# A Smart Peephole on the Cloud

Maria De Marsico$^{(\boxtimes)}$, Eugenio Nemmi, Bardh Prenkaj, and Gabriele Saturni

Sapienza University of Rome, Rome, Italy
demarsico@di.uniroma1.it, eugenio.nemmi@uniroma1.it,
prenkaj.1602894@studenti.uniroma1.it, gabriele.saturni@gmail.com

**Abstract.** This paper does not present a novel technique for biometric recognition, but rather a novel way to use it. The proposal is to exploit cloud computing in order to support everyday applications. These are not necessarily bound to security, but span a wide range of possible useful tasks. This work presents a smart peephole able to recognize the person at the door, possibly automatically allowing entrance according to rules decided by the home keeper. The novelty is that very little processing is carried out locally, and biometrics is implemented as a service. The system relies on Microsoft Cognitive Services, a suite of remote services included in Microsoft Azure platform. The single user has to install nothing but a camera with a sound capture facility in correspondence to the peephole, and a lightweight software. A movement detector module triggers the capture/recognition activity. The captured audio and video samples are sent to the service. Most processing and recognition are carried out via the remote suite, and a final result is sent back to possibly trigger a response action. The present prototype includes face, speech and emotion recognition. It does not completely cover all system aspects. The aim is to demonstrate the feasibility of the approach.

**Keywords:** Biometrics · Microsoft cognitive services · Cloud computing

## 1 Introduction

Nowadays, more and more sophisticated devices and software applications have to address both traditional as well as new challenges. These arise with the increasing use of digital resources and information. Having to cope with enormous quantity of daily data a person downloads, uploads and reads, private personal data and information are exposed to potential, usually hidden threats. The possible weaknesses of password-based and token-based approaches led to the genesis and continuous development of biometric security systems [3]. In principle, any human trait (DNA, face, retina, palm, iris, fingerprints, ears, etc.) which differs from one person to another, is universally available and measurable, and permanent over time, can be used as a personal identification key [2]. A biometric system collects, processes and records this kind of data (enrollment) in order to verify the identity or identify a person at a later time (recognition).

Biometrics technology finds application in many diverse fields, which are divided by access control type, either physical or logical. When the size of data or the size of possible users increases, both processing and storage resources may become an issue. This is the joining point between biometrics and cloud computing. When the single user cannot afford maintaining locally the needed resources, it is possible to rely on a smart remote distribution and deployment of storage, resources, and processing tools, which is the basic mechanism underlying cloud technology. In general, more and more kinds of resources can be remotely provided as a service. Microsoft Cognitive Services (MCS) is used here as an example platform to demonstrate the feasibility of biometrics as a service.

*Smart Peephole* can be deployed as a house security system. As the name suggests, it has the capability to observe incoming objects (humans, animals and so on) at the house doorstep. It carries out the basic function of a classic door peephole; however, rather than the person looking through the peephole, the system will take the decision whether to open the door, thus accepting the object in front of the door, or to remain in a closed state, rejecting the objects entrance. The system is composed of four main modules, implementing different biometric recognition tasks and combining their results, which will be described in detail in the following sections. These modules rely on MCS, and are combined together to make up a robust control mechanism, able to discriminate between members of the allowed group of people (*members*), and *intruders.*

**Peephole module** is responsible of detecting movements and to recognize whether, at the moment a movement is detected, there is a face in front of the peephole.

**Face detection module** can both register a family member (or other allowed subject) to the system by uploading a number of photos of that person's face and facial expressions, and identify the person later by matching the incoming face image with the gallery set (enrolled faces).

**Speech verification module** uses a microphone to record a person' s voice and upload it to the system; when a user desires to enter the house, he/she has to speak in through an appropriate microphone to be recognized as a member of the allowed group.

**Emotion detection module** aims at catching the mood of the person in front of the peephole, by extracting and analyzing features from facial expression.

The rest of the paper continues as follows. Section 2 presents the general ideas underlying this biometric security system, and its key requirements. Also, it discusses the decisions taken to implement efficiently the system and to guarantee a high responsiveness to the user. Microsoft Cognitive Services are also briefly introduced. Section 3 describes *Smart Peephole* system and its design. Section 4 synthesizes the system performance, and finally Sect. 5 draws conclusions and briefly summarizes future developments.

## 2  System Requirements and Implementation Choices

### 2.1  The Functions of *Smart Peephole* and the remote choice

The aim of *Smart Peephole* is to realize an intelligent peephole for the recognition of family members or house mates, such as extended family members or family friends. The purpose of this project is to create a smart system able to grant access to the house in a secure and easy way and, possibly, to eliminate the need for a key to open the house door. This is achieved by a combination of biometric measurements. As a particular security system, *Smart Peephole* should address the three principles of data security: Confidentiality, Integrity and Availability (CIA triad [14]). A detailed discussion about related topics is out of the scope of this paper, and also implementation of possible countermeasures was out of the scope of the present design.

