physiological signals, such as, for example, heart sounds [26] and lung sounds [27]. SVM is a supervised learning method that has a sound theoretical basis, is robust to overfitting (loss of generalization on noisy or incomplete data [28]) and capable of producing very complex decision boundaries.

### A. Feature extraction by Wavelet Packet Decomposition

First, the sound stream was split into a series of overlapping epochs with fixed duration  $D$  and step  $S$ . A Hanning window was applied to each epoch. Second, a time-frequency decomposition of each epoch was obtained using Wavelet Packet Decomposition creating  $2^N$  wavelet packets (where  $N$  is the level of decomposition) [29]. A packet on the previous level is decomposed into two packets on the next level as

$$w_{2n}(t) = \sqrt{2} \sum_k h_k w_n(2t - k), w_{2n+1}(t) = \sqrt{2} \sum_k g_k w_n(2t - k)$$

where  $h_k$  is the low-pass Finite Impulse Response (FIR) filter and  $g_k$  is the high-pass FIR filter such as  $g_k = (-1)^k h_{1-k}$ . The WPD was computed using Coiflet C4 wavelet. Advantages of the Coiflet wavelet include near linear phase, good amplitude response and fast computation [30]. WaveLab [31] package for Matlab was used to perform WPD. Third, each wavelet packet was converted into a scalar feature forming a feature vector  $f_i$  of length  $2^N$  for each epoch. The chosen feature was the unbiased estimate of entropy [32]. Fourth, to account for the time-varying structure of a swallow, a time-lagged feature vector was produced by merging feature vectors of the  $K$  adjacent epochs:  $f_i' = \{f_{i-K}, f_i, f_{i+K}\}$ .

### B. Feature extraction by mel-scale Fourier transform

First, segmentation of the sound signal into overlapping epochs was performed identically to the one used for WPD. Second, the Fourier amplitude spectrum  $F(k)$  of length  $L$  was computed for every epoch. Third, a mel-scale triangle filter bank  $M_i(k)$  [33] was used to compute  $2^N$  point feature vector  $f_i$  (where  $N$  is an equivalent to WPD's level of decomposition) defined as

$$f_i = \log \left( \sum_{k=0}^{L/2} F(k) M_i(k) \right), i = 0, \dots, N - 1. \text{ Finally, the time-lagged vector } f_i' \text{ was obtained in the same way as for WPD. Fig. 1 shows a segment of the sound recording containing a swallow and its respective representation obtained by WPD and msFS processing with decomposition level } N = 8, \text{ epoch duration } D = 1.5 \text{ s and step } S = 0.2 \text{ s.}$$

### C. Support Vector Machines

The time-lagged feature vectors  $f_i'$  obtained either through WPD or msFS processing were used as inputs for training and validation of an SVM classifier [28]. The choice of the SVM as a classifier was defined by sound theoretical foundation and robust performance of SVM classifiers. A comparison of SVM performance to performance of 16 classification and 9 regression methods on 21 data sets for classification and 12 data sets for regression [34] ranked SVM as one the best techniques on most data sets, especially for classification. LibSVM package for Matlab [35] was used for training the SVM classifier using the Gaussian radial basis kernel function. Optimal parameters of the SVM classifier were found by a grid search procedure.

#### D. Optimal epoch duration and decomposition level

Optimal epoch duration  $D$ , epoch step size  $S$ , decomposition level  $N$  and number of lags  $K$  were determined in a grid search procedure. The epoch duration and step size ( $D/S$ ) were taken from a set  $\{3.0/0.4, 1.5/0.2, 0.75/0.1, 0.375/0.05\}$  seconds which represents progressively finer time resolutions. Decomposition level both for WPD and msFS was taken as  $N \in \{5, 6, 7, 8, 9\}$  thus producing from 32 to 512 features for each epoch. The number of lags  $K$  was either 0 or 1 since a higher number of lags produced long feature vectors which substantially slowed the classifier. Since a grid search procedure is time consuming, it was performed on randomly selected two visits that included noise and talking during the meal and thus presented a harder classification case. The grid search procedure repeatedly trained classifiers defined by various combinations of  $D/S$ ,  $N$  and  $K$ . The validation accuracy was used to evaluate the goodness of parameters. Training and validation were performed with 34% of the data (one fold) used for training and 66% (two folds) used for validation. The accuracy of swallowing detection was estimated as described further.

#### E. Training and Validation

The pairs of feature vectors and class labels to be used as inputs for the SVM classifier were obtained in the following way: if any part of the epoch belonged to a swallow marked in the manual score the epoch label was set as '1' (swallow epoch), otherwise it was set as '-1' (non-swallow epoch). Individual intra-visit models were built for 70 visits of 20 subjects. The training and validation sets were formed by taking into account the highly non-homogeneous structure of each visit. For example, a period of quiet resting with no talking and no food intake will not have enough variability in the data to train a classifier that would work reliably if talking or food intake is introduced. Since talking, food intake and external noise are introduced at various times in each visit, a longitudinal segmentation was used. Each visit was divided into 55 segments of equal duration, each segment 1 minute in duration on average. Three-fold cross-validation was performed with two folds used for training and one fold used for validation.

