Mariusz SZWOCH
Gdansk University of Technology, Faculty of Electronics, Telecommunication
and Informatics

# ACQUISITION AND INDEXING OF RGB-D RECORDINGS FOR FACIAL EXPRESSIONS AND EMOTION RECOGNITION[1]

**Summary**. In this paper KinectRecorder comprehensive tool is described which provides for convenient and fast acquisition, indexing and storing of RGB-D video streams from Microsoft Kinect sensor. The application is especially useful as a supporting tool for creation of fully indexed databases of facial expressions and emotions that can be further used for learning and testing of emotion recognition algorithms for affect-aware applications. KinectRecorder was successfully exploited for creation of Facial Expression and Emotion Database (FEEDB) significantly reducing the time of the whole project consisting of data acquisition, indexing and validation. FEEDB has already been used as a learning and testing dataset for a few emotion recognition algorithms which proved utility of the database, and the KinectRecorder tool.

**Keywords**: RGB-D data acquisition and indexing, facial expression recognition, emotion recognition

## AKWIZYCJA ORAZ INDEKSACJA NAGRAŃ RGB-D DO ROZPOZNAWANIA MIMIKI I EMOCJI

**Streszczenie**. W pracy przedstawiono kompleksowe narzędzie, które pozwala na wygodną i szybką akwizycję, indeksowanie i przechowywanie nagrań strumieni RGB-D z czujnika Microsoft Kinect. Aplikacja jest szczególnie przydatna jako narzędzie wspierające tworzenie w pełni zaindeksowanych baz mimiki i emocji, które mogą być następnie wykorzystywane do nauki i testowania algorytmów rozpoznawania emocji użytkownika dla aplikacji je uwzględniających. KinectRecorder został z powodzeniem wykorzystany do utworzenia bazy mimiki i emocji FEEDB, znacznie skracając czas całego procesu, obejmującego akwizycję, indeksowanie i walidację nagrań. Baza FEEDB została już z powodzeniem wykorzystana jako uczący i testują-

cy zbiór danych dla kilku algorytmów rozpoznawania emocji, co wykazało przydatność zarówno jej, jak również narzędzia KinectRecorder.

**Słowa kluczowe**: akwizycja oraz indeksowanie nagrań RGB-D, rozpoznawanie mimiki, rozpoznawanie emocji

## 1. Introduction

Human emotions can have great impact on the way of using software, learning or training results, and also overall experience received from using computers. Therefore, robust and accurate recognition of human emotional states can play an important role in many software systems in different application fields. Affective computing is one of the emerging research areas on human behavior that develop emotion recognition, interpretation, and processing methods to create affective and affect-aware software to better adapt its behavior to user needs. Such applications can have a significant impact in many fields such as healthcare, education, entertainment, software engineering, e-learning, etc. [1]-[3]. Affect recognition methods can use different information channels, such as video [4], audio [5], standard input devices [6, 7], physiological signals [8, 9], depth information [10], and others.

The most popular source of non-invasive and non-intrusive information is a video camera as it provides for the recognition of facial expressions, gestures, posture, and body movements, which are typically used by human to analyze emotions of other people [11]. Unfortunately, algorithms that use video camera are sensitive to face illumination conditions, especially, to low light or unevenly illumination of the scene. One of possible solutions is to exploit additional channel of scene depth information. As many depth sensors use additional, infrared light source the information is to a large degree insensitive to different light conditions. Rapid development and high availability of relatively cheap RGB-D consumer sensors, such as Microsoft Kinect, makes it possible to create of real-time systems for recognition of human emotions.

As most reported approaches use supervised training to create an emotion classifier, preparing of proper training and testing datasets is of great importance. Unfortunately, many researchers use different and relatively small data sets, what limits estimation of their generalizability and real recognition effectiveness, and makes it difficult to compare different approaches. The lack of popular, convenient and commonly available tools for data acquisition and indexing seems to be one of the reason of this state.