  *Smart Peephole* has to recognize subjects belonging to two groups: *members* and *intruders*. *Members* are persons that are registered by an administrator, who will install the equipment and set up the system. The *intruders*, on the other hand, are all other persons that are not in the members set. *Smart Peephole* recognizes persons via face and voice recognition. As any biometric system, it includes two phases:

(1) In the first phase, or the *enrollment* phase, a person is registered by memorizing an arbitrary number (typically twenty) of photos into the system and by recording his/her voice. The system searches for face-like silhouettes in the stored photos: if a photo does not contain a face, the system responds with a rejection of that photo, thus the user has to send another picture of his/her face. The sound quality is verified too: if there is noise in the background, the user is required to resend another clearer voice recording. For security matters, a person is required to insert a secret password just in case that the first two barriers will be broken by an impostor. At the end of this procedure, the person will result enrolled and can be recognized correctly during the identification phase.
(2) The second phase, or the *identification* phase, requires for a person to stand in front of the camera in order to capture a face picture. If the picture matches with one of the set of photos in the system memory, then the user is prompted to speak to an appropriate microphone; when this is completed, the system checks whether the audio recorded corresponds to the same person's voice whose audio lies in the systems. If these two verification steps lead to a positive result then the person is allowed to enter the house classifying the person as a *member*; otherwise if one of these two steps does not end leading to a "true result", then the person is required to provide his/her secret password memorized at enrollment time. If the password provided is wrong, then the person is refused to enter the house and is classified as an *intruder*.

The face recognition task could have relied on a vast variety of methods, including the popular Local Binary Pattern (LBP) with its variations ([1,11]).

In the same way, voice recognition could have been implemented as suggested in [7] where feature extraction relies on an algorithm based on Mel Frequency Cepstral Coefficients (MFCC), widely exploited for voice processing, while Dynamic Time Warping(DTW) is used to measure the similarity between two time series which may vary in time or speed. Starting from the wheel when building an application is nowadays substituted by the massive use of library functions made available in many contexts. As for image processing, it is sufficient mentioning OpenCV [6]. The functions provided can be integrated seamlessly within one owns' code. The next step is to rely on complete frameworks, where functions are executed by a completely detached module according to a black-box model, where only interfaces or API are exposed. When the detached module is a remote one, it is possible to consider this at large as an example of programming by services or cloud computing. The advantage is the ready availability of a rich set of pre-coded functions, which are possibly updated and optimized in a way transparent to the final user (in this case, the programmer of an application). Given that APIs stay unchanged, the algorithms and their implementations may change without affecting the existing software.

Along the preceding line, *Smart Peephole* was implemented by choosing suited API procedure calls from Microsoft Cognitive Services (MCS). Such services are part of Microsoft Azure[1] [9], a collection of cloud computing integrated services that can be used by developers to create, deploy, and manage applications across a network. Azure can be distributed by connecting the cloud environment and the local environment through consistent hybrid cloud capabilities and using open source technologies to achieve maximum portability. While it was born as a solution for enterprise settings [10], it can also be used on a smaller scale for consumer applications like the one presented here. API method calls were really efficient and demonstrated a high responsiveness during the testing phase, thus leading us to exploit them to build *Smart Peephole*. In particular, the application exploits MCS Faces API, which identifies and recognizes faces, Speech Recognition API, which recognizes voice, and an additional module that interfaces with Emotion API. The last mentioned module was added to interact with the person standing at the doorstep in a "more human way". This is obtained by catching up with the person's mood and deciding whether it is worth to improve it if his/her emotions are negative (anger, sadness, fear, disgust, contempt), or to favor it if his/her emotions are positive (happiness, neutral, surprise). All of the core biometric functions planned for the application now are reduced only to API function calls that execute on highly efficient machines requiring a minimal temporal cost. In fact, the temporal costs are mostly reduced to the transmission of the algorithmic request to the server and the arrival of the response: an error code or the answer related to request. The reverse of the medal for this service (in practice Biometrics As A Service) is the cost due for full functionality and storage capacity. However, this is often the case for commercial cloud-distributed resources. Moreover, given the commercial

---

[1] https://azure.microsoft.com/it-it/.

nature of the services, the algorithmic details are not available. In any case, they are out of the scope of this work.

## 2.2   Microsoft Cognitive Services

Microsoft Cognitive Services work across devices and platforms such as iOS, Android and Windows, keep improving and are easy to set up. The framework contains a lot of APIs to facilitate programmers dealing with difficult problems by exploiting remote calls in a remote access machine. The APIs used for the implementation of *Smart Peephole* are: Face API, Speech Recognition API and Emotion API which we describe in the following.