#### F. Accuracy of detecting swallowing instances

Predicted class labels represent accuracy of the classifier on epoch level and do not correspond well to the accuracy of detection of swallowing events. Transition from the epochs to swallowing events was done by identifying all situations where either Manual Score (MS) or Automatic Score (AS) indicated presence of a swallow and calculating the numbers of true positives, false positives and false negatives in terms of swallowing events. A true positive ( $T_+$ ) was counted if both MS and AS contained continuous sequences of epochs marked as swallows intersecting at one or more epochs or on the sequence boundary (Fig. 2, *a*). A false positive ( $F_+$ ) was identified if the AS marked a swallow which was not present in the MS (Fig. 2, *b*). A true negative ( $T_-$ ) was counted if both MS and AS contained continuous sequences of epochs marked as non-swallows intersecting at one or more epochs (Fig 2, *c*). A false negative ( $F_-$ ) was counted if the MS marked a swallow which was not present in the AS (Fig 2, *d*). The accuracy of swallowing events detection was then estimated using weighted accuracy =

$$\frac{T_+ + T_-}{T_+ + T_- + F_+ + F_-}, \text{ Sensitivity} = \frac{T_+}{T_+ + F_-} \text{ and Specificity} = \frac{T_-}{T_- + F_+}.$$

#### V. Results

The graphs obtained by the grid search of optimal epoch duration, decomposition level and number of lags on a subset from 2 visits are shown in Fig. 3 which suggests the best parameters for WPD processing: 9<sup>th</sup> level of decomposition on 1.5 s epochs. For msFS processing the best parameters are at 7<sup>th</sup> level of decomposition on 1.5 s epochs. These parameters with and without lagging were used to process throat microphone signal collected in 70 visits. SVM training was performed with misclassification penalty  $C = 10$  and Gaussian kernel width parameter  $\gamma$

= 0.05. Results obtained in per-epoch recognition and detection of swallowing events are presented in Table I.

The best average weighted accuracy in terms of epochs and swallows was produced by msFS-7 with 3 lags and found to be  $96.8 \pm 1.4\%$  for epochs and  $84.7 \pm 6.9\%$  for swallows. The distribution of average weighted accuracy in classification of epochs and swallowing events versus the subject's BMI and corresponding linear fit of the data are presented in Fig. 4. To assess the impact of sound artifacts on accuracy of identifying swallowing events the average weighted swallowing accuracy was also computed individually for the four non-overlapping parts of the validation set corresponding to the following categories: periods of no food intake and no talking (88.0%), periods of no food intake with talking (86.4%), periods of food intake and no talking and background noise (86.2%), and periods of food intake with talking and background noise (82.9%).

## VI Discussion

One of the goals of this work was to determine the optimal duration of an epoch since durations reported in existing literature [18,19,23] varied over a wide range of 12.5–250 ms. As Fig. 3 shows that the epoch duration of 1.5s clearly demonstrates the highest recognition accuracy both for msFS and WPD with or without lagging. This corresponds well with the mean duration of swallow which in our study was found to be 1.15s with a standard deviation of 0.29s (based on analysis of 10,686 swallows), comparable to previously reported duration of 0.86s [24]. Thus, the epoch duration of 1.5s is sufficient to completely include an average swallow. We believe that such choice of the epoch duration is one of the reasons that our recognition rate is substantially higher than per epoch accuracy of 70% reported in [23] where the authors used a 0.25s epoch which cannot cover a complete swallow. Although our study excluded dysphagic subjects, we may anticipate that recognition of longer-than-normal dysphagic swallows [36] may benefit from a longer epoch.

Fig. 3 also demonstrates accuracy growth with increase in the level of decomposition. The most pronounced increase for msFS is observed up until the 7<sup>th</sup> level of decomposition. As Table 1 shows a lagged version produces higher overall accuracy due to better preservation of accuracy during transitioning from epochs to swallows on some of the visits. Lagging takes the feature evolution over time into account and thus produces more accurate results. The non-lagged version of WPD processing clearly peaks at the 9<sup>th</sup> level of decomposition and trends toward further growth. Unfortunately, higher levels of decomposition result in unacceptably long processing times both for feature extraction and classification. The lagged version of WPD behaves somewhat erratically (which may be attributed to a limited dataset used in the grid search procedure) but clearly peaks at 8<sup>th</sup> level of decomposition, confirming that 8<sup>th</sup>–9<sup>th</sup> levels are probably near the optimum for WPD. Fig. 3 and Table 1 also demonstrate that msFS time-frequency decomposition clearly outperforms WPD resulting in higher recognition accuracy. A possible explanation is non-linear scaling of the frequencies by msFS which allows for a better representation of the lower frequencies which contain most of the energy of a swallowing sound.