In this paper, KinectRecorder comprehensive tool is described which provides for convenient and fast acquisition, indexing and storing of RGB-D video streams from Microsoft Kinect sensor. The application can be used for creation of fully indexed databases of facial expressions and emotions that can be further exploited for learning and testing of emotion

recognition algorithms for affect-aware applications. In the following section a short background is given on the Microsoft Kinect sensor, human face indexing for facial emotion recognition and the problems of RGB-D data acquisition and storing. Next, the concept of KinectRecorder tool is presented followed by a report on its application in the creation of Facial Expression and Emotion Database [12, 13]. Finally, some discussions and conclusions are presented.

## 2. Background

Although there are many available databases with face images for face recognition [14], only small part of them can be used to support recognition of facial expressions and emotions [15]-[18]. Most databases offer only 2D images or video recordings [15, 16], but some databases contain also 3D images acquired by the means of stereography or specialized 3D scanners [17, 18]. Additional information about the third dimension in the image, that is its *depth*, can significantly improve efficiency of some algorithms for face detection and recognition of facial expressions or human emotions.

Depth sensors have become very popular 3D acquisition devices recently, that makes it possible to develop new algorithms for recognition of human face, facial expressions, gestures, and body posture. Enriching traditional algorithms based on RGB cameras with additional information channel can lead to better localization of parts of human's body, and improve recognition efficiency in difficult illumination conditions.

Microsoft Kinect is one of most interesting depth sensors nowadays, that is why many research in the area of face detection or recognition of facial expressions and emotions is based on this device. Kinect is able to provide any application for three raw data streams containing color frames (RGB), scene depth frames (D), and audio information (A). Though the color and depth channels are separated and their fields of view are slightly shifted each other, they can be merged together in order to create a complete three dimensional image of any object within the sensor's field of view. Microsoft offers its propriety XED format for storing the RGB-D streams from the Kinect sensor. There are two important advantages of this format. Firstly, it stores both video streams in one file, which is very useful from a data administration point of view. Secondly, any XED file may simulate the direct on-line input from the Kinect sensor, thus making the repetitive tests possible. Unfortunately, there are also several drawbacks of using XED format. Firstly, this is Microsoft proprietary file format with no possibility of accessing a file content without a Microsoft Kinect Studio (MKS) application. Secondary, using XED file demands that a Kinect sensor is attached to the system, which is sometimes inconvenient. The last problems concern the MKS software which provides some

basic editing operations but the final file is lossy compressed, which makes such recordings sometimes unsuitable for further face processing. Unfortunately, there is currently no option available to switch off or to reduce this compression. For these reasons storing RGB-D data in other, open format seems very useful for many researchers.

Generally, recordings of human facial expressions and emotions may be divided into two categories: containing posed and spontaneous session. Posing allows for recording extreme emotions such as despair, envy, and others, that are hardly observed in real life in computer using situations. Drawbacks of posed recording are human problems to perform them and their artificial character. On the other hand recording of spontaneous emotions is much harder and requires a lot of time to catch the desired emotion or expression. That is why databases of facial expressions and emotions should contain images or recordings of both types.

Correct indexing, or labeling of the data set is the key issue for the possibility of its use as a training set for creating pattern recognition algorithms. The information about the actual expression as well as specific information about the face geometry or its state are required for recognition of facial expressions and emotions. As observed emotions are sometimes different from intended ones two or more labels are usually included to each expression. A name of an emotion is often extended by additional description of a human emotional state expressed in some emotional space, such as PAD (pleasure, arousal, dominance) [19]. Self-assessment manikin (SAM) scale [20] is a very convenient way of declaring somebody's feelings in the PAD emotional space.

The location of facial characteristic points, or facial landmarks, is another very important information for facial expressions recognition algorithms. In the case of video recordings (instead of still images) such information should be available for at least few frames representing different stages of facial expressions. Facial action coding system (FACS) is another way of labeling video facial recordings that is frequently used [21]. FACS contains a set of nearly 100 action units (AU) that describes movements of different muscles or parts of a human face.

Finally, some additional information about the scene is also useful, such as illumination conditions or appearance of other persons in the background.

