# Meal-Related Activities Classification from Bone-Conducted Signal

Archit Jain[2], Takumi Kondo[1], Anna Yokokubo[1] and Guillaume Lopez[1]

[1] Department of Integrated Information Technology,
Aoyama Gakuin University, Sagamihara, Japan

[2] University Jean Monnet, Saint-Etienne, France

July 31, 2019

## ABOUT LAB

At the Wearable Environment and Information System Laboratory (W.I.L.), we aim to support a healthier and more comfortable life, together with extending at most smartphones capabilities, we are doing research and development on next-generation multimedia devices that will turn both information and environment wearable. Our research is applied in mainly three fields: healthcare support, skill science, and daily-life augmentation through wearables. Current projects are 'Healthy Eating Habit Support System Using Wearable Sensors and IoT Devices', 'Integrated Thermal Comfort Control System', 'Mental Health Monitoring', and 'Skill Improvement Support using Wearables'.

**Abstract**

    **Increasing the number of mastication can help suppress obesity, but it is difficult to keep constant awareness of it in everyday life. Besides, the conventional mastication number measurement apparatus is large and non-portable such that it is difficult to use it in daily life. This research proposes a system that may support both consciousness improvement feedback in real-time and accurate quantified monitoring of mastication, swallowing and such other meal-related activities. It is composed of a cheap and small bone-conducted microphone to collect intra-body sound signal, bone a smartphone that can process it provides feedback. Though meal-related activities quantification using wearable devices has been studied for many years, proposed systems have only been tasted in experimental environment. In this study, we collected dietary sound data in a natural eating environment, evaluated the performance of different models, and proposed an optimization for classification of chewing swallowing and speaking activities.**

# I. INTRODUCTION

    Obesity may cause lifestyle diseases such as diabetes and heart disease. The Japanese Ministry of Health, Labor and Welfare has taken measures for preventing these diseases, but the number of obese patients has not decreased compared to 10 years ago [1]. As measures against obesity, it is known to be useful to improve eating habits and exercise moderately, but it is also possible to significantly prevent it by increasing the number of mastications [2, 3]. As a concrete example, when attempting to improve mastication activity for young Chinese men with obesity, it was possible to reduce the intake of energy in all the subjects consistently [4]. The same study also demonstrated that there is a useful possibility of measures against obesity by the activity of increasing the number of chewing.

    Improvement in the mastication amount is also crucial since healthcare experts always check the number of chewing as well as meal duration and food type as an indispensable factor in assessing dietary habits. In addition to the above, to prevent obesity, the nervous system and chewing activities are closely related. This relation is because chewing activities are closely related. This relation is because chewing repetition stimulates the satiety centre and sympathetic nervous system, which can reduce obesity by secreting hormones that suppress appetite [5]. Notably, several past works have reported that people with fast-eating have higher tendency to be obese,

which is partly because lowering secretion of hormones by eating fast causes an increase in dietary amount [6,7]. In addition to this, it is considered desirable to encourage utterances during meals. Indeed, Kishida et al. have reported that making conversation during meals is related to good health [8].

Though chewing and swallowing processes depend on many factors both human and food property dependants [9], recent research suggest that self-quantification is strongly associated with the will to optimise or improve own's performance or behaviour [10,11]. Our research aims at proposing a system that may support both consciousness improvement feedback in real-time and accurate quantified monitoring of meal-related activities. It is composed of a cheap and small bone conduction microphone to collect sound intra-body sound signal, and a smartphone that can process sound and provides feedback in real-time so that it can be used conveniently in daily life.

## II. STATE OF THE ART

A decade ago, studies started to focus on chewing as an improvement of dietary habits, mainly proposing various methods and devices to quantify mastication activity with little burden [12, 13]. They proposed to use mainly devices that measure myoelectric potential from the masseter muscle can count bite. Another technique using infrared sensor can detect small changes in temporal muscle tension, but one can consider that this method is not applicable in the sense that it bothers users during meals due to the sensing medium and the appearance ([14]). Tanigawa et al. explored the use of the Doppler effect in their system to sense the Doppler signal of mastication produced from vertical jaw movements [15]. However, some calibration is required. Recently Keum et al. proposed a novel multimodal sensing strategy combining accelerometer and range sensing, implemented into a discreet and lightweight instrumented necklace that captures head and jawbone movements without direct contact with the skin [16]. Though eating episodes could be detected with high precision and recall (95.2% and 81.9%) in controlled conditions, the performances drop considerably in free-living conditions (78.2% precision and 72.5% recall). All these methods also have the inconvenience that wearing the apparatus in daily life is a significant burden for the user.

On the other hand, analysis of internal body sounds spectra has attracted attention as a way to differentiate between biting and speaking activities, and to classify several types of food with less

burden [13, 17, 18]. Uno et al. proposed a system to detect the chewing frequency and bite fidelity using bone conduction microphones [19]. Paying attention to the amplitude during chewing, it is a system that judges chewing when amplitude magnitude exceeds a certain level, and the judgment accuracy was about 89%. However, activity discrimination method is limited to specific ailments and robustness to other sounds not evaluated. Similarly, Nishimura et al.[20] and Faudot et al. [21] proposed to measure the chewing frequency using a wireless and wearable in-ear microphone. However, to estimate the number of chewing operations, still, some parameters need to be adjusted by the user each time, which is a severe constraint in practical use. Recently, Bi et al. developed a wearable device that can automatically recognize eating behavior in free-living conditions using an off-the-shelf contact microphone placed behind the ear [22]. Though they achieved accuracy exceeding 92.8% and F1 score exceeding 77.5% for eating event detection, the accuracy of detailed meal-related activities such as number of mastications or swallowing are not assessed.

As summed-up above, eating activity monitoring systems using wearable devices have been proposed Despite numerous efforts by researchers over the last decade, an objective and usable method for detailed tracking of dietary intake behavior in natural meal environment remains unrealized.