One of major advantages of the current study is that it was designed to be close to real life conditions and include sound artifacts originating from chewing of food of different textures, talking, head movements, occasional intrinsic sounds (for example, coughing), and background noise of various origins. Thus, the classifier had to deal with a significantly more complex problem than previous studies while achieving a comparable (vs. 93% in [19]) or better performance (vs. 79% in [18]). The closest study that allows direct comparison is [23] which achieved the epoch-accurate average recognition rate of 70%. For comparison, the methodology proposed here yielded the average weighted epoch accuracy of 96.8% that relates

to 84.7% average weighted accuracy in detection of swallowing events. Furthermore, our study utilized a wider variety of solid foods (cheese pizza, an apple, and a peanut butter sandwich) with varying physical properties that directly impact the sounds of mastication [37] and subsequently influence swallowing recognition. Another advantage is unrestricted consumption of liquids which were limited in [23] to 5ml and 15ml of volume at a time. Liquid consumption is characterized by a very high swallowing frequency [9] in which identification of individual swallows is difficult due to the fact that consecutive swallows may be recognized as one. The results show that artifact sounds negatively impact the recognition accuracy but not to a degree that would render the method unusable. As expected, the highest recognition accuracy is observed for quiet periods of no food intake (88.0%) and the lowest recognition accuracy is observed for periods of food intake combined with talking and background noise (82.9%). Thus, application of noise cancellation techniques may further improve on the classification accuracy.

As Fig. 4 suggests the accuracy of detecting swallowing events is very likely not to be dependent on the subject's BMI. While more data are needed to appropriately test the effect of obesity, this may be an important advantage of the acoustical approach of detecting swallowing events. The highest BMI of a volunteer in the study was 42.1 which is considered severe (morbid) obesity. Even for this volunteer the swallowing identification accuracy was greater than 80%. Thus, these results suggest that the proposed methodology could be used for monitoring of food intake in obese individuals.

The reported experimental results were obtained on the dataset containing 64.5 hours of data with 9,966 swallows collected from 20 subjects with the experimental conditions resembling those of food consumption in free living. To our knowledge this is largest dataset collected to date. In addition, the manual score of swallows used for training of the classifiers has known reliability metrics [25]. Overall, the proposed methodology showed good performance in testing on a more complicated dataset than any of the previous studies. The next step in the development of the acoustical method of the detection of swallowing is development of inter-visit individual and group models that could be practically applied for automatic scoring of the swallowing sound recordings. The desired accuracy of the identification of swallowing is another question that needs further investigation. However, the methods for detection of food intake and prediction of ingested mass [9] should offers some tolerance to the errors in detection of swallowing instances since they rely on multiple swallows and relatively long time windows (up to 2 minutes).

The results of this study also have important implications for the original intent to use automatic recognition of swallowing sounds in a wearable device for monitoring of ingestion. The time sequence of swallows detected by the proposed method can be further processed by algorithms in [9] to detect and characterize food intake to achieve real-time monitoring of ingestion. With the rapid progression of computing power available in modern ubiquitous platforms (cell phones, PDA) the proposed MIB methodology can be implemented as a wearable device allowing for real-time biofeedback to individuals. Such a wearable device may potentially find numerous applications in research, clinical nutrition and self-monitoring of food intake by general population.

## VII Conclusion

In this paper we describe two automatic acoustical swallowing detection methods for use in MIB applications. The methods were based on combination of mel-scale Fourier Spectrum (msFS) or Wavelet Packet Decomposition (WPD) and Support Vector Machines. The proposed methodology was tested on the data collected from 20 human subjects with 35% of the subjects being obese with Body Mass Index (BMI) of at least 30 and the average BMI of 28.53 using



a multimodal data collection system designed for non-invasive monitoring of chewing and swallowing. The total duration of data used for training and validation was 64.5 hours including 9,966 swallows which makes it the largest dataset to date. Average weighted epoch classification accuracy of 96.8% resulted in 84.7% average weighted accuracy in detection of swallowing events. Optimal duration of a sound time slice was found to be 1.5s which corresponds well to statistics of swallowing duration. The msFS decomposition with 3 lags clearly outperformed WPD in recognition accuracy. A study of impact of food intake, talking and background noise on accuracy of swallowing detection suggests robustness of the proposed methodology to such events as well as its ability to accurately separate swallowing sounds from sound artifacts that originate in respiration, talking, head movements, food ingestion, and ambient noise. The method was also demonstrated to work equally well for both obese and non-obese subjects. The described methodology and sensors may be implemented in a wearable monitoring device, thus enabling MIB applications in free living individuals.

## Acknowledgments

This work was supported in part by the National Institute of Heart, Lung and Blood under Grant R21HL083052-02.

The authors would like to acknowledge the anonymous reviewers for their suggestions on improving this manuscript.