## 3. KinectRecorder

KinectRecorder is a comprehensive tool for convenient and fast acquisition, labeling and storing of RGB-D video streams from the Kinect sensor. The application meets all the requirements presented in the previous section. One of its most important feature is the supervision over whole recording experiment. The application leads the user by hand indicating next

steps, allowing to repeat some of them, and formally validating the results. This assures that it is possible to start using the application without any special training.

The application was created in C++ and C# languages using Microsoft Visual Studio. Its capabilities are described in more detail in the following subsections.

### 3.1. Acquisition of RGB-D video streams

KinectRecorder enables acquisition of RGB-D video streams and writes them in a single AVI video container. Though, storing of two video channels in one file is not a very common solution it enables easier access and synchronization of both streams. The maximum resolution of both channels from Kinect is relatively high and is 640×480 at 30 fps for depth channel and 1280×960 at 12 fps for the color one. As KinectRecorder merges both streams together it uses VGA resolution for both channels. These two video streams are not compressed to avoid losing details in particular frames. Another advantages of this project decision are simple frame access and avoiding problems with codec distribution. Though, the resulting video files are larger without compression, their size still permits for efficient processing. Unfortunately, due to the specific nature of the Kinect sensor, frames from the color and depth channel are not synchronized. Moreover, in low light conditions, the color camera delivers frames at a significantly lower speed. In order to produce two synchronized streams KinectRecorder fills the missing color frames with the nearest ones, which enables easier displaying and processing of both video streams.

All recordings are named according to the names given in XML configuration file supported for each recording session. For each experiment attendee and for each recording session a separate folder is automatically created basing on a user's identifier, which helps to avoid confusion with file naming.

### 3.2. Acquisition scenarios

KinectRecorder allows to acquire recordings in three different modes which are prepared for different experiment scenarios, namely *posed*, *evoked spontaneous*, and *natural spontaneous* expressions. Though, the acquiring interface of the application remains very similar for all scenarios (Fig. 1), they slightly differ in the assumptions of the experiment and its course. For posed expressions, experiment attendees are asked to record a certain number of expressions according to instructions displayed on screen. The user has enough time between subsequent recordings to adequately prepare for them and can repeat them as many times as it is needed. The assumption is that each recording starts and ends with a neutral face expression and last only a few seconds (e.g. 3-5). Although this assumption is not forced in any way by the application it seems sensible as gives the researchers a full facial expression

as it appears in a real life. Similarly, too short or too long expressions look rather unnatural and are not suitable for automatic analysis. The number of required expressions, their names, the multimedia source file and additional text displayed on-screen, the preparing and recording time, and the output file name are fully customizable by an XML configuration file. This solution makes it possible to adapt application behavior to the needs of the particular acquisition experiment.
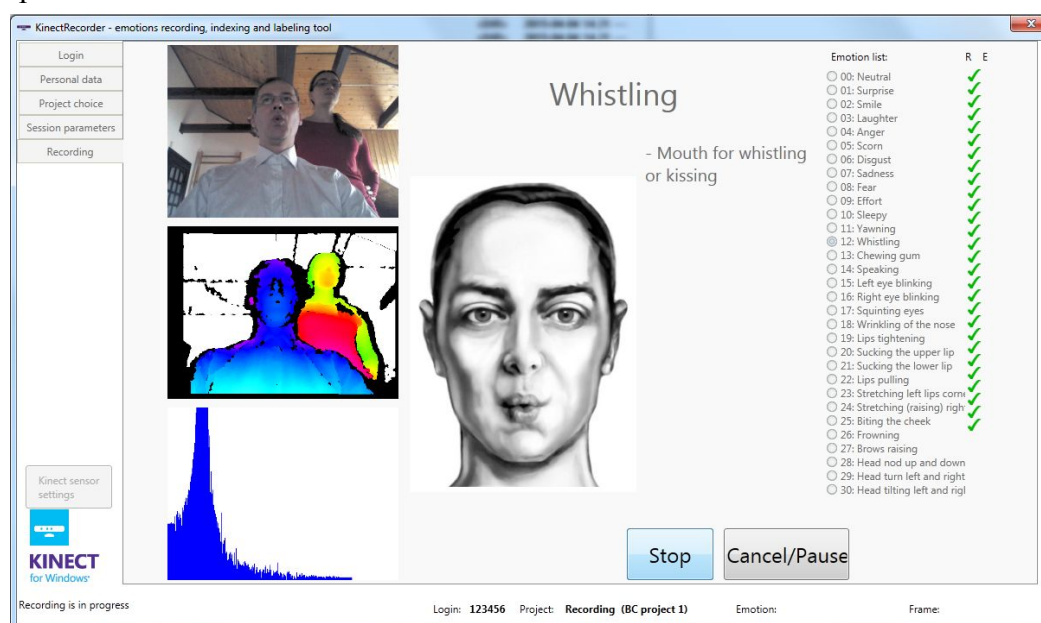