# III. PROPOSED SYSTEM

*A. System Outline*

Figure 1 shows the outline of proposed system. Based on the above review of related research, it has been decided to use a bone-conduction microphone to enable the collection of both chewing, swallowing, and utterance activities information. Some of the hands-free headsets available on the market are integrating such specific microphone, making it easily accessible to everyone. Moreover, Fontana et al. have shown earlier that even a strain sensor to detect chewing events and a throat microphone to detect swallowing sounds present enough comfort levels, such the presence of the sensors does not affect the meal [23]. A smartphone is used to deal with the real-time processing of the collected sound signal. Processing task may be whether a simple data transfer to some online computation and storage resource for detailed monitoring purpose, or onboard data analysis and feedback for behavior awareness improvement [24].
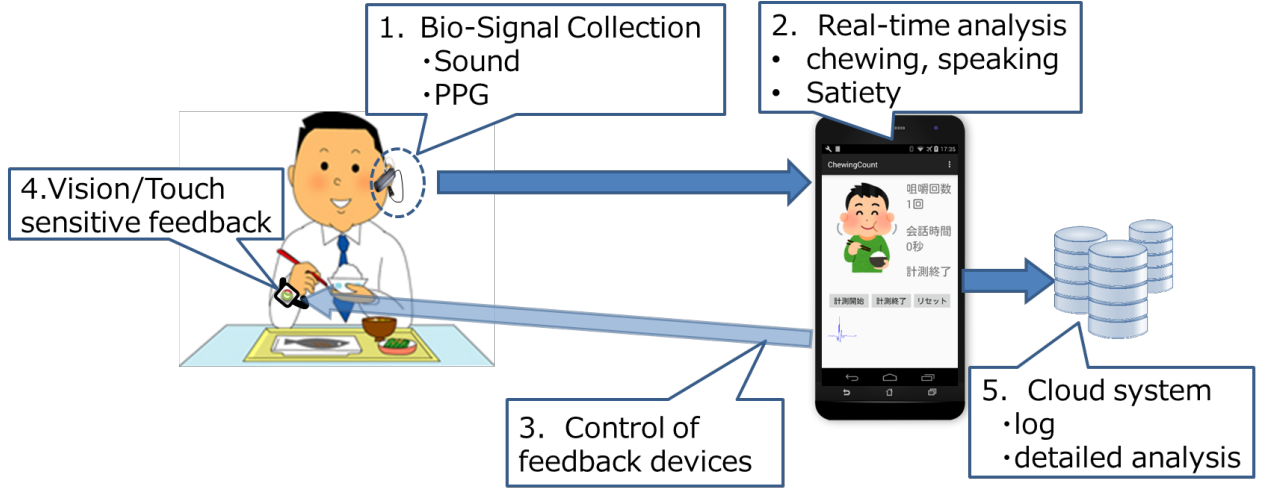
Fig 1. Outline schematic of proposed eating habits monitoring and support system

In this study, we have collected dietary sound data in a natural eating environment. Each sound data sequence corresponding to a single chew, swallow, and speaking events have been labeled. The obtained dataset has been used to evaluate the performance of different models for chewing, swallowing, and speaking activities classification depending on the amount of features and the type of classifier. Finally, the most efficient model has been optimized.

## IV. COLLECTIONS OF DAILY MEAL SOUND

### A. *Experimental conditions*

To discriminate mastication, swallowing, and utterance from sound collected in a natural eating environment, eating sound data collection was carried out in a completely free meal environment, that is, not in a laboratory environment. For example, some data were collected in a dining room and a standard household table with other family members, or at the university cafeteria with friends, such we can assume that represents different noisy conditions. The meal content was also totally free, and subjects ate whatever they wanted as usual in daily life, such various food types were mixed unpredictably during the same meal.

To collect dietary sound data, we used a commercial bone conduction microphone (Motorola Finiti HZ800 Bluetooth Headset, Motorola co. Ltd.), attached to one ear of the subject, that can operate Bluetooth communication with a smartphone (Google Pixel 3, Google co. Ltd.) and collected dietary voice data using a dedicated Android OS application. The sound signal sampling from the microphone was 8KHz. After collection, data were transferred to a computer for la-

belling and analysis. Besides, since data were collected in a totally free environment, it was necessary to perform labelling afterwards.

To label sound segments corresponding respectively to chewing, swallowing, and utterance after collecting the data, a video was taken together with sound data to assist the labelling work. The video shooting was performed so that the mouth and throat of the subject were reflected. Figure 2 shows a picture of the data collection conditions (for privacy, it is a photograph that reproduces the actual environment).



Fig 2. Picture reproducing the data collection conduitions

*B. Collected data labelling*

From the collected data it was necessary to extract the activity values associated with each data set. For that purpose the collected sound data were labeled in sections corresponding to one of the targeted three activities. To obtain the best labelling accuracy as possible, recorded video were synchronized with audio data and both were used as references. Labelling of audio data was done using "Praat", which is a software frequently used for speech analysis [25]. Labels were set to: "chew", "talk", "drink" and "swallow". Labelling in praat is shown in Fig3
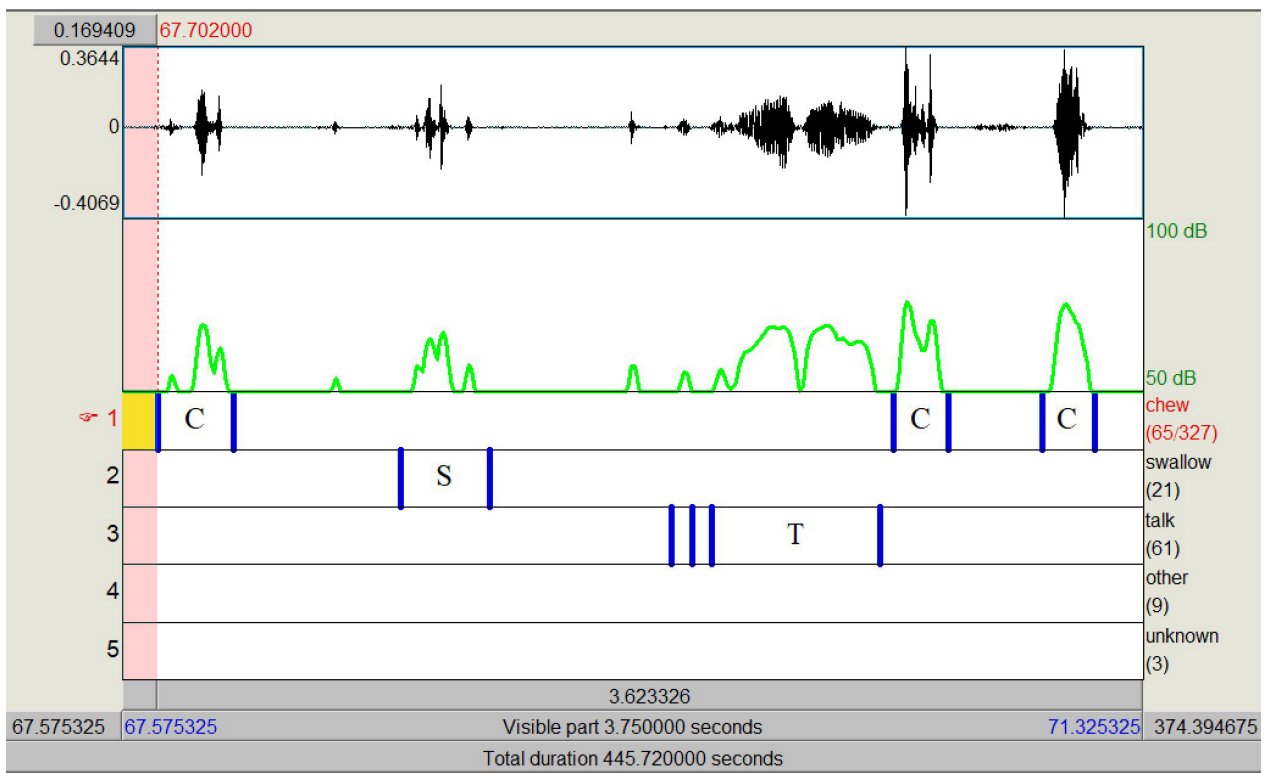
Fig 3: Labelling on "Praat"