## References

1. World Health Organization. Global Infobase for Overweight and Obesity. Oct. 2008  
<http://www.who.int/mediacentre/factsheets/fs311/en/index.html>
2. Olshansky SJ, Passaro DJ, Hershow RC, Layden J, Carnes BA, Brody J, Hayflick L, Butler RN, Allison DB, Ludwig DS. A Potential Decline in Life Expectancy in the United States in the 21st Century. *N Engl J Med* Mar;2005 352:1138–1145. [PubMed: 15784668]
3. Flatt JP. Substrate utilization and obesity. *Diabetes Rev* 1996;4:433–449.
4. Hill JO, Wyatt HR, Reed GW, Peters JC. Obesity and the Environment: Where Do We Go from Here? *Science* Feb;2003 299:853–855. [PubMed: 12574618]
5. Ainslie P, Reilly T, Westerterp K. Estimating human energy expenditure: a review of techniques with particular reference to doubly labelled water. *Sports Medicine (Auckland, N.Z)* 2003;33:683–98.
6. Weber JL, Reid PM, Greaves KA, DeLany JP, Stanford VA, Going SB, Howell WH, Houtkooper LB. Validity of self-reported energy intake in lean and obese young women, using two nutrient databases, compared with total energy expenditure assessed by doubly labeled water. *European Journal of Clinical Nutrition* Nov;2001 55:940–50. [PubMed: 11641742]
7. De Castro JM. Methodology, correlational analysis, and interpretation of diet diary records of the food and fluid intake of free-living humans. *Appetite* Oct;1994 23:179–92. [PubMed: 7864611]
8. Kaczkowski CH, Jones PJ, Feng J, Bayley HS. Four-day multimedia diet records underestimate energy needs in middle-aged and elderly women as determined by doubly-labeled water. *The Journal of Nutrition* Apr;2000 130:802–5. [PubMed: 10736333]
9. Sazonov E, Schuckers S, Lopez-Meyer P, Makeyev O, Melanson E, Neuman M, Hill J. Toward objective monitoring of ingestive behavior in free living population. *Obesity*. 2009 advance online publication. 10.1038/oby.2009.153
10. Ward, C. *Compulsive Eating: The Struggle to Feed the Hunger Inside*. The Rosen Publishing Group; 1998.
11. Stunkard AJ, Grace WJ, Wolff HG. Night eating syndrome. *Eating Disorders and Obesity. A Comprehensive Handbook* 2002:183–188.
12. Kant AK, Ballard-Barbash R, Schatzkin A. Evening eating and its relation to self-reported body weight and nutrient intake in women, CSFII 1985–86. *Journal of the American College of Nutrition* Aug;1995 14:358–63. [PubMed: 8568112]
13. Haines PS, Hama MY, Guilkey DK, Popkin BM. Weekend eating in the United States is linked with greater energy, fat, and alcohol intake. *Obesity Research* Aug;2003 11:945–9. [PubMed: 12917498]

14. Ertekin C, Aydogdu I, Seçil Y, Kiylioglu N, Tarlaci S, Ozdemirkiran T. Oropharyngeal swallowing in craniocervical dystonia. *Journal of Neurology, Neurosurgery, and Psychiatry* Oct;2002 73:406–11.
15. Stellar E, Shrager EE. Chews and swallows and the microstructure of eating. *The American Journal of Clinical Nutrition* Nov;1985 42:973–82. [PubMed: 4061369]
16. Pehlivan M, Yüceyar N, Ertekin C, Celebi G, Ertaş M, Kalayci T, Aydoğdu I. An electronic device measuring the frequency of spontaneous swallowing: digital phagometer. *Dysphagia* 1996;11:259–64. [PubMed: 8870354]
17. Lazarek LJ, Moussavi ZK. Automated algorithm for swallowing sound detection. *Proc Canadian Med and Biol Eng Conf* 2002:1–4.
18. Aboofazeli, M.; Moussavi, Z. Automated classification of swallowing and breath sounds. *Conference Proceedings: Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conference*; 2004. p. 3816-9.
19. Aboofazeli, M.; Moussavi, Z. Automated extraction of swallowing sounds using a wavelet-based filter. *Conference Proceedings: Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conference*; 2006. p. 5607-10.
20. Aboofazeli M, Moussavi Z. Analysis of swallowing sounds using hidden Markov models. *Medical and Biological Engineering and Computing* Apr;2008 46:307–314. [PubMed: 18000695]
21. Das A, Reddy NP, Narayanan J. Hybrid fuzzy logic committee neural networks for recognition of swallow acceleration signals. *Computer Methods and Programs in Biomedicine* 2001;64:87–99. [PubMed: 11137191]
22. Makeyev, O.; Sazonov, E.; Schuckers, S.; Lopez-Meyer, P.; Baidyk, T.; Melanson, E.; Neuman, M. Recognition of Swallowing Sounds Using Time-Frequency Decomposition and Limited Receptive Area Neural Classifier. *Proceedings of AI-2008, The Twenty-eighth SGAI International Conference on Innovative Techniques and Applications of Artificial Intelligence*; Cambridge. 2008. p. 33-46.
23. Amft O, Tröster G. Recognition of dietary activity events using on-body sensors. *Artificial Intelligence in Medicine* 2008;42:121–136.
24. Palmer JB, Rudin NJ, Lara G, Crompton AW. Coordination of mastication and swallowing. *Dysphagia* 1992;7:187–200. [PubMed: 1308667]
25. Sazonov E, Schuckers S, Lopez-Meyer P, Makeyev O, Sazonova N, Melanson EL, Neuman M. Non-invasive monitoring of chewing and swallowing for objective quantification of ingestive behavior. *Physiological Measurement* 2008;29:525–541. [PubMed: 18427161]
26. Turkoglu I, Arslan A, Ilkay E. An intelligent system for diagnosis of the heart valve diseases with wavelet packet neural networks. *Computers in Biology and Medicine* Jul;2003 33:319–31. [PubMed: 12791405]
27. Liu Y, Zhang C, Peng Y. Neural Classification of Lung Sounds Using Wavelet Packet Coefficients Energy. *PRICAI 2006: Trends in Artificial Intelligence* 2006:278–287.
28. Cristianini, N.; Shawe-Taylor, J. *An Introduction to Support Vector Machines and Other Kernel-based Learning Methods*. Cambridge University Press; 2000.
29. Addison, N. *The Illustrated Wavelet Transform Handbook*. Taylor & Francis; 2002.
30. Fu, S.; Liu, X.; Muralikrishnan, B.; Raja, J. Wavelet Analysis with Different Wavelet Bases for Engineering Surfaces. *Proceedings of the Sixteenth Annual Meeting of The American Society for Precision Engineering*; Raleigh, NC. 2001. p. 249-252.
31. Buckheit, JB.; Donoho, DL. *LECTURE NOTES IN STATISTICS. NEW YORK: SPRINGER VERLAG*; 1995. WaveLab and Reproducible Research; p. 55-55.
32. Moddemeijer R. On estimation of entropy and mutual information of continuous distributions. *Signal Processing* Mar;1989 16:233–248.
33. WU G, LIN C. Word boundary detection with mel-scale frequency bank in noisy environment. *IEEE transactions on speech and audio processing* 2000;8:541–554.
34. Meyer D, Leisch F, Hornik K. The support vector machine under test. *Neurocomputing* 2003;55:169–186.
35. Chang, C.; Lin, C. *LIBSVM : a library for support vector machines*. Software. 2001. available at: <http://www.csie.ntu.edu.tw/~cjlin/libsvm>