Fig. 1.   KinectRecorder in the first scenario mode
Rys. 1.   KinectRecorder w trybie pierwszego scenariusza

The main assumption of the second scenario is to present the user with a series of multimedia stimuli with a defined timing of preparation, stimulation and calming phases. Also this time the course of the scenario may be freely configured with and external XML file. KinectRecorder is able to use many multimedia formats including still images in .jpg and .png formats, raster animations in .gif files, vector animations in Flash .swf format, and many others.

Finally, the third scenario provides for recording of fully natural emotions that are caused by another application or other stimuli source. In this mode, only one recording is usually done but it can be much longer of up to several minutes. The goal of this scenario is to make the user to forget about the experiment conditions and try to his or her natural reaction to other application, e.g. very stimulating video game.

### 3.3.  Data labeling

KinectRecorder offers an advanced replaying and indexing tools. Labeling can be performed on three description levels. At the top level, users introduce some meta-data about themselves (e.g. age, gender, etc.) and about the experiment conditions (e.g. lighting condi-

tion, scene background, etc.).  At the medium level they can describe their emotion for each recording using self-assessment manikin (SAM) scale [19] and action units (AU) from facial action coding system (FACS) [21]. The interface of this indexing tool is presented in Fig. 2.
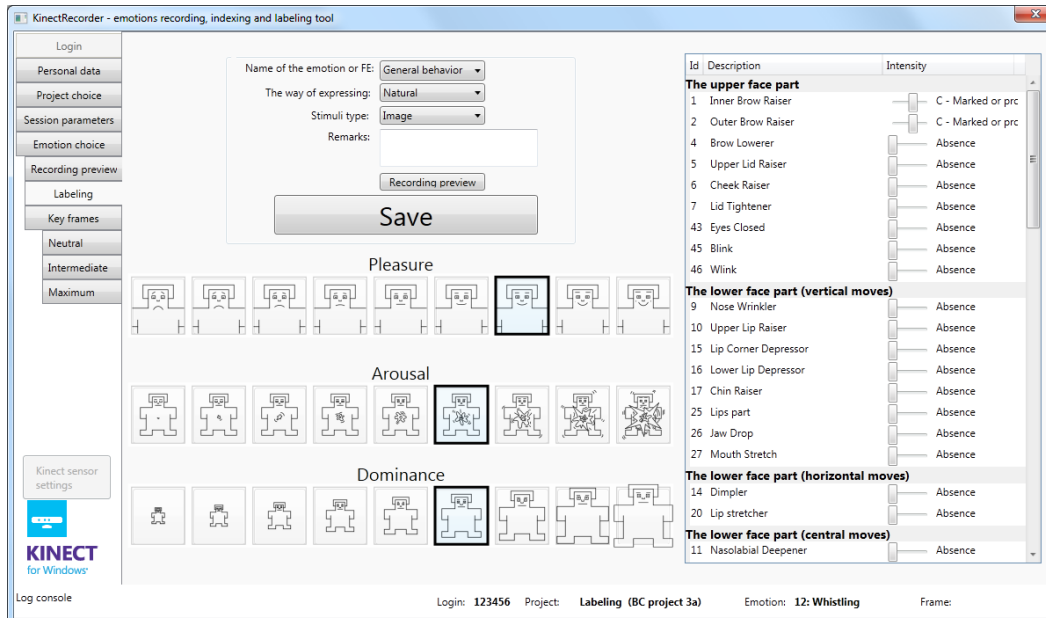


