

# Web-based visualisation of head pose and facial expressions changes: monitoring human activity using depth data

Grigorios Kalliatakis

School of Computer Science and  
Electronic Engineering  
University of Essex, UK  
Email: gkallia@essex.ac.uk

Nikolaos Vidakis

Department of Informatics Engineering  
Technological Educational  
Institute of Crete, Greece  
Email: nv@ie.teicrete.gr

Georgios Triantafyllidis

Mediology Section, AD: MT  
Aalborg University  
Copenhagen, Denmark  
Email: gt@create.aau.dk

**Abstract**—Despite significant recent advances in the field of head pose estimation and facial expression recognition, raising the cognitive level when analysing human activity presents serious challenges to current concepts. Motivated by the need of generating comprehensible visual representations from different sets of data, we introduce a system capable of monitoring human activity through head pose and facial expression changes, utilising an affordable 3D sensing technology (Microsoft Kinect sensor). An approach build on discriminative random regression forests was selected in order to rapidly and accurately estimate head pose changes in unconstrained environment. In order to complete the secondary process of recognising four universal dominant facial expressions (happiness, anger, sadness and surprise), emotion recognition via facial expressions (ERFE) was adopted. After that, a lightweight data exchange format (JavaScript Object Notation-JSON) is employed, in order to manipulate the data extracted from the two aforementioned settings. Such mechanism can yield a platform for objective and effortless assessment of human activity within the context of serious gaming and human-computer interaction.

## I. Introduction

Automatic and effective estimation of head pose is a challenging problem of computer vision systems. Since it is considered as a key element of human behaviour analysis, many applications would benefit from automatic and robust head pose estimation systems such as: (i) face recognition; (ii) human activity analysis; (iii) human-computer interaction and (iv) robotic vision. As a result, head pose estimation has drawn great attention from academia and a variety of techniques have been reported in the literature [1]–[4].

Likewise, the field of facial expression analysis is still regarded as an enthusiastic issue in the latest research works [5]–[7]. Due to its various purposes and applications, such as: (i) designing better human/ma-

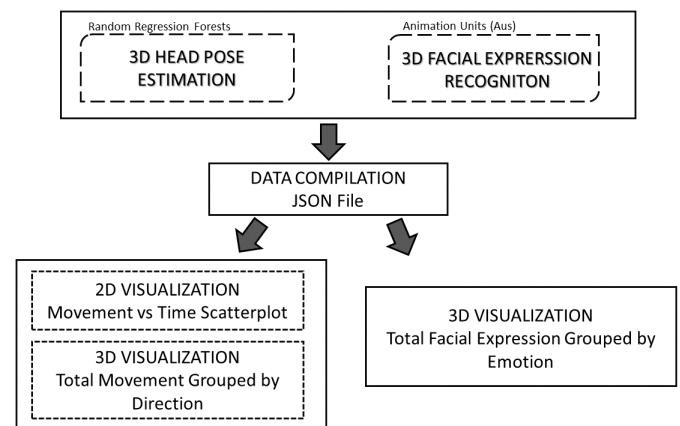


Fig. 1. An overview of the proposed system.

chine interfaces; (ii) video gaming; (iii) computer generated animations and (iv) identification, facial expression analysis plays a key role in emotion recognition and thus contributes to the development of human-computer interaction systems [8].

With the recent technological advancements in the area of depth sensors, it is now feasible to perform large-scale data collection for subsequent analysis. Providing an objective assessment and evaluation, findings such as head pose changes and facial expression variations, can lead to valuable conclusions regarding the overall experience of users in many applications. One such application is the evaluation of the player's training and ludology<sup>1</sup> experience in the case of serious games such as [9], [10]. In this context, accessible visualisations can play a major part in that kind of assessment by creating encodings of data into visual channels

<sup>1</sup> A borrowing from Latin word "ludus" (game), combined with an English element; The term has historically been used to describe the study of games.

that people can view and understand comfortably. The process of data visualisation is suitable for externalizing the facts and enabling people to understand and manipulate the results at a higher level. Additionally, visualisations can be used in several distinct ways to help tame the scale and complexity of the data so that it can be interpreted effortlessly.