36. Vaiman M, Nahlieli O. Oral vs. pharyngeal dysphagia: surface electromyography randomized study. *BMC Ear, Nose and Throat Disorders* 2009;9:3.
37. De Belie N, Sivertsvik M, De Baerdemaeker J. Differences in chewing sounds of dry-crisp snacks by multivariate data analysis. *Journal of Sound and Vibration* Sep;2003 266:625–643.

## Biographies



**Edward Sazonov** (M'02) received the Diploma of Systems Engineer from Khabarovsk State University of Technology, Russia, in 1993 and the Ph.D. degree in Computer Engineering from West Virginia University, Morgantown, WV, in 2002.

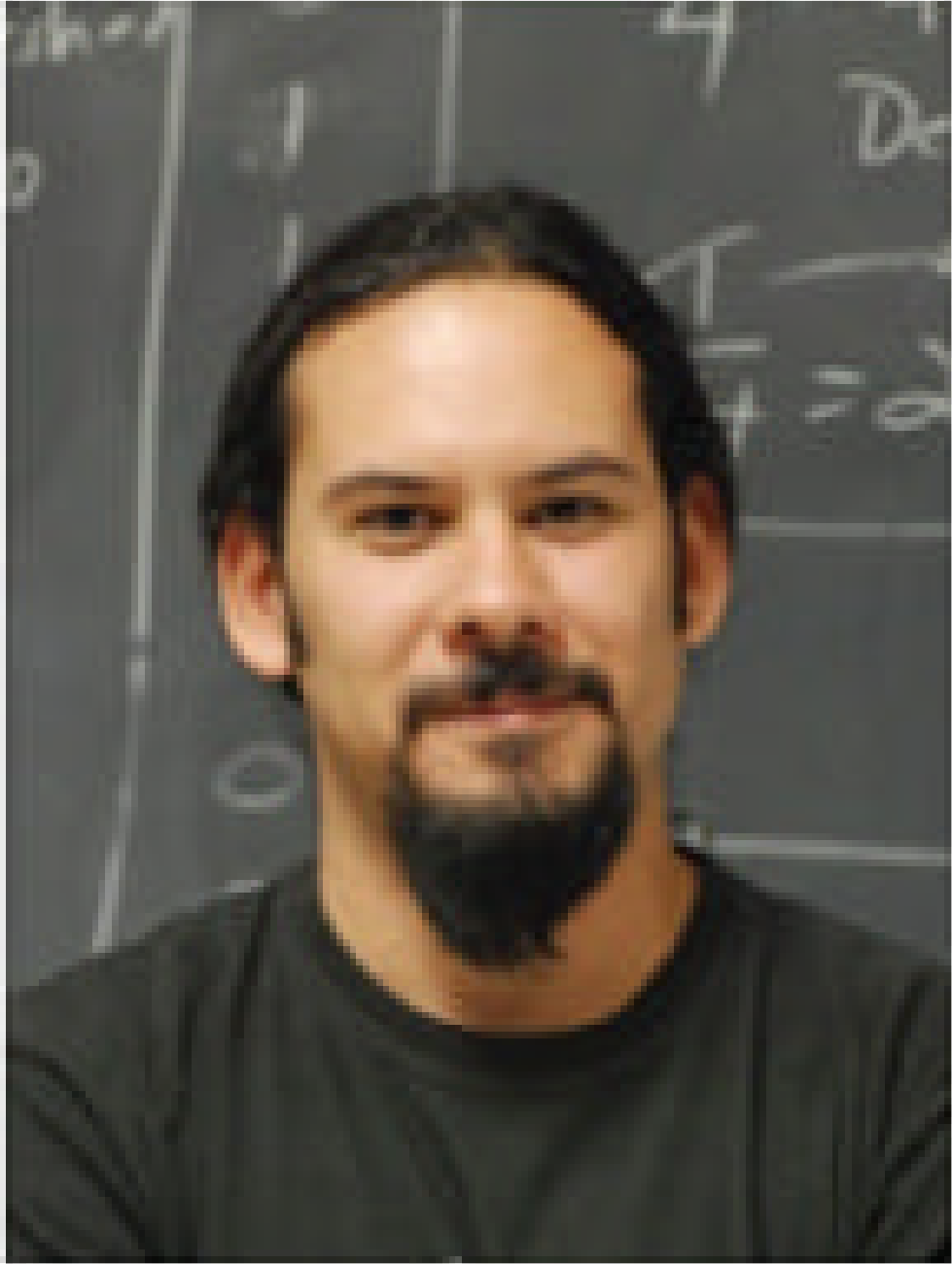
Currently he is an Associate Professor in the department of Electrical and Computer Engineering, Potsdam, NY and the head of the Clarkson Laboratory of Ambient and Wearable Devices. His research interests span bioengineering, computational intelligence, wireless, ambient and wearable devices. Applications include development of methods and wearable sensors for non-invasive monitoring of ingestion; methods and devices for monitoring of physical activity and energy expenditure; wearable platforms for rehabilitation of stroke patients and monitoring of the risk of falling in elderly; and, self-powered ambient sensors for structural health monitoring. His work has been supported by national (National Science Foundation, National Institutes of Health, National Academies of Science) and state agencies, and private industry.