Fig. 2.  KinectRecorder labeling tool for each recording (SAM and FACS)
Rys. 2.  KinectRecorder – narzędzie indeksujące dla każdej ramki (SAM i FACS)

The access to the lowest indexing level is possible from the player tool that allows to re-play each recording in continuous and frame-by-frame modes in both directions. Both color and depth streams are displayed synchronously in adjacent windows (Fig. 3). The application allows to indicate up to three special frames that could for example denote the beginning, increase and maximum of the recorded expression. Each of such frames is additionally stored in two BMP images one for the color channel and the second for the depth one.
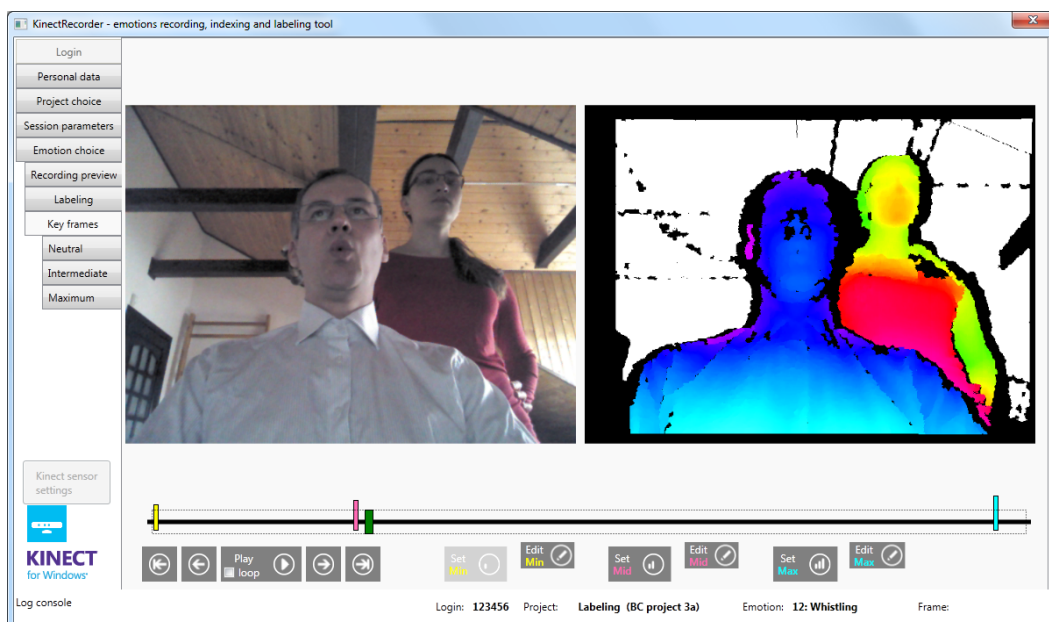
Fig. 3.    KinectRecorder in a recording viewing mode
Rys. 3.    KinectRecorder w trybie przeglądania nagrań