In this paper, we address the problem of human activity monitoring through head pose and facial expressions changes and, in the same time, we introduce two innovative web-based visualisations for evaluating purposes with respect to those data. The proposed system consists of three distinctive components as shown in Fig.1. First the real-time head pose estimation and facial expression events are separately obtained for different users sitting and moving their head without restriction in front of a Microsoft Kinect sensor for specified intervals. Experimental results on 20 different users show that the proposed system can achieve 83.95% accuracy for head pose changes and 76.58% accuracy for facial expressions recognition when validated against manually constructed ground truth data. Then the data for every user session are stored in a JSON file for offline manipulation. In the last step, two different types of visualisations are exploited: (i) a scatter plot for demonstrating head pose changes and intensities; (ii) 3D columns for presenting the players movement grouped by direction and the players facial expressions grouped by the dominant emotion respectively. As noted in our previous work [11], which was limited to head pose changes, *the principal objective of this work is to acquire efficient and user-friendly visualisations in order to improve the understanding and the analysis of the captured data*, easily accessed through a web-page.

The remainder of the paper is structured as follows: Section II describes the head pose estimation framework that was employed in our approach, while Section III explains the method for real-time emotion recognition via facial expressions. The data compilation phase is presented in Section IV, while Section V contains a brief description of the libraries that were used for presenting the data on the web, before presenting the actual web-based visualisations which were made for user activity assessment. Finally Section VI consists of a summary and concluding remarks.

## II. Head Pose Estimation Framework

Systems relying on 3D data have demonstrated very good results for the task of head pose estimation, compared to 2D systems that have to overcome ambiguity in real time applications. We partly followed the approach of Fanelli *et al.* [12], as it is regarded

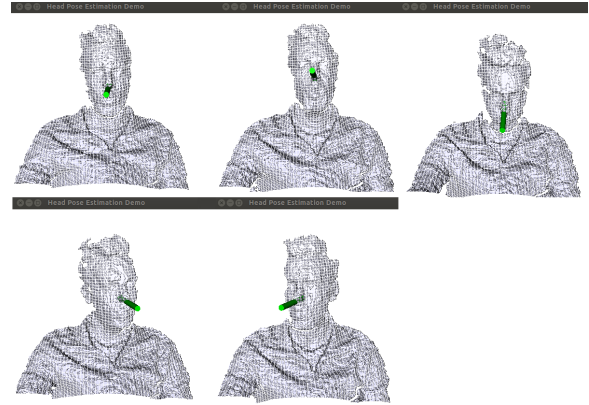


Fig. 2. Head pose estimation results.

to be suitable for real time 3D head pose estimation, considering its robustness to the poor signal-to-noise ratio of current consumer depth cameras like Microsoft Kinect sensor.

For this reason, regression forests are being extended in such a manner that depth patches belonging to a head can be discriminated and solely used for the prediction of the pose resulting in solving both the classification and regression problems respectively. While several works in the literature contemplate the case where the head is the only object present in the field of view [13], the proposed method concerns depth images where other parts of the body might be visible at the same time, and therefore need to be disjointed into image patches either belonging to the head or not. The system is able to perform on a frame-by-frame basis while it runs in real time without the need of initialization. Forests of randomly trained trees are less sensitive to over-fitting and generalize better than decision trees independently. In the proposed setup [14], depth patches are annotated with class label and a vector containing the offset between the 3D points falling on the patch's center and the head center location, plus the Euler rotation angles describing the head orientation. Randomness is imported in the training process, either in the set of training examples provided to each tree or in the set of tests used for optimization at each node, or even in both. When the pair of classification and regression are engaged, the aggregation of trees which simultaneously separate test data into positive cases (they represent part of the object of interest) are labeled as Discriminative Random Regression Forests (DRRF). This signifies that an extracted patch from a depth image is sent through all trees in the forest. The patch is evaluated at each node according to the stored binary test and passed either to the right or left child until a leaf node is