**Oleksandr Makeyev** (M'03) received B.Sc. in mathematics and M.Sc. in statistics from Taras Shevchenko National University of Kyiv, Kiev, Ukraine, in 2003 and 2005 respectively. Currently he is working towards Ph.D. in engineering science at the Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY, USA. His broad research interests include development and application of computational intelligence based pattern recognition methods to engineering problems. His current research involves signal processing and pattern recognition for monitoring of ingestive behavior in humans.



**Stephanie Schuckers** is an associate professor in the Department of Electrical and Computer Engineering at Clarkson University. Schuckers received the B.S. in electrical engineering from the University of Iowa in 1992. As a Whitaker Foundation Graduate Fellow, she received the M.S. and Ph.D. degree in electrical engineering from the University of Michigan in 1994 and 1997, respectively. Her research focuses on processing and interpreting signals which arise from the human body. Signals include the electrocardiogram, biometric signals like fingerprints, respiration, and electroencephalograms. Methods involve classic signal processing, statistical techniques, pattern recognition, algorithm development and evaluation, data mining, and image processing. Much of her work involves analysis of real data collected from human, cadaver, and animal studies. Her work is funded from various sources, including National Science Foundation, American Heart Association, National Institute of Health, Department of Homeland Security, the Center for Identification Technology, and private industry, among others. She has over 30 journal publications, as well as many conference papers and book chapters.



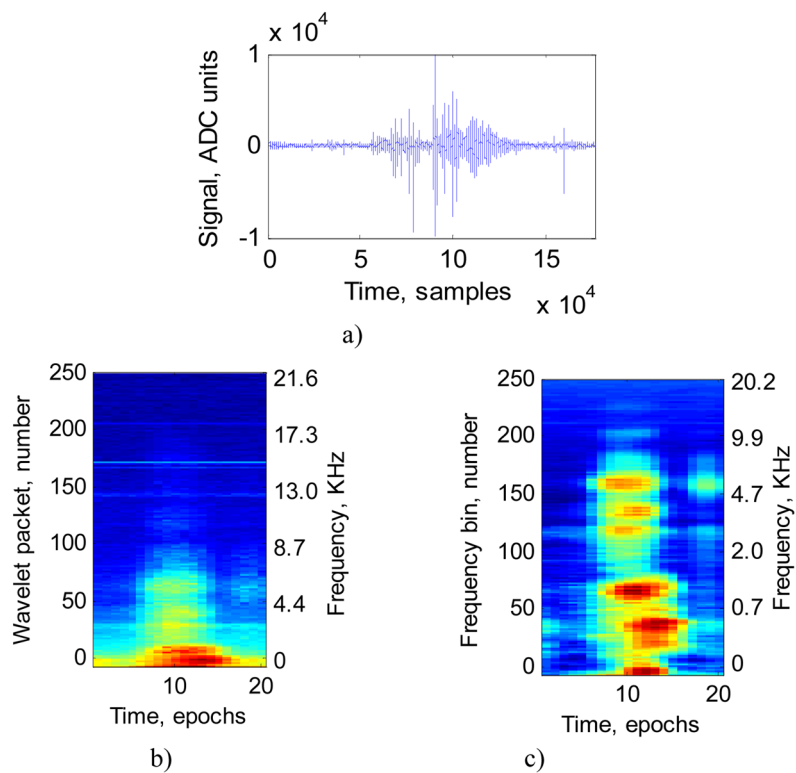
**Paulo Lopez-Meyer** received his Bachelors in Telecommunications Engineering in 2003, and his Masters in Instrumentation Engineering in 2005, both at the National Autonomous University of Mexico. He is currently a PhD candidate in Biomedical Engineering at Clarkson University. His research interests are the applications of Machine Learning and Pattern recognition in the solution of real life problems.