Choosing any of these three frames allows to select 27 facial landmarks such as mouth and eyes corners, pupils, nasal root and tip, brows etc. This indexing tool is very easy and fast to use thanks to additional illustration of the facial characteristic points, and a handy magnifier (Fig. 4). The location of each facial landmark may be corrected at any time. It is also possible to visualize those points against the depth image.
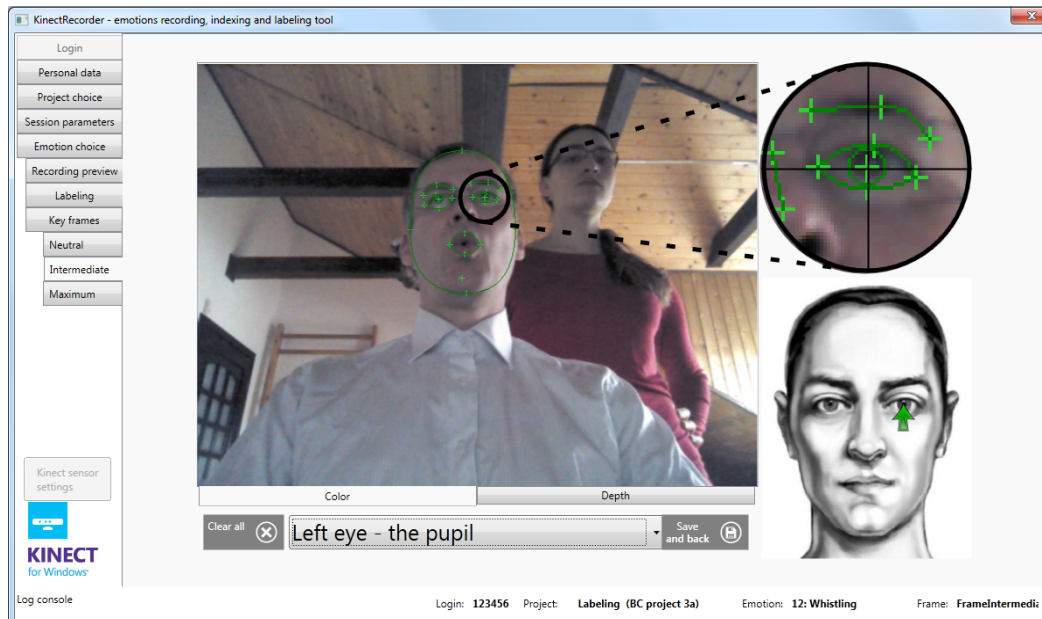


Fig. 4.    KinectRecorder indexing tool for facial landmarks in chosen frames
Rys. 4.    KinectRecorder – narzędzie indeksujące charakterystyczne punkty twarzy w wybranych klatkach

## 4. Application

KinectRecorder was successfully applied to create the second version of Facial Expression and Emotion Database [13]. As the first version of FEEDB was created using Microsoft Kinect Studio [12] it was possible to directly compare both tools in terms of their functionality and ease of use. Except the introduction of the new file format, KinectRecorder has many advantages over MKS, which has been proved in many experiments carried out through the last two years. During these experiments, over 100 students, aged 19 to 23, of both gender were asked to use KinectRecorder in two scenarios. In the first scenario attendees were asked to express 31 facial expressions and emotions according to displayed instructions. In the second scenario students were automatically stimulated by 10 multimedia materials and their reactions were recorded. All these recordings were then indexed by the recorded persons ac-

cording to description given in the previous section. Finally, 50 complete sets of recordings of the best quality were chosen to create FEEDB v.2 [13].

One of additional tasks for attendees of the experiments was to evaluate KinectRecorder. In general, the application received very good marks, especially for its user friendly interface, effective supervision of the experiment scenarios and very convenient tools for indexing and file naming. All students praised the comfort of the situation when both recording controls as well as displaying the stimuli multimedia are managed within the same application, contrary to MKS. Very important were also the possibility to repeat failed recordings and some additional time to prepare for a new recording. KinectRecorder verifies also current completeness of recordings and their indices, which is very handy to ensure the completeness of the structure of all saved files both for experiment attendees as well as for their supervisor.

FEEDB data sets were successfully used in a few research projects aiming in recognition of human emotions [22]. Additionally, to enable easy viewing and searching of the database content FEEDB Digital Library was developed basing on dLibra engine. This required developing of additional mass import procedures and a specialized plug-in for synchronous replaying of color and depth video streams. Both usage of FEEDB proved that introducing of new format for storing color and depth streams is not a limitation to the researchers.