Fig. 3. Facial expression recognition (FER) results.

reached [15], at which point it is classified and only if this classification outcome is positive (head leaf), a Gaussian distribution is recaptured and then used for casting a vote in a multidimensional continuous space which is stored at the leaf. Fig.2 shows some processed frames regarding two DOF (*pitch* and *yaw*). Starting from left to right, the first row estimations displayed are: *still*, *up*, *down*. The second row estimations are *left* and *right* correspondingly. All calculations derived from the difference between the exact previous frame and the current frame, at each iteration of the program. The green cylinder encodes both the estimated head center and direction of the face.

### III. Facial Expression Recognition Framework

Emotion recognition via facial expressions (ERFE) is a growing active research field in computer vision compared to other emotion channels, such as body actions and speech, primarily because superior expressive force and a larger application space is provided.

A similar to [16] approach was followed for real-time emotion recognition. Video sequences acquired from the Kinect sensor are regarded as input. Then face detection and feature extraction are performed on each frame of the stream. The Face Tracking SDK [17], which is included in Kinect's Windows Developer toolkit, is used for tracking human faces with RGB and depth data captured from the sensor. Furthermore, facial animation units and 3D positions of semantic facial feature points can be computed by the face tracking engine, which can lead to emotion recognition via facial expressions.

Face tracking results are expressed in terms of weights of six animation units, which belong to a subset of what is defined in the Candide3 model [18]. Each AU, that is deltas from the neutral shape, is expressed as a numeric weight varying between  $-1$  and  $+1$ , and the neutral states of AUs are normally assigned to 0. The AU's feature of each frame can be written in the form of a 6-element vector:

$$\bar{a} = (A_1, A_2, A_3, A_4, A_5, A_6) \quad (1)$$

where  $A_1, A_2, A_3, A_4, A_5$ , and  $A_6$  refer to the weights of *lip raiser*, *jaw lower*, *lip stretcher*, *brow lower*, *lip corner depressor*, and *brow raiser*, respectively. For the purpose of this paper four different emotions were tested: *anger*, *happiness*, *sadness* and *surprise* as shown in Fig.3.

### IV. Data Compilation

In this section, the data compilation process based on the aforementioned frameworks is presented. Given the pitch  $pitch_t$  and yaw  $yaw_t$  intensities of the ongoing streaming frame, and the exact previous frame's pitch  $pitch_{t-1}$  and yaw  $yaw_{t-1}$  intensities, the system operates in three steps as follows: (a) the differences regarding pitch and yaw are calculated by (2)–(3);

$$pitchDiff = pitch_{t-1} - pitch_t \quad (2)$$

$$yawDiff = yaw_{t-1} - yaw_t \quad (3)$$

(b) then a threshold value was experimentally set around 4 in order for our system to ignore negligible head movements in all four directions tested; (c) finally, the changes with respect to the four different directions are given by (4)–(7).

$$up = pitchDiff > THRESH \quad (4)$$

$$down = pitchDiff < THRESH \quad (5)$$

$$left = yawDiff > THRESH \quad (6)$$

$$right = yawDiff < THRESH \quad (7)$$

Concerning the detection of emotions, boundaries for each Animation Unit had to be created in order to associate the vector obtained by the AU feature, as defined by (1), with the four main emotions. For example,  $(0.3, 0.1, 0.5, 0, -0.8, 0)$  corresponds to a happy face, which means showing teeth slightly, lip corner raised and stretched partly, and the brows are in the neutral position. We experimentally assembled the following equations (8)–(11) for our test sessions:

$$sadness = A_6 < 0 \wedge A_5 > 0 \quad (8)$$

$$surprise = (A_2 < 0.25 \vee A_2 > 0.25) \wedge A_4 < 0 \quad (9)$$

```

1 [
2   {"SessionDate": "2/Mar/16",
3     "SessionData": [
4       {
5         "time": 2.23,
6         "direction": "RIGHT",
7         "intensity": 6.78485
8       }
9     ]