Following the above described procedures we could prepare a dataset that details are described in Table 1. Data were collected from Japanese men and women (three each) from 11 to 23 years old. Data was collected from ten different meals. Though the number of subjects looks too small, the dataset represents 79 minutes of eating sound from unconstrained meal of various types of food, resulting in 1706 chewing samples, 99 swallowing samples, 29 drinking samples and 424 utterance (talk) samples.  Such, we consider the dataset is sufficient for subject independent activities classification. We also see an imbalance in dataset, which we'll look into after

TABLE 1: DETAIL OF AMOUNT OF LABEL RESPECTIVE TO THE EACH MEAL

| Meal time (Minutes:sec) | No of Chew | No of Drink | No of Swallow | No of Talk |
|---|---|---|---|---|
| 18:12 | 156 | 5 | 14 | 73 |
| 7:26 | 181 | 0 | 10 | 39 |
| 12:07 | 119 | 6 | 8 | 10 |
| 8:21 | 206 | 0 | 8 | 65 |
| 5:43 | 500 | 14 | 14 | 63 |
| 3:17 | 110 | 0 | 7 | 44 |
| 10:07 | 156 | 0 | 21 | 50 |

7

| Meal time (Minutes:sec) | No of Chew | No of Drink | No of Swallow | No of Talk |
|---|---|---|---|---|
| 4:28 | 83 | 0 | 6 | 0 |
| 5:33 | 110 | 4 | 7 | 70 |
| 3:53 | 85 | 0 | 4 | 10 |
| **79:07** | **1706** | **29** | **99** | **424** |

# V. MEAL-TIME ACTIVITIES CLASSIFICATION

*A. Feature extraction*

Features extraction has been performed from the dataset labeled according to the previous section before operating machine learning models for meal-related activities classification. A total of 26 features were extracted to aim at obtaining high classification accuracy. Table 2 sums-up an outline description of extracted features.

Table 2: Feature and the respective number extracted

| Sr no | Feature | No. of features |
|---|---|---|
| 1 | Mean of Chroma vector | 1 |
| 2 | Root mean square energy | 1 |
| 3 | Spectral centroid | 1 |
| 4 | Spectral bandwidth | 1 |
| 5 | Spectral roll off | 1 |
| 6 | Zero crossing rate | 1 |
| 7 | Mel frequency cepstral coefficients (MFCC) | 20 |

*1. Chroma Vector*

A **chroma vector** is a typically a 12-element feature vector indicating how much energy of each pitch class {C, C#, D, D#, E, F, F#, G, G#, A, A#, B} is present in the signal. One main property of chroma features is that they capture harmonic and melodic characteristics of music, while being robust to changes in timbre and instrumentation. It is also used for audio-matching making it a useful feature.

*2. Root mean square energy (RMSE)*

The **energy** of a signal corresponds to the total magnitude of the signal. For audio signals, that roughly corresponds to how loud the signal is. The energy in a signal is defined as

$$\sum_n \left| x(n) \right|^2$$

The root-mean-square energy (RMSE) in a signal is defined as

$$\sqrt{\frac{1}{N} \sum_n \left| x(n) \right|^2}$$

*3. Spectral centroid*

The **spectral centroid** indicates at which frequency the energy of a spectrum is centered upon, or in other terms it indicates where the centre of mass of the spectrum is located. This is like a weighted mean

$$f_c = \frac{\sum_k S(k) f(k)}{\sum_k S(k)}$$

Where *S(k)* is the spectral magnitude at frequency bin k, *f(k)* is the frequency at bin k.

*4. Spectral bandwidth*

computes the order-p spectral bandwidth

$$\left( \sum_k S(k) \left( f(k) - f_c \right)^p \right)^{\frac{1}{p}}$$

9

where $S(k)$ is the spectral magnitude at frequency bin $k$, $f(k)$ is the frequency at bin $k$, and $f_c$ is the spectral centroid. When $p = 2$, this is like a weighted standard deviation.
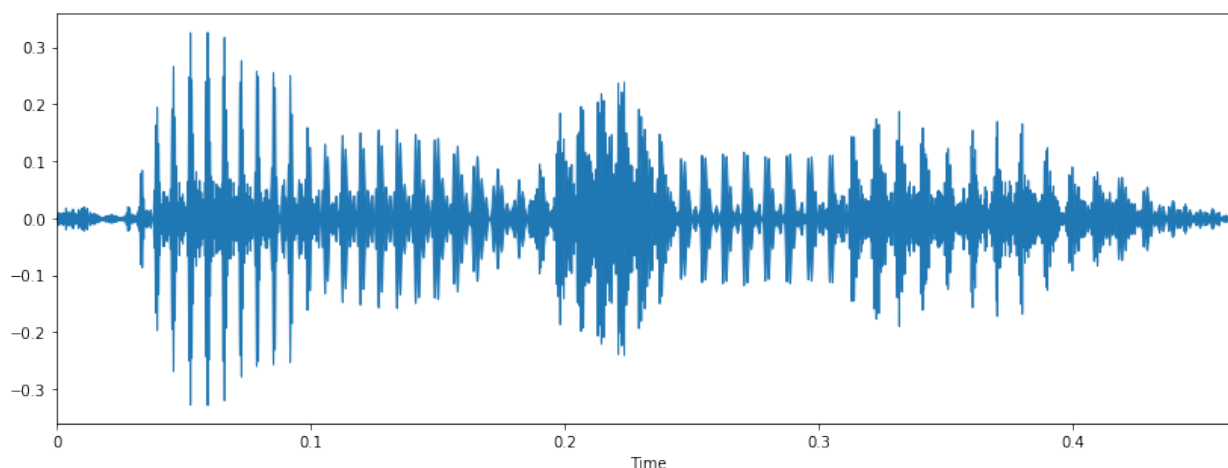
## 5. *Spectral roll off*

**Spectral rolloff** is the frequency below which a specified percentage of the total spectral energy, e.g. 85%, lies.
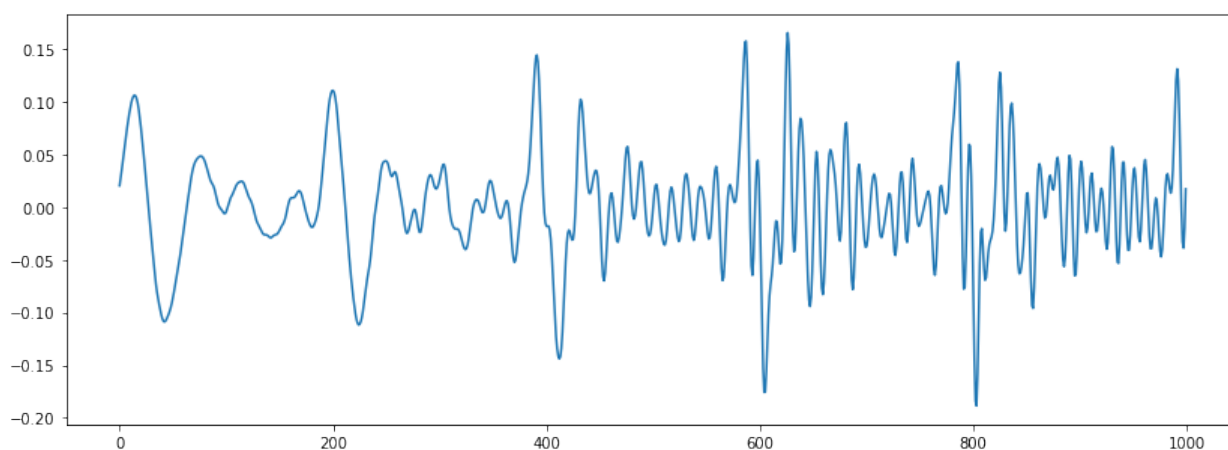
## 6. *Zero crossing rate*