## 5. Conclusion

Microsoft Kinect Studio (MKS), that is distributed with each Kinect Developer Toolkit, seems to be an obvious solution in case of developing a dataset of recordings stored in the XED format. Unfortunately, this is Microsoft propriety format with several limitations mentioned earlier. This cause that researchers have to develop their own software to get access to particular frames of color and depth streams, anyway. Moreover, MKS does not offer any additional functionality supporting labeling or indexing of acquired recordings. These processes are often done by hand which can lead to many mistakes and errors. Developing additional independent tools that would combine acquisition, organizational and labeling functionality is the natural alternative way to overcome above limitations. The main drawback of this solution is necessity of introducing a new file format for storing video recordings. Though this new format forces researchers to create their own procedures for stream opening, this is very similar to limitations associated with the use of XED.

Described KinectRecorder is a comprehensive tool which allows for convenient and fast acquisition, labeling and storing of RGB-D video streams from Microsoft Kinect sensor. The application provides functionality that is missing in popular Microsoft Kinect Studio and is especially useful as a supporting tool for creation of fully labeled databases of facial

expressions and emotions that can be further used for learning and testing emotion recognition algorithms for affect-aware applications. KinectRecorder was successfully used for creation of the second version of Facial Expression and Emotion Database. Taking into account experiences within this project it can be said that KinecRecorder significantly reduces the time of the whole process of data acquisition, labeling and validation. Kinect recorder is available for free for other researchers.

**BIBLIOGRAPHY**

1.   Picard R.: Affective Computing: From Laughter to IEEE. IEEE Transactions On Affective Computing, Vol. 1, No. 1, 2010, p. 11÷17.
2.   Landowska A.: Affect-awareness Framework for Intelligent Tutoring Systems. The 6th Int. Conf. on Human System Interaction, IEEE, 2013, p. 540÷547.
3.   Wrobel M. R.: Emotions in the software development process. The 6th Int. Conf. on Human System Interaction, IEEE, 2013, p. 518÷523.
4.   Gunes H., Piccardi M.: Affect Recognition from Face and Body: Early Fusion vs. Late Fusion. IEEE Int. Conf. on Systems, Man and Cybernetics, Vol. 4, IEEE, 2005, p. 3437÷3443.
5.   Zeng Z., Pantic M., Roisman G., Huang T.: A survey of affect recognition methods: Audio, visual, and spontaneous expressions. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 31, No. 1, IEEE, 2009, p. 39÷58.
6.   Vizer L. M., Zhou L., Sears A.: Automated stress detection using keystroke and linguistic features. Int. Journal of Human-Computer Studies, Vol. 67, 2009, p. 870÷886.
7.   Kolakowska A.: A review of emotion recognition methods based on keystroke dynamics and mouse movements. The 6th Int. Conf. on Human System Interaction, IEEE, 2013, p. 548÷555.
8.   Bailenson J. N., Pontikakis E. D., Mauss I. B., Gross J. J., Jabon M. E., Hutcherson C. A. C., Nass C., Oliver J.: Real-time classification of evoked emotions using facial feature tracking and physiological responses. Int. Journal of Human-Computer Studies, Vol. 66, No. 5, 2008, p. 303÷317.
9.   Szwoch W.: Using Physiological Signals for Emotion Recognition. The 6th Int. Conf. on Human System Interaction, IEEE, 2013, p. 556÷561.
10.  Burgin W., Pantofaru C., Smart W. D.: Using depth information to improve face detection. The 6th Int. Conf. on Human-Robot Interaction, IEEE, 2011, p. 119÷120.

11.    Castrillón-Santana M., Déniz-Suárez O., Antón-Canalís L., Lorenzo-Navarro J.: Face and Facial Feature Detection Evaluation. Int. Conf. on Computer Vision and Computer Graphics Theory and Applications, 2008, p. 167÷172.

12.    Szwoch M.: FEEDB: a Multimodal Database of Facial Expressions and Emotions. The 6th Int. Conf. on Human System Interaction, IEEE, 2013, p. 524÷531.

13.    Szwoch M.: On Facial Expressions and Emotions RGB-D Database. 10th Int. Conf. Beyond Databases, Architectures, and Structures, CCIS, Vol. 424, Springer, 2014, p. 384÷394.