```

Fig. 4. JSON structure for head pose changes.

```

1 [
2   {"SessionDate": "10/Mar/16",
3     "SessionData": [
4       {
5         "time": 7.98,
6         "emotion": "ANGRY",
7       }
8     ]

```

Fig. 5. JSON structure for facial expressions changes.

$$happiness = A_3 > 0.4 \vee A_5 < 0 \quad (10)$$

$$anger = ((A_4 > 0 \wedge (A_2 > 0.25 \vee A_2 < -0.25)) \vee (A_4 > 0 \wedge A_5 > 0)) \quad (11)$$

Regarding the storage of the obtained data, JavaScript Object Notation (JSON) format was used mainly because of its lightweight nature, convenience in writing and reading and more importantly, as opposed to other formats such as XML, its suitability in generating and parsing tasks in various Ajax applications as described in [19]. A record in an array was created for each user session, while an extra array was inside it, preserving three variables: *time*, *direction* and *intensity* for each movement that was detected as shown in Fig.4. For facial expressions, a similar array was created, but in this case only two variables were required: *time* and *emotion*, as shown in Fig.5.

## V. Visualisation

Many different approaches have been proposed in the literature to solve the problems of head pose estimation and facial expression recognition. However, very few focus on how those data must be presented in order to deliver a useful meaning conveniently. In this section, the final step of the proposed system is presented alongside the actual visualisation instances that can be found directly on the web. Furthermore, in the following subsections two JavaScript libraries, for data-driven document manipulation, are briefly exposed.

### A. D3: Data-driven Documents

*Data-Driven Documents* is a novel representation-transparent concept for web-based visualisations. This JavaScript library assists users at bringing data to life using varied technologies such as HTML for page content, CSS for aesthetics, JavaScript for interaction, SVG for vector graphics and so on. As claimed by Michael Bostock *et al.* in [20], D3's emphasis on web standards provides full capabilities of modern browsers while it combines powerful visualisation components and a data-driven approach to a shared representation of the page called the document object model (DOM). However D3 must not be considered as a traditional visualisation framework because rather than introducing a novel graphical grammar, D3 solves the problem of efficiently manipulating documents based on data. Therefore D3's fundamental contribution is a visualisation kernel, closer to other document transformers like *jQuery* [21] and *CSS*, rather than a framework.

### B. Highcharts

*Highcharts* is a charting library written in pure JavaScript which suggests an easy way of adding interactive charts to web applications. Currently many different chart types are supported by this library such as box plot, pie charts, column charts etc. Many of these can be combined in one chart. This library was first released in late 2009, more details can be found in [22]. One big advantage of this library lies in being packed with adapters, which means that it does not rely on one particular framework, but instead is pluggable to different frameworks. The default framework implementation of High-charts uses *jQuery* [21]; hence the only requirement for users is to load the *jQuery* library before *Highcharts*.

### C. Player Movement vs Time Graph

In the final step of the proposed system, two visualisations were established for the desirable web-based data interpretation regarding head pose changes. The first one is a 2D scatterplot displaying the head movement of the player over specified time period as shown in Fig.6.

In more details, x-axis represents the time scale in seconds during which the tests take place (Fig.5 shows only a zoomed portion of the whole scatterplot graph), while each label in y-axis symbolizes each different user performing the test. Four different arrows imitate the movement of the human's head in two DOF. Furthermore an additional feature is displayed when the mouse is hovering an arrow, showing the respective time each movement occurred and the intensity, which