**Edward Melanson** began his research career as a post-doctoral fellow in the Center for Human Nutrition in 1998, and received a faculty appointment in the Division in 2003. His research interests are the effects of diet, exercise, and obesity on substrate metabolism and energy expenditure. Currently, he is performing studies on the effects of different intensities of exercise and manipulations in dietary fat on fat oxidation. These studies are performed using whole-room indirect calorimetry. He also is currently performing an exercise intervention aimed at determining the effects of exercise on the different components of energy expenditure, particularly non-resting energy expenditure. In these studies, a variety of approaches are used to assess energy expenditure including doubly labeled water, indirect calorimetry, and accelerometry. His research is funded by the NIH.

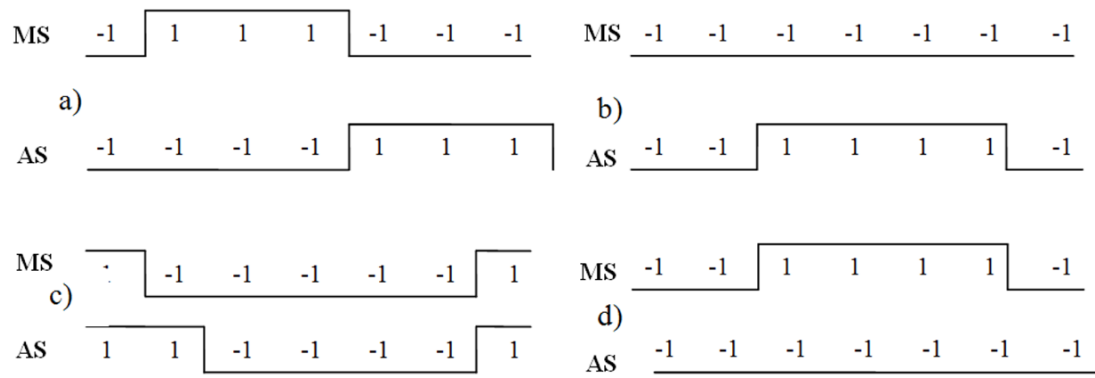


**Michael R. Neuman** joined the Department of Biomedical Engineering at Michigan Technological University in August, 2003, as Professor and Chair. He previously held the Herbert Herff Chair at the Memphis Joint Program in Biomedical Engineering and served for thirty-two years at Case Western Reserve University in the Departments of Biomedical and Electrical Engineering and the Departments of Reproductive Biology and Obstetrics and Gynecology. He received the BS, MS and PhD degrees in electrical engineering from Case Institute of Technology in Cleveland, Ohio and the MD degree from Case Western Reserve University. His primary research interests have been in the application of microelectronic technology to problems in clinical medicine. He has served as President of the International Society on Biotelemetry and was Editor in Chief of the *IEEE Transactions on Biomedical Engineering*, 1989–1996. He also served as Editor in Chief of the international journal, *Physiological Measurement*, 2002–2007, and he is currently Editor in Chief of the *IEEE Engineering in Medicine and Biology Magazine*. He lives with his wife, Judy, on a hobby farm just outside of Houghton in Michigan’s Upper Peninsula.