**Zero crossing rate** indicates the number of times that a signal crosses the horizontal axis. This is a key feature for percussive sounds and hence used for distinguishing whether human voice is present in audio or not. Lets see if it works or not

6.1 First of all lets import a sound file of **talk.** With 'time (sec)' on x axis and amplitude on y
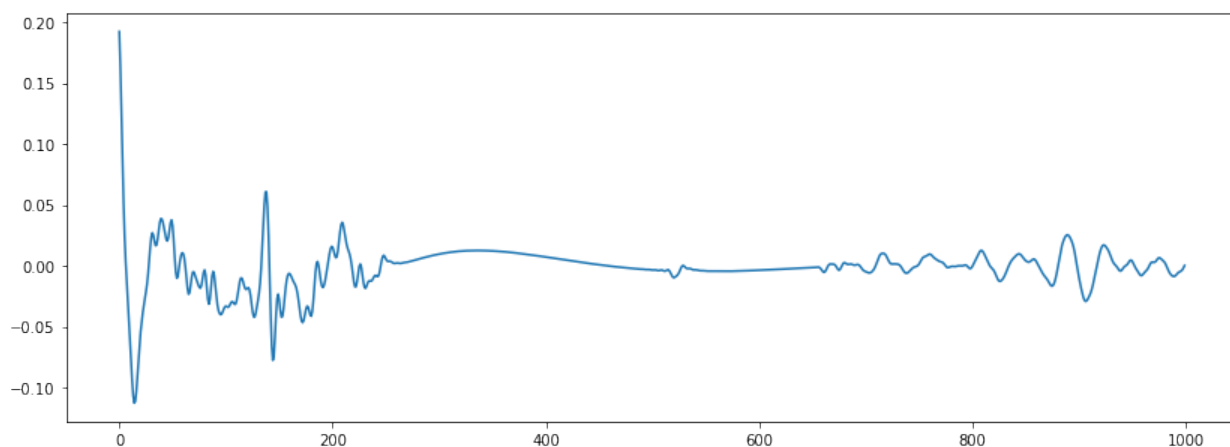


6.2 Zoom further to see crossing rate closely.

6.3 Now lets import some other file like **swallow,** time (sec) on x-axis and amplitude on y.
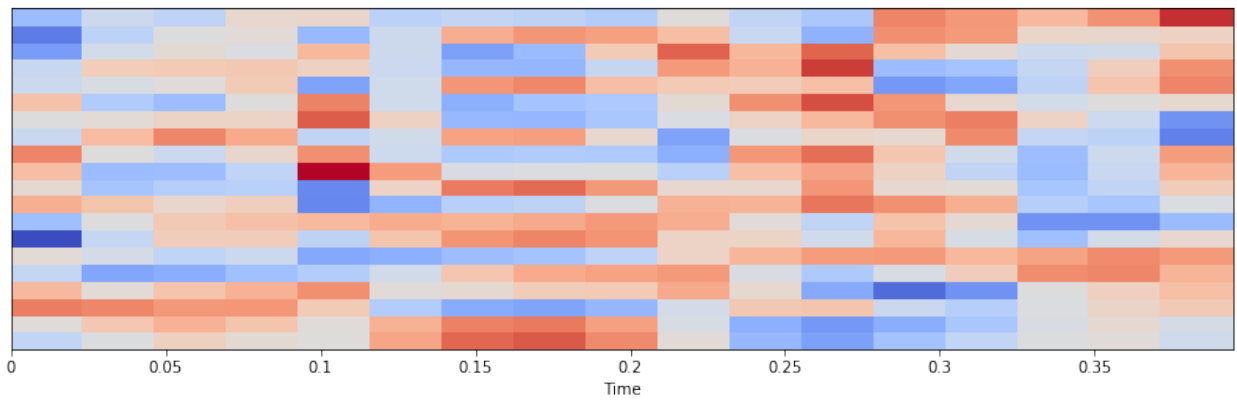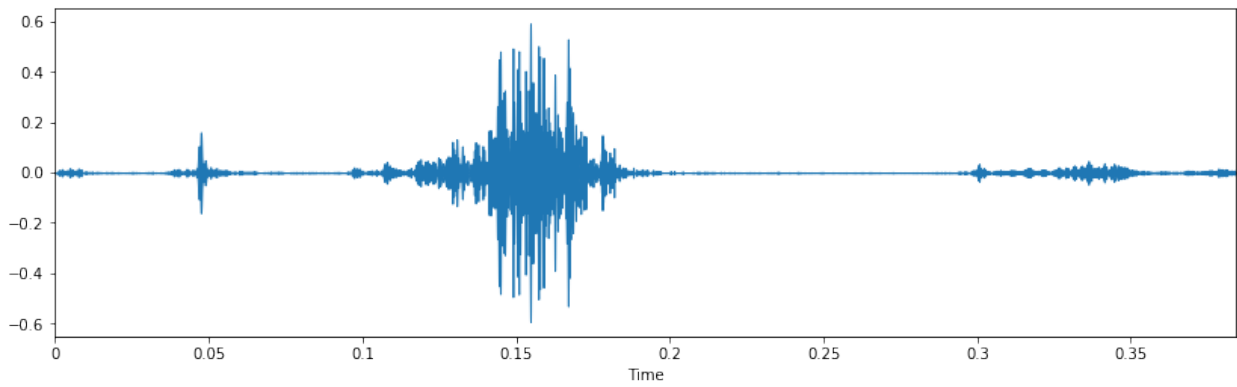