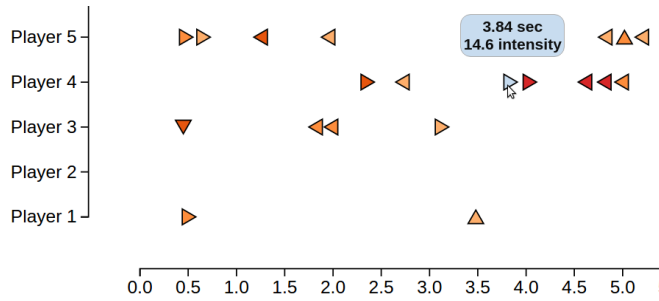


Fig. 6. 2D Scatterplot of head pose changes.

is based on how large the difference between the previous and the current frame was, as explained in Section IV. Apart from those elements, a color fluctuation is also evident which serves as an intensity indicator for each movement (the closer to red color the arrow is, the higher the intensity of the movement). One can easily examine the motion of the player that way, alongside its intensity, which adds a different dimension to the knowledge gained from the visualisation. The full version of this visualisation is available at: <http://83.212.117.19/HeadPoseScatterplot/>.

#### D. Overall Movement Grouped by Direction

The second visualisation consists of a 3D column diagram which illustrates the aggregation of all head movements grouped by direction every two seconds as shown in Fig.7. The four different directions are imitated by four different colors. In one hand, x-axis represents the time scale which is divided every two seconds until the end of the test. On the other hand, y-axis displays the number of movements for all the users that take part in the tests. Furthermore, when hovering above a column, the number of the corresponding direction summary is displayed. In this fashion, the dominant direction amongst all users every time interval is effortlessly assumed. Moreover, not so evenly distributed movements (e.g. columns between 2-4 seconds in Fig.7) can lead into practical conclusions taking into account the nature of the test as well. The full version of the overall head movement visualisation is available at: <http://83.212.117.19/HeadPose3D/>.

#### E. Overall Facial Expressions Grouped by Emotion

The visualisation regarding the recognised emotions via facial expressions is assembled in the same fashion as the previous one. In this case the facial expressions are grouped by the recognised emotions. Fig.8 displays only one emotion, *happiness*. However the rest of the recognised emotions can be set visible by clicking the corresponding check-box. The four different emotions are represented by four different

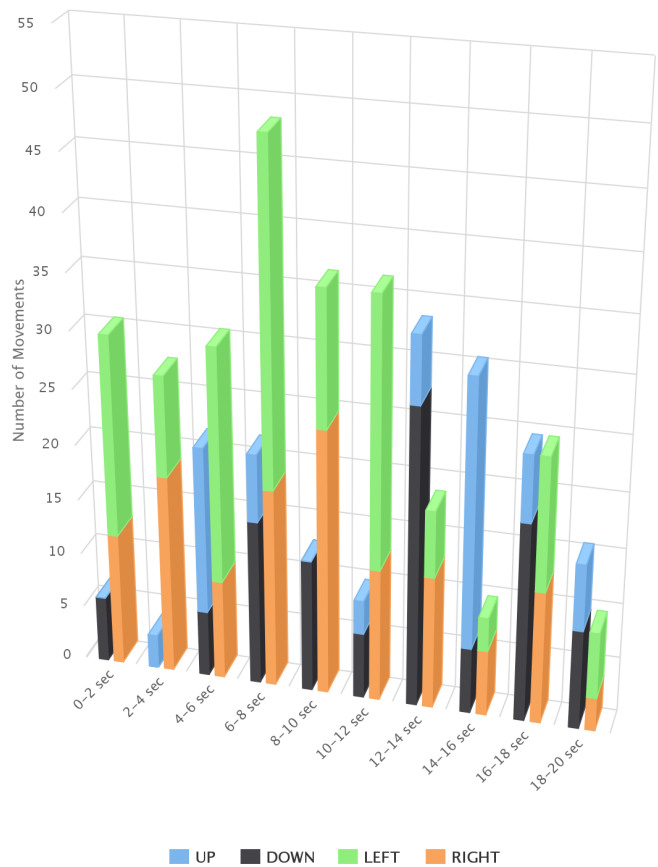


Fig. 7. 3D column visualisation of head pose changes.

colors. In one hand, x-axis represents the time scale which is divided every two seconds until the end of the test. On the other hand, y-axis displays the number of recognised emotions for all the users that take part in the tests. Furthermore, when hovering above a column, the number of the corresponding emotion summary is displayed. The full version of the overall facial expressions visualisation is available at: <http://83.212.117.19/FacialExpression3D/>.