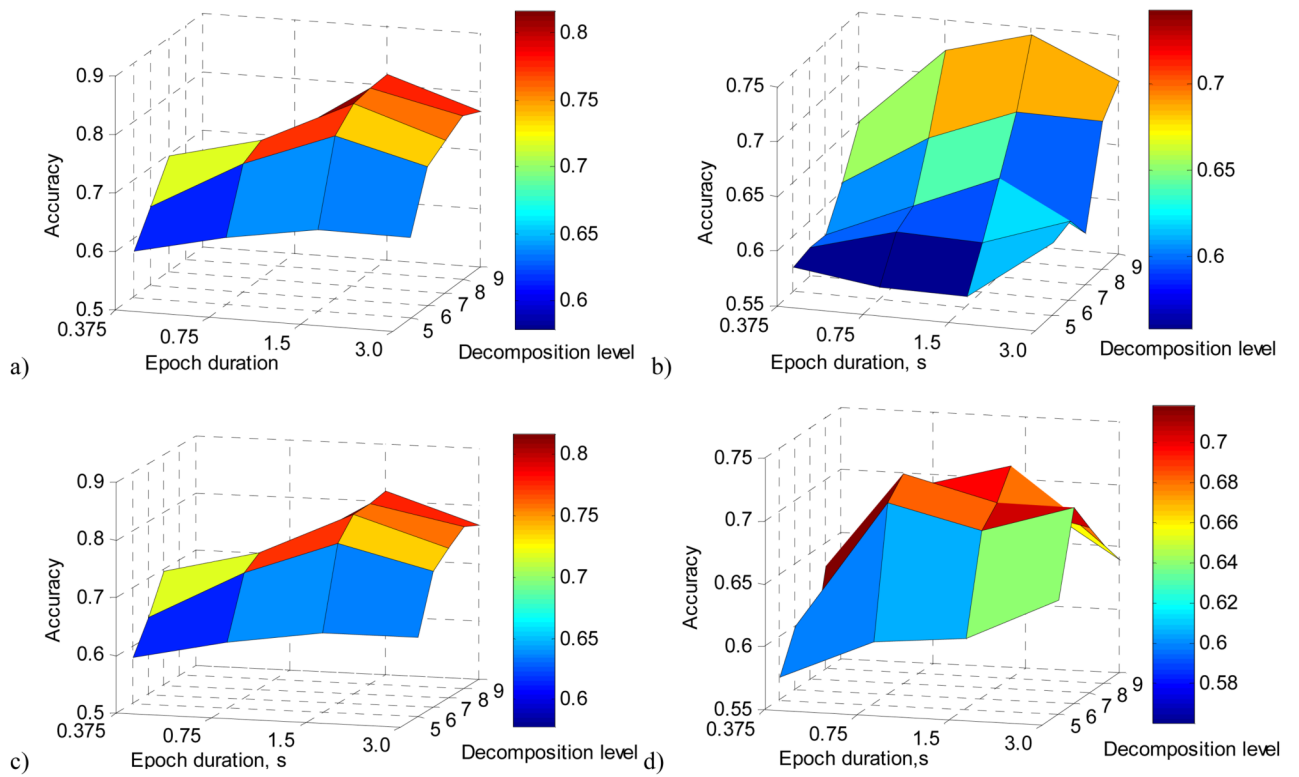


**Fig. 1.** a) a 4.0s fragment of a sound recording including a swallow; b) features extracted by WPD processing; c) features extracted by msFS processing. Frequencies are shown for the center of each packet or bin.





**Fig. 2.** Examples of: a) true positive, b) false positive, c) true negative, d) false negative. Each number represents a class label for an epoch ('-1' – non-swallow epoch, '1' – swallow epoch).



**Fig. 3.** Accuracy of swallowing sound recognition as a function of epoch duration and decomposition level a) msFS with no lags b) WPD with no lags c) msFS with  $K=1$  (3 lags) d) WPD with 3 lags.

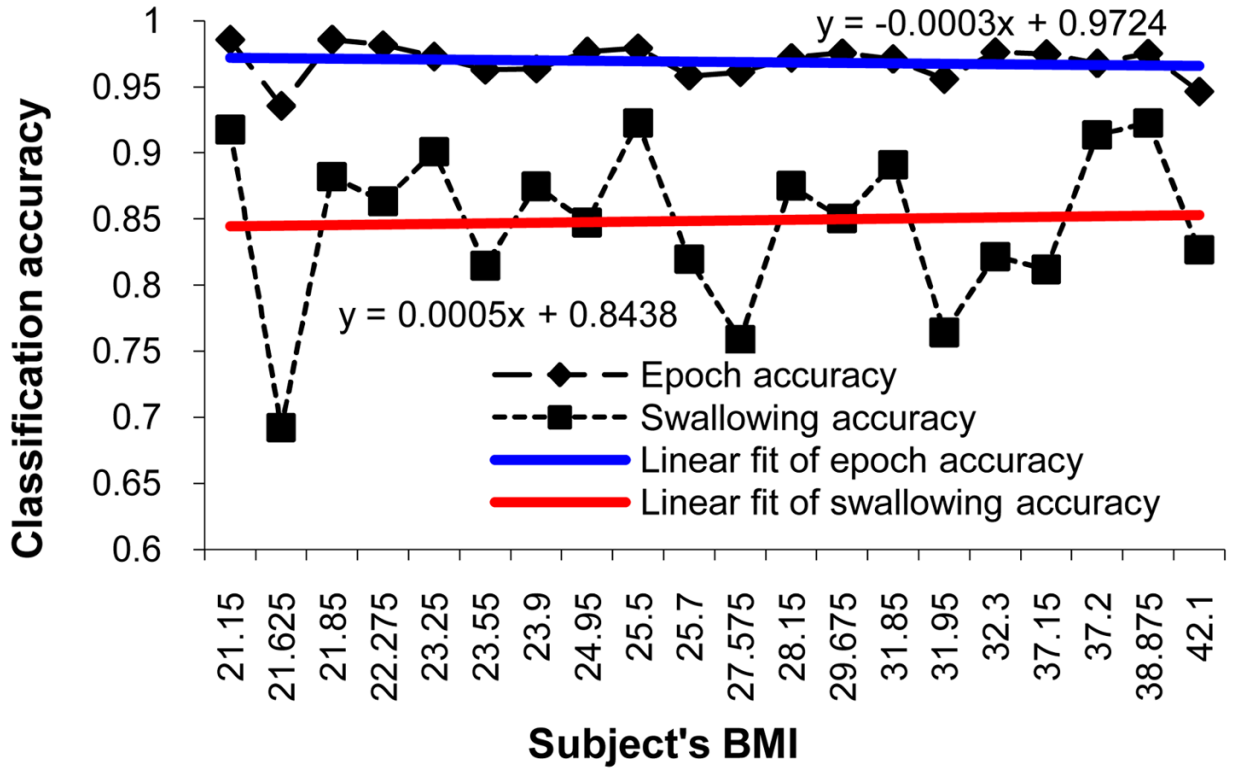


Fig. 4. Distribution of average weighted accuracy in classification of epochs and swallowing events versus subject's BMI and corresponding linear fit of the data.

**TABLE I**

Accuracy obtained in detection of swallowing events for threefold cross-validation

Feature	WPD 9	WPD 9	msFS 7	msFS 7
Number of lags	1	3	1	3
Average per-epoch accuracy (%)	95.9	96.4	96	96.8
Average per-swallow accuracy (%)	79.4	79.4	79	84.7