6.4 Further zoom to see the zero crossings



6.5 So it can be concluded that ZCR is a strong feature in recognising human voice (or more per‑ cussion) in an audio, as its obvious from the plots the talk has more crossing rates
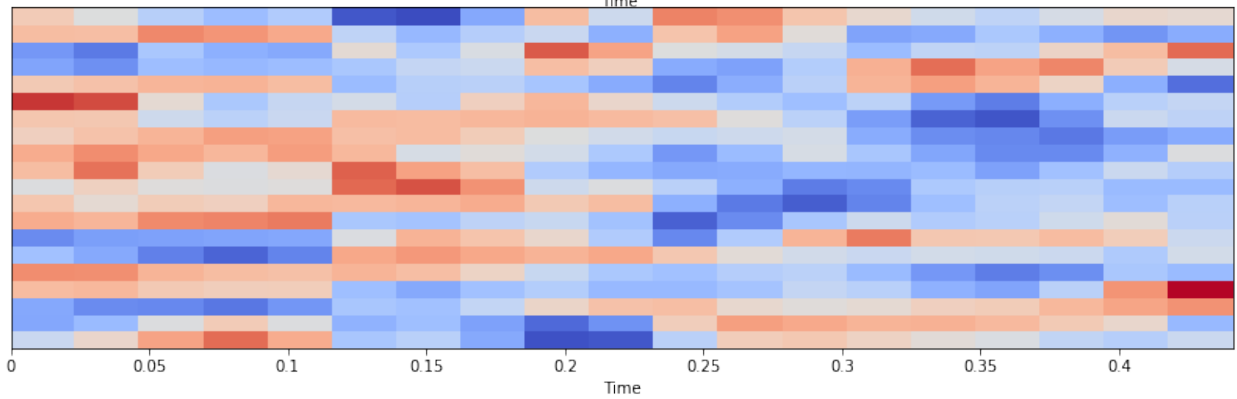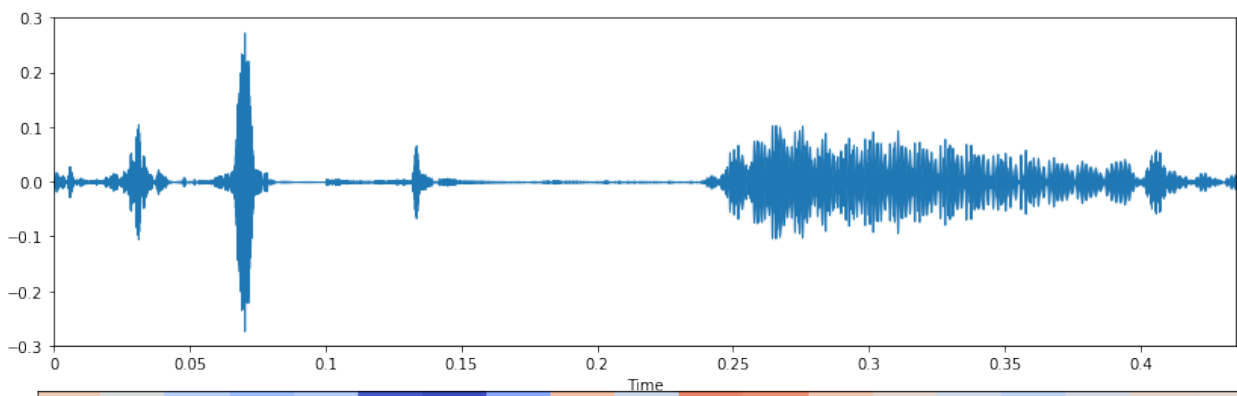
*7. Mel frequency cepstral coefficients (MFCCs)*

The most famous and used features for speech recognition are **MFCCs,** they are even used in speech recognition, and also because of them being powerful they are even used in CNN as pic‑ tures, for classifying further. Here, 20 order of MFCC are used in for features [27]
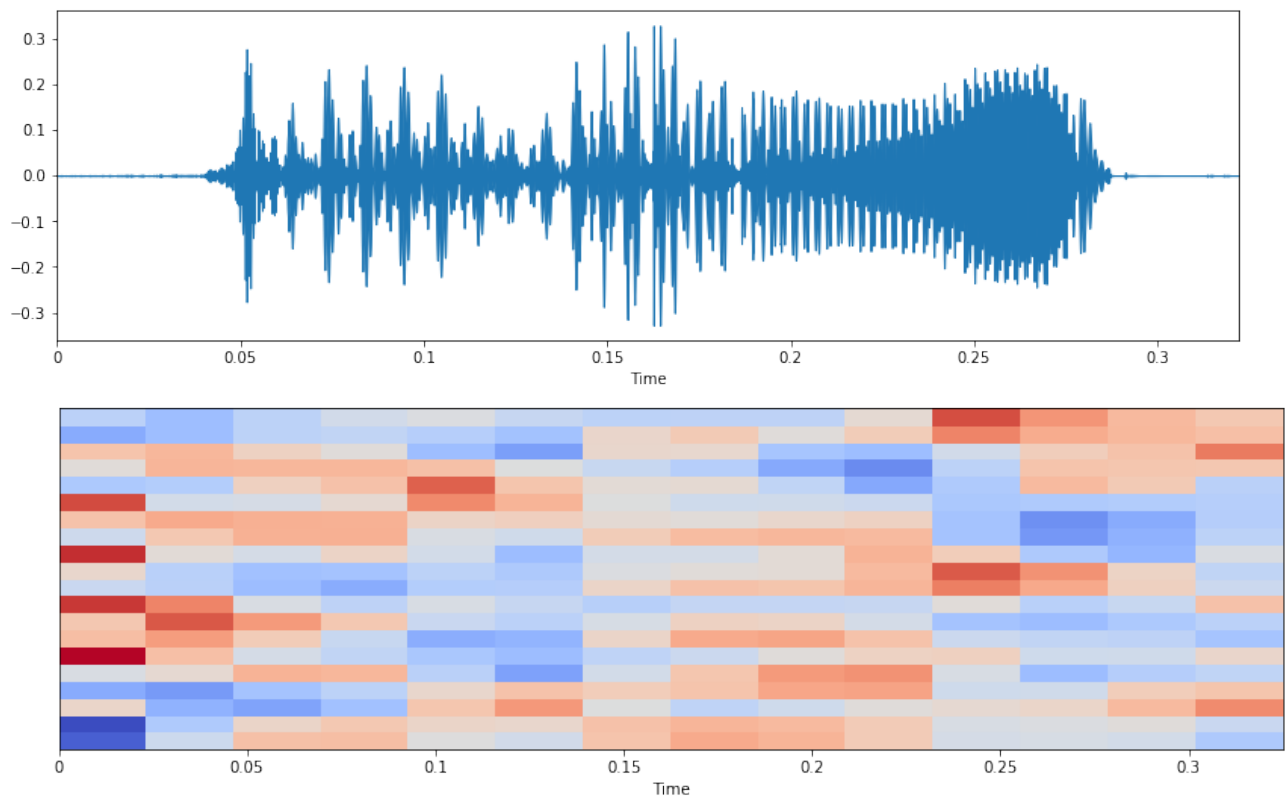
11