## VI. Conclusion

In this paper we presented a system for generating efficient and user-friendly visualisations of head pose and facial expressions changes in the direction of human activity monitoring, utilising a consumer depth camera. Two selected approaches were extended in an applicable way for collecting and storing the required data using a key-value style lightweight exchanging format, JSON. Finally a 2D and two different 3D visualisations derived from the data compilation stage that can be easily accessed on the web. Our intention was to demonstrate that easily operated visualisations can provide a scaffold for objective and straightforward

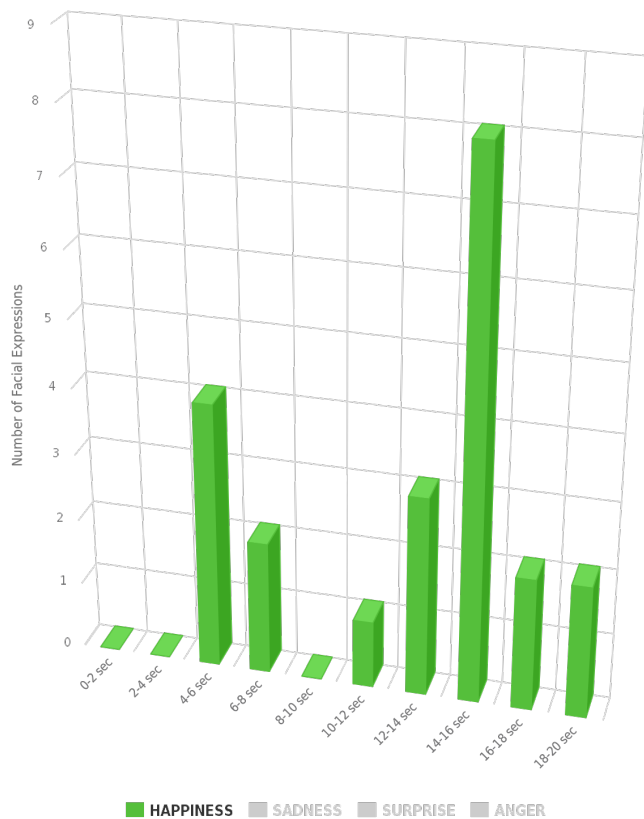


Fig. 8. 3D column visualisation of "HAPPINESS".

understanding of human activity across diverse applications, in the hope that it would provide a functional mechanism for future evaluations, and good baselines for human activity monitoring research. Interesting future directions will be to investigate whether other visualisation techniques can be more explanatory and effective, while fear and disgust could be included alongside the four main emotions.

### References

[1] M. D. Breitenstein, D. Kuettel, T. Weise, L. van Gool, and H. Pfister, "Real-time face pose estimation from single range images," in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, June 2008, pp. 1–8.

[2] L. P. Morency, J. Whitehill, and J. Movellan, "Generalized adaptive view-based appearance model: Integrated framework for monocular head pose estimation," in *Automatic Face Gesture Recognition, 2008. FG '08. 8th IEEE International Conference on*, Sept 2008, pp. 1–8.

[3] P. Paderleris, X. Zabulis, and A. A. Argyros, "Head pose estimation on depth data based on particle swarm optimization," in *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, June 2012, pp. 42–49.

[4] L.-P. Morency, "3d constrained local model for rigid and non-rigid facial tracking," in *Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, ser. CVPR '12, 2012, pp. 2610–2617.

[5] M. F. Valstar, T. Almaev, J. M. Girard, G. McKeown, M. Mehu, L. Yin, M. Pantic, and J. F. Cohn, "Fera 2015 - second facial expression recognition and analysis challenge," in *Automatic*

*Face and Gesture Recognition (FG), 2015 11th IEEE International Conference and Workshops on*, vol. 06, May 2015, pp. 1–8.