14.    Face Recognition Homepage: http://www.face-rec.org/databases/, 31.03.2015.

15.    Lucey P., Cohn J. F., Kanade T., Saragih J., Ambadar Z., Matthews I.: The extended Cohn-Kanade (CK+): A complete dataset for action unit and emotion-specified expression. Conf. on Computer Vision and Pattern Recognition, IEEE, 2010, p. 94÷101.

16.    Wang S., Liu Z., Lv S., Lv Y., Wu G., Peng P., Chen F., Wang X.: A Natural Visible and Infrared Facial Expression Database for Expression Recognition and Emotion Inference. IEEE Transactions on Multimedia, Vol. 12, No. 7, IEEE, 2010, p. 682÷691.

17.    Savran A., Alyüz N., Dibeklioğlu H., Çeliktutan O., Gökberk B., Sankur B., Akarun L.: Bosphorus Database for 3D Face Analysis. Lecture Notes in Computer Science, Vol. 5372, Springer, 2008, p. 47÷56.

18.    Colombo A., Cusano C., Schettini R.: UMB-DB: A Database of Partially Occluded 3D Faces. IEEE Int. Conf. on Computer Vision Workshops, IEEE, 2011, p. 2113÷2119.

19.    Mehrabian A: Comparison of the PAD and PANAS as models for describing emotions and for differentiating anxiety from depression. Journal of psychopathology and behavioral assessment, Vol. 19, No. 4, Springer, 1997, p. 331÷357.

20.    Bradley M. M., Lang P. J.: Measuring Emotion: The Self-Assessment Manikin and the Semantic Differential. Journal Behav. Ther. &Exp. Psychiat., Vol. 25, No. 1, Elsevier, 1994, p. 49÷59.

21.    Ekman P., Friesen W.: Facial Action Coding System: A Technique for the Measurement of Facial Movement. Consulting Psychologists Press, Palo Alto, 1978.

22.    Szwoch M., Szwoch W.: Emotion Recognition for Affect Aware Video Games. Image Processing & Communications Challenges 6 Advances in Intelligent Systems and Computing, Vol. 313, Springer, 2015, p. 227÷236.

**Omówienie**

W pracy przedstawiono istotną rolę, jaką pełnią właściwie zaindeksowane bazy danych uczących i testujących w rozwoju algorytmów rozpoznawania nagrań mimiki i emocji. Pod-

kreślono również dodatkowe korzyści, jakie daje w tych algorytmach wykorzystanie dodatkowych kanałów informacji, a w szczególności kanału głębi. Wskazano również, że w chwili obecnej najpopularniejszym sensorem dostarczających danych o kolorze i głębi (RGB-D) jest Microsoft Kinect, jednak wciąż brakuje baz danych z nagraniami mimiki i emocji z tego sensora. Do tworzenia tego typu baz zaproponowano stworzone w tym celu oryginalne narzędzie KinectRecorder, które, oprócz zarządzania samym proces pozyskiwania nagrań, pozwala na samodzielne przygotowywanie scenariuszy eksperymentów, zarządza nazewnictwem plików nagrań, weryfikuje ich kompletność oraz pozwala na wygodną ich indeksację na poziomie metadanych, pojedynczych nagrań oraz charakterystycznych punktów twarzy dla wybranych klatek nagrań.

Następnie przedstawiono skuteczne wykorzystanie aplikacji do stworzenia drugiej wersji bazy nagrań mimiki i emocji (FEEDB v.2), zwracając uwagę na bardzo pozytywne opinie użytkowników. Przedstawiono również możliwość prezentacji tak stworzonej bazy danych w formie biblioteki cyfrowej z wykorzystaniem silnika dLibra. W podsumowaniu podkreślono oryginalność zaprezentowanego narzędzia oraz jego przewagę nad standardową aplikacją Microsoft Kinect Studio, które wykorzystuje do przechowywania plików zamknięty format XED.

**Address**

Mariusz SZWOCH: Politechnika Gdańska, Wydział Elektroniki, Telekomunikacji i Informatyki, ul. Narutowicza 11/12 80-233 Gdańsk, Poland, szwoch@pg.gda.pl.