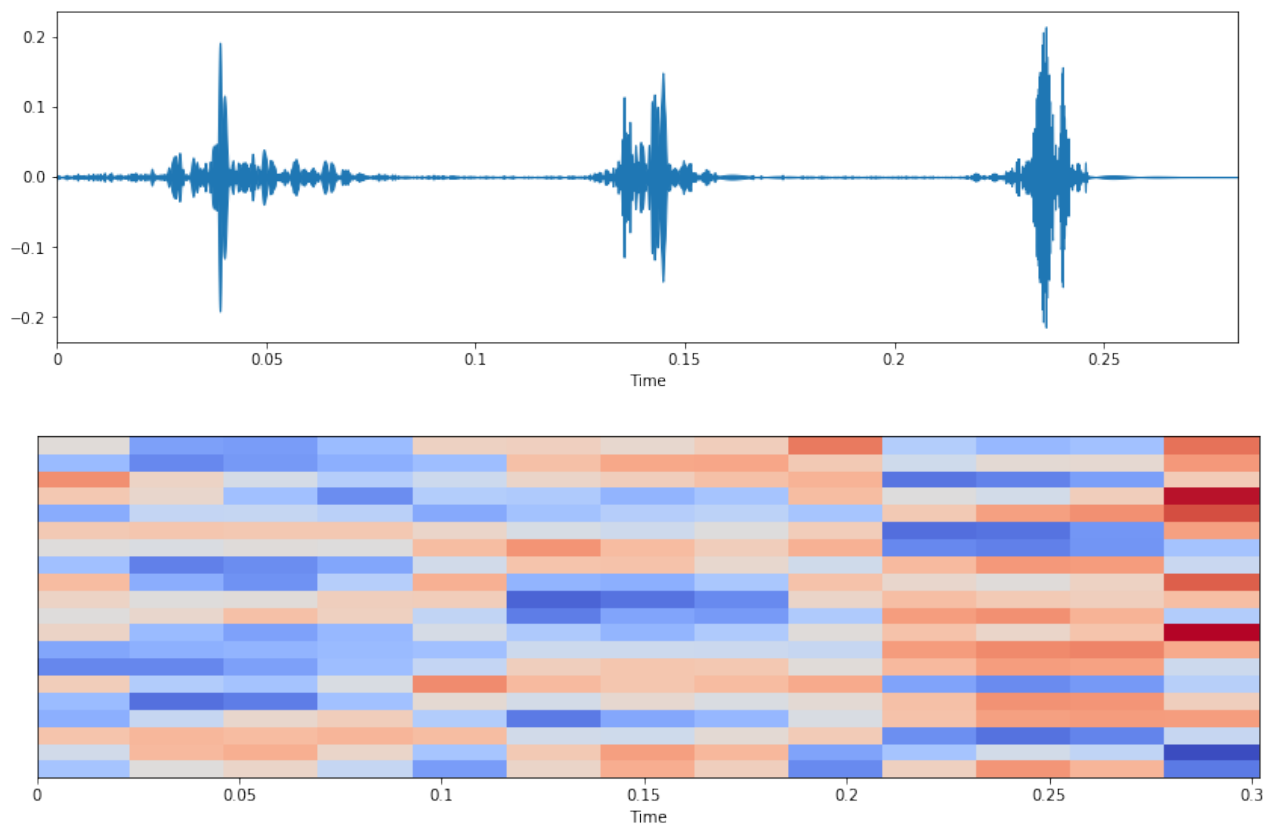## 7.1 For **chew**, applying the MFCC and then further scaling it