[6] G. Sandbach, S. Zafeiriou, M. Pantic, and L. Yin, "Static and dynamic 3d facial expression recognition: A comprehensive survey," *Image Vision Comput.*, vol. 30, no. 10, pp. 683–697, Oct. 2012.

[7] N. Hesse, T. Gehrig, H. Gao, and H. K. Ekenel, "Multi-view facial expression recognition using local appearance features," in *Pattern Recognition (ICPR), 2012 21st International Conference on*, Nov 2012, pp. 3533–3536.

[8] T. Fang, X. Zhao, O. Ocegueda, S. K. Shah, and I. A. Kakadiaris, "3d facial expression recognition: A perspective on promises and challenges," in *Automatic Face Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*, March 2011, pp. 603–610.

[9] N. Vidakis, E. Syntychakis, K. Kalafatis, E. Christinaki, and G. Triantafyllidis, "Ludic educational game creation tool: Teaching schoolers road safety," in *Universal Access in Human-Computer Interaction. Access to Learning, Health and Well-Being*. Springer, 2015, pp. 565–576.

[10] N. Vidakis, E. Christinaki, I. Serafimidis, and G. Triantafyllidis, "Combining ludology and narratology in an open authorable framework for educational games for children: the scenario of teaching preschoolers with autism diagnosis," in *Universal Access in Human-Computer Interaction. Universal Access to Information and Knowledge*. Springer, 2014, pp. 626–636.

[11] G. Kalliatakis, G. Triantafyllidis, and N. Vidakis, "Head pose 3d data web-based visualization," in *Proceedings of the 20th International Conference on 3D Web Technology*, ser. Web3D '15. ACM, 2015, pp. 167–168.

[12] G. Fanelli, T. Weise, J. Gall, and L. Van Gool, *Pattern Recognition: 33rd DAGM Symposium, Frankfurt/Main, Germany, August 31 – September 2, 2011. Proceedings*. Springer Berlin Heidelberg, 2011, ch. Real Time Head Pose Estimation from Consumer Depth Cameras, pp. 101–110.

[13] G. Fanelli, J. Gall, and L. V. Gool, "Real time head pose estimation with random regression forests," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, June 2011, pp. 617–624.

[14] G. Fanelli, J. Gall, and L. van Gool, "Real time 3d head pose estimation: Recent achievements and future challenges," in *5th International Symposium on Communications, Control and Signal Processing (ISCCSP)*, May 2012, pp. 1–4.

[15] G. Fanelli, M. Dantone, J. Gall, A. Fossati, and L. Van Gool, "Random forests for real time 3d face analysis," *International Journal of Computer Vision*, vol. 101, no. 3, pp. 437–458, 2013.

[16] Q.-r. Mao, X.-y. Pan, Y.-z. Zhan, and X.-j. Shen, "Using kinect for real-time emotion recognition via facial expressions," *Frontiers of Information Technology & Electronic Engineering*, vol. 16, no. 4, pp. 272–282, 2015.

[17] Microsoft. (2016, apr) Microsoft kinect sdk documentation – face tracking. [Online]. Available: <https://msdn.microsoft.com/en-us/library/jj130970.aspx>

[18] J. Ahlberg, "Candide-3 - an updated parameterised face," Tech. Rep., 2001.

[19] B. Lin, Y. Chen, X. Chen, and Y. Yu, "Comparison between json and xml in applications based on ajax," in *Computer Science Service System (CSSS), 2012 International Conference on*, Aug 2012, pp. 1174–1177.

[20] M. Bostock, V. Ogievetsky, and J. Heer, "D3: Data-driven documents," *IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis)*, 2011.

[21] jQuery. (2016, apr) jquery project. [Online]. Available: <https://jquery.com/>

[22] J. Kuan, *Learning Highcharts*. Packt, December 2012.