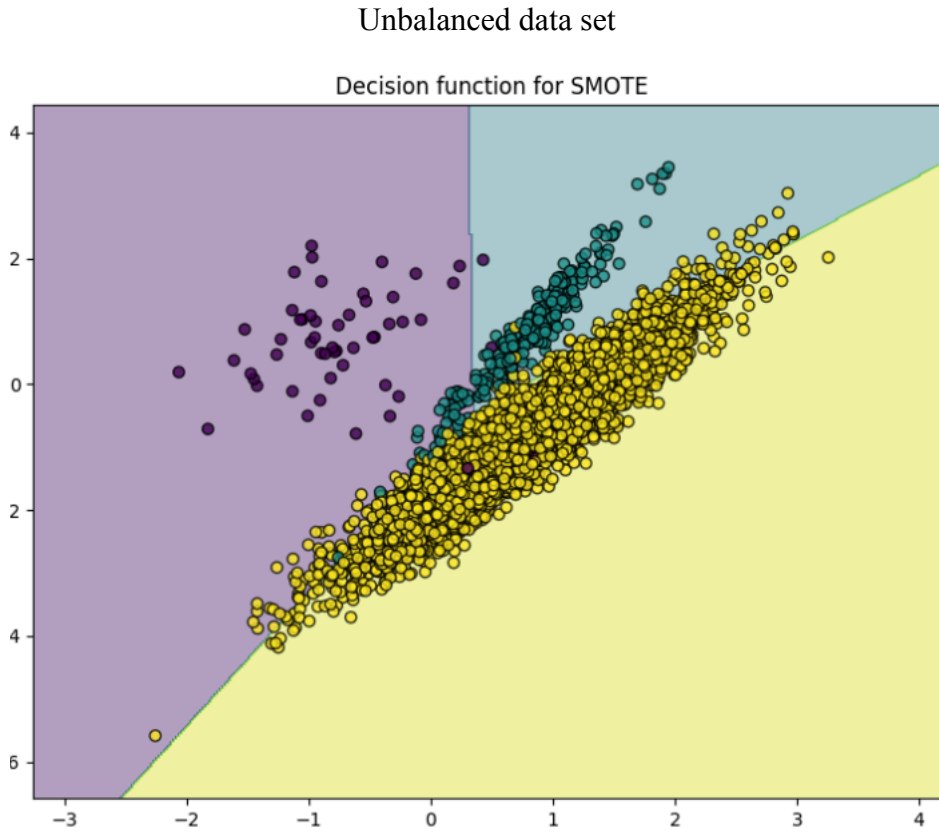

## 7.2 For **swallow**

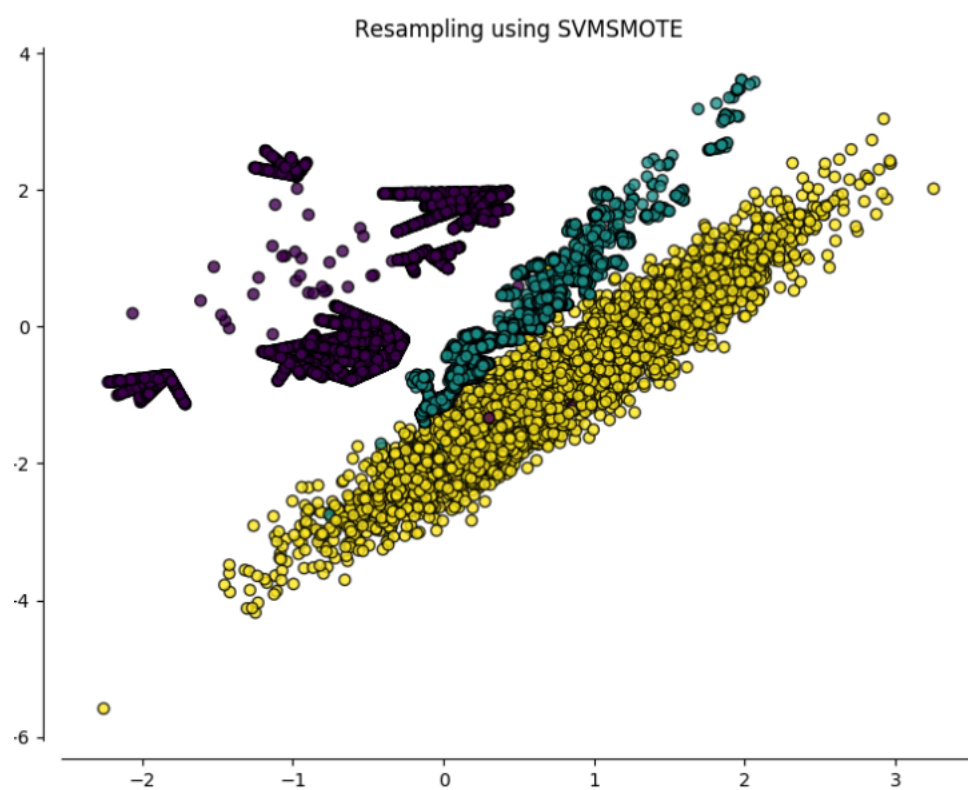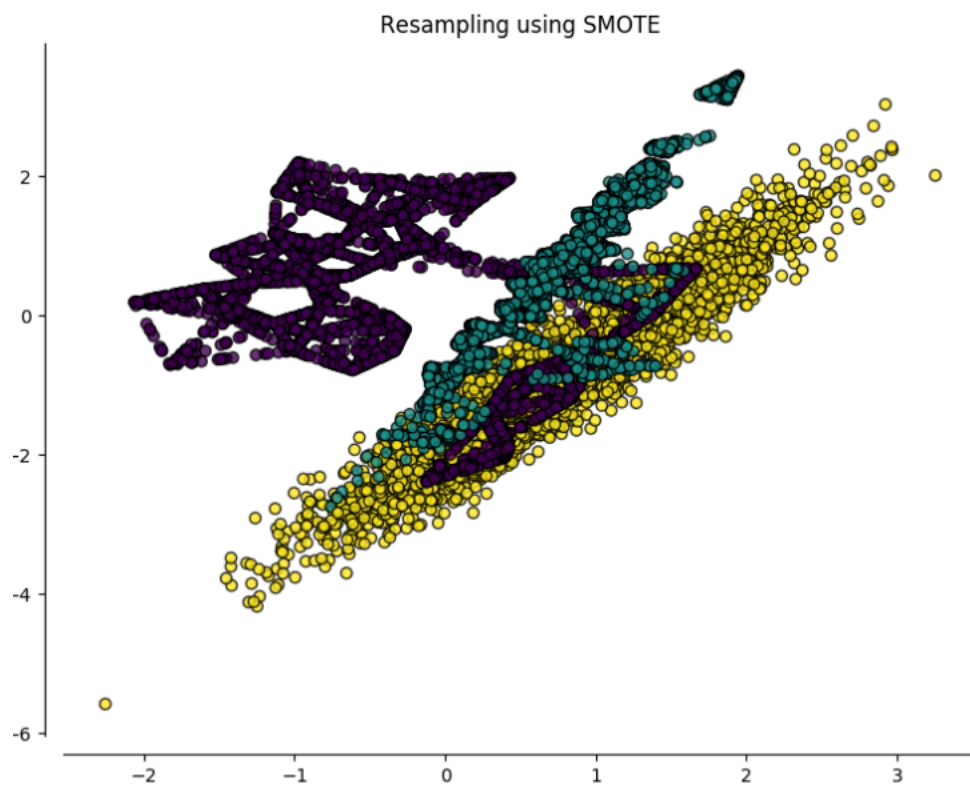## 7.3 For *talk*





## 7.4 For *drink*

*B. Balancing unbalanced datasets*

As previously show in Table 1, the obtained dataset have a huge number of "chewing" labeled data, and an only minimal number of other labels data to compare. Therefore, to equilibrate this unbalanced dataset, we applied SMOTE (Synthetic Minority Oversampling Technique), a library often used in imbalanced-learning tasks [28]. Simple replication of less representative classes samples would result in over-training and so over-learning. SMOTE, by creating a new specimen using neighbourhood interpolated data, enables to avoid this bias.

However, even after applying the conventional SMOTE, the results were still over-learning, because of interpolation technique used by traditional SMOTE. So, further ahead we tried different SMOTE technique, SVM-SMOTE, as it was using the SVM technique (maximizing the margins) for interpolating the data. The entire scenario can be depicted below. [The images are given for the illustrative purpose and do not come from original data]

Unbalanced data set



14

Resampling using SMOTE


Resampling using SVMSMOTE

We obtained a dataset with the different number of samples for each activity label as shown in Table 3 using SVM-SMOTE. The obtained balanced dataset was divided randomly for each label into training data (80%) and test data (20%).

Table 3: after applying SVM-SMOTE on the data

| Label name | Total | Train | Test |
|------------|-------|-------|------|
| chew | 1680 | 1344 | 336 |
| drink | 1090 | 872 | 218 |
| swallow | 1705 | 1364 | 341 |
| talk | 1675 | 1340 | 335 |

*C. Classifier Selection*

Finally, we performed classification of meal related activities using various supervised learning. Automatic machine learning and validation using a five-fold cross-validation method was performed to evaluate the average performance of the classifiers. A total of 5 classification models have been built based on the following four classifiers: Decision Tree, Support Vector Machine (SVM)and Nearest Neighbours classifier (KNN). The average accuracy of each output model is shown in Table 4. Though Fine tree and linear SVM were able to achieve accuracy of 85.1% and 86.2%, the best result was obtained for the medium Gaussian SVM (rbf kernel) with 96.7%) average accuracy.

Table 4 : Classification accuracies by different models (after tuning them)

| Classification models | | Accuracy (%) |
|-----------------------|---|--------------|
| Decision tree | Fine tree | 85.1 |
| | Coarse tree | 71.4 |
| Support Vector Machine | Linear SVM | 86 |
| | Gaussian | 96.7 |
| Nearest neighbour | KNN | 82.2 |

# VI. OPTIMIZATION OF SELECTED MODEL & VALIDATION WITH TEST DATA

The purpose of the evaluation is to propose a classification method of chewing, swallowing, drinking and utterance (talk) in natural meal environment, from sound data collected by a bone-conduction microphone, that is not only highly accurate but also lightweight enough to have the possibility to run real-time a smartphone. Such, we performed both features selection and optimization of SVM parameters.

## A. Optimisation of parameters

We used "grid search" method to optimize SVM parameters "C" and "Gamma". Six different values were tested for each parameter using a five-folds cross-validation. The cross-validation accuracy result for each parameters values combination is shown in Figure 3. A very high accuracy of 96%, which is almost identical as when using all 26 features (97.6%), were obtained with the parameters being 10 for "C" and 0.1 for "Gamma".

Score from the rbf model

| Label | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| Chew | 0.96 | 0.90 | 0.93 | 336 |
| Swallow | 0.99 | 0.99 | 0.99 | 215 |
| Drink | 0.93 | 0.98 | 0.95 | 346 |
| Talk | 0.98 | 0.98 | 0.98 | 333 |
| Accuracy | | | | |
| Macro avg | 0.96 | 0.96 | 0.96 | 0.96 |
| Weighted egg | 0.96 | 0.96 | 0.96 | 0.96 |

## B. Feature reduction

Among the 26 features used, we assumed some features are highly related to the classification model performances. To reduce the number of features to an easier to handle amount, we reduced the features to just MFCCs and Zero crossing rate, and for the same the accuracy was more than 80%

# CONCLUSION AND FUTURE WORK

In this study, we proposed a classification method of chewing, swallowing, drinking, and speaking activities using bone conduction sound corresponding to natural diet environment. We classified chewing, swallowing, and speaking activities by SVM using Gaussian kernel. 26 features were extracted and reduced feature set after feature selection also investigated. Generalization performance of optimized model using only the top 15 features confirmed its high accuracy, since the precision, recall, and F1 value all exceeded 90% both at macro level and for each activity. Though most of the extracted 26 features have been chosen from the state-of-the-art.

As a prospect, we plan to use the proposed classification model to classify mastication, swallowing, and utterance in real time using bone conduction microphone and smartphone. In realizing this, it is necessary to design a system that automatically extracts sound data segments that can be considered to be whether chewing, swallowing, drinking or utterance in real-time. Besides, it is also necessary to add the other sounds such as noises in the model so that it is more robust to natural meal environment.

# REFERENCES

[1] MHLW, in *The National Health and Nutrition Survey in Japan, 2014*. Japan Ministry of Health Labor and Welfare, 2016, (in Japanese).

[2] Y. Ando, N. Hanada, and S. Yanagisawa, "Does eating slowly lead to prevent obesity? a literature review," *Health Science and Health Care*, vol. 8, no. 2, pp. 54–63, 2008.

[3] T. Nicklas, T. Baranowski, K. Cullen, and G. Berenson, "Eating patterns, dietary quality and obesity," *Journal of the American College of Nutrition*, vol. 20, no. 6, pp. 599–608, 2001.

[4] J. Li, N. Zhang, L. Hu, Z. Ki, R. Li, C. Li, and S. Wang, "Improvement in chewing activity reduces energy intake in one meal and modulates plasma gut hormone concentrations in obese and lean young chinese men," *American Journal of Clinical Nutrition*, vol. 94, no. 3, pp. 709– 716, 2011.

[5] Kao, "The effect of chewing well, tasting and eating - preventive measures against metabolic syndrome and obesity," in *Kao Health Care Report N.19*, 2007, pp. 4–5, (in Japanese).

[6] E. Denney-Wilson and K. J. Campbell, "Eating behaviour and obesity," *BMJ*, vol. 337, pp. 73–75, 2008.

[7] D. Gaul, W. Craighead, and M. Mahoney, "Relationship between eating rates and obesity," *Journal of Consulting and Clinical Psychology*, vol. 43, no. 2, pp. 123–125, 1975.

[8] N.KishidaandY.Kamimura,"Relationshipofconversationduringmeal and health and dietary life of school children," *The Japanese Journal of Nutrition and Dietetics*, vol. 51, no. 1, pp. 23–30, 1993.

[9] J. Logemann, "Critical factors in the oral control needed for chewing and swallowing," *Journal of texture studies*, vol. 45, no. 3, pp. 173–179, 2014.

[10] M. Ruckenstein and M. Pantzar, "Beyond the quantified self: Thematic exploration of a dataistic paradigm," *New Media & Society*, vol. 19, no. 3, pp. 401–418, 2017.

[11] G. Neff and D. Nafus, *The Self-Tracking*. MIT Press, 2016.

[12] K. Kohyama, L. Mioche, and P. Bourdio, "Influence of age and dental status on chewing behavior studied by emg recordings during consump- tion of various food samples," *Gerontology*, vol. 20, no. 1, pp. 15–23, 2003.

[13] O. Amft, M. Stäger, P. Lukowicz, and G. Tröster, "Analysis of chewing sounds for dietary monitoring," in *UbiComp 2005, 7th Int. Conf. on Ubiquitous Computing*, 2005, pp. 52–72.

[14] K. Obata, T. Saeki, and Y. Tadokoro, "No contact-type chewing number counting equipment using infrared sensor," *Transactions of the Society of Instrument and Control Engineers*, vol. 38, no. 9, pp. 747–752, 2002.

[15] S. Tanigawa, H. Nishihara, S. Kaneda, and H. Haga, "Detecting mas- tication by using microwave doppler sensor," in *PETRA 2008, 1st international conference on Pervasive Technologies Related to Assistive Environments, article 88, 7 pages*, 2008.

[16] S. C. Keum, B. Sarnab, and T. Edison, "Detecting eating episodes by tracking jawbone movements with a non-contact wearable sensor," *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 2, no. 1, pp. 1–21, 2018. [Online]. Available: http://doi.acm.org/10.1145/3191736

[17] M. Shuzo, S. Komori, T. Takashima, G. Lopez, S. Tatsuta, S. Yanag- imoto, S. Warisawa, J.-J. Delaunay, and I. Yamada, "Wearable eating habit sensing system using internal body sound," *Journal of Advanced Mechanical Design Systems and Manufacturing*, vol. 4, no. 1, pp. 158– 166, 2010.

[18] H. Zhang, G. Lopez, M. Shuzo, and I. Yamada, "Analysis of eating habits using sound in- formation from a bone-conduction sensor," in *e- Health 2011, Int. Conf. on e-Health*, 2011, pp. 18–27.

[19] S. Uno, R. Ariizumi, and S. Kaneda, "Advising the number of mas- tication by using bone-conduction microphone," in *The 24th Annual Conference of the Japanese Society for Artificial Intelligence*. Japanese Society for Artificial Intelligence, 2010, pp. 1–4, (in Japanese).

[20] J. Nishimura and T. Kuroda, "Eating habits monitoring using wireless wearable in-ear mi-crophone," in *3rd International Symposium on Wire- less Pervasive Computing*, 2008, pp. 376–381.

[21] T. Faudot, G. Lopez, and I. Yamada, "Information system for mapping wearable sensors into healthcare services: Application to dietary habits monitoring," in *WIVE 2010, 2nd International Workshop on Web Intell- gence and Virtual Enterprises*, 2010, pp. 1–9.

[22] S. Bi, T. Wang, N. Tobias, J. Nordrum, S. Wang, G. Halvorsen, S. Sen, R. Peterson, K. Odame, K. Caine *et al.*, "Auracle: Detecting eating episodes with an ear-mounted sensor," *Pro-ceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 2, no. 3, p. 92, 2018.

[23] J. Fontana and E. Sazonov, "Evaluation of chewing and swallowing sensors for monitoring ingestive behavior," *Sensor letters*, vol. 11, no. 3, pp. 560–565, 2013.

[24] G. Lopez, H. Mitsui, J. Ohara, and A. Yokokubo, "Effect of feedback medium for real-time awareness increase using wearable sensors," in *HEALTHINF 2019, The 12th International Con-ference on Health Informatics*, 2019.

[25] P. Boersma, "Praat, a system for doing phonetics by computer," *Glot International*, vol. 5, no. 9/10, pp. 341–345, 2001.

[26] H. Zhang, G. Lopez, M. Shuzo, J.-J. Delaunay, and I. Yamada, "Mastication counting method robust to food type and individual," in *HEALTHINF 2012, The 5th International Confer-ence on Health Informatics*, 2012.

[27] J. Ando, T. Saio, S. Kawsaki, M. Katagiri, D. Ikeda, H. Mineno, and M. Nishimura, "Con-versational and eating behavior recognition by leveraging throat sound," in *DICOMO 2017, Mul-timedia, Distributed, Cooperative, and Mobile Symposium of the Information Processing Society of Japan*, 2017, (in Japanese).

[28] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of artificial intel- ligence research*, vol. 16, pp. 321–357, 2002.