**Face API** - The Face API has two main functions: face detection with attributes and face recognition. The first function detects up to 64 human faces in the same image with high precision location. The input can be specified by a file name or by a valid URL. The face rectangle (left, top, width and height) indicating the face location in the image is returned along with each detected face. Optionally, the face detection function extracts a series of face related characteristics such as pose, gender, age, head position, facial hair and glasses. Four face recognition functions are provided: face verification, finding similar faces, face grouping, and person identification. At present *Smart Peephole* exploits the Face Identification function. It can be used to match a probe subject against a people database (a *person group*) which must be created in advance and can be edited over time. Each group may contain up to one thousand person objects. Each person can have one or more faces registered. Later on, identification can be carried out against the created *person group*. If the probe face is identified as a person object in the group, the person object with be returned.

**Speech Recognition API** - A voice has unique characteristics that can be used to identify a person, just as in the case of a fingerprint. Using voice for access control and authentication scenarios has emerged as a tool that can be used to offer a level up in security, that also simplifies the authentication experience for customers. MCS offer Speaker Recognition following two modalities: speaker verification and speaker identification. *Smart Peephole* only exploits the Speaker Verification API. It can be used to automatically verify and authenticate users using their voice or speech. The implemented approach is text dependent, i.e. the user is asked to pronounce a predetermined text. The Speaker Verification is divided into two submodules as follows:

(1) Enrollment - each enrolling speaker has to choose a specific pass phrase to use during both enrollment and verification phases. Thus, the speaker's voice is recorded saying a specific phrase, then a number of features are extracted that will be used to recognize the subject to match by the chosen phrase. The extracted voice features for the chosen phrase form a unique signature.

(2) Verification - an input spoken phrase is compared against an enrollment voice signature and phrase in order to verify whether or not they are from the same person, and if this person is saying the correct phrase. In *Smart Peephole*

voice and text verification is used in cascade after face identification, so that only the template of the identified person is verified by speech.

**Emotion API** - The Emotion API takes a facial expression in an image as an input, and returns the confidence across a set of emotions for each face in the image, as well as bounding boxes for such faces, using the Face API. If the application has already called the Face API, it can submit the face rectangle as an optional input. The emotions detected are the following: anger, contempt, disgust, fear, happiness, neutral, sadness, surprise.

The emotions are understood to be cross-culturally and universally communicated with particular facial expressions, even if with a different level of sharpness. The Emotion API uses world class machine learning techniques to provide the results. In interpreting them, the emotion detected should be interpreted as the one with the highest normalized score. Users may choose to set a higher confidence threshold within their application depending on their needs. The Emotion API works also with videos. However, this functionality is out of the scope of *Smart Peephole* and consequently of this paper.

## 3    Architecture of *Smart Peephole*

The `Peephole Module` is the core of the proposed application, which interacts with every other module, exploiting their features to the maximum extent in order to achieve system requirements. Among its functions, it is worth mentioning movement detection, from which a recognition action is triggered. The module uses Dense Optical Flow (DOPTFlow) based on Gunner Farnebäck's algorithm [4] that computes the optical flow for all the points in the frame. In this way, setting a threshold chosen empirically, it is possible to infer when a person is in front of the peephole. Optical flow is the pattern of apparent motion of image objects between two consecutive frames caused by the movement of objects or camera. It is a 2D vector field where each vector is a displacement vector showing the movement of points from the first frame to the second one. The optical flow is always calculated from the previous frame to the current one, making the algorithm resistant to the light changes during the day. Figure 1 shows an example of the movement caught by the vectors during the execution of the DOPTFlow algorithm. If the `Peephole Module` recognizes some movement, face detection is triggered, using the popular Viola-Jones algorithm [13] implemented in OpenCV. It is worth noticing that this task is carried out locally, in order to improve performance and avoid a bottleneck to the MCS Servers to compute face detection and most likely respond with a negative result. To improve the overall efficiency, if a face is detected, eyes are searched for. This helps decreasing false positives and executing a call to the API only if the system is confident enough to have retrieved a face-like contour.

The `Face Recognition Module` exploits the Face Identify function from MCS. It identifies unknown faces from a person group. For each face in the set of detected faces, it will compute similarities between the query face and all the faces in the enrolled person group, and will return candidate person(s) for

**Fig. 1.** An example of result from DOPTFlow algorithm.

face(s) ranked by similarity confidence. The algorithm allows more than one face to be identified independently at the same request, but no more than 10 faces. Identification works well for frontal faces and near-frontal faces, but, given the application context, this is not a true limitation. As a matter of fact, people that in front of a peephole do not maintain such pose are likely to be trying to avoid recognition, and therefore they could be automatically marked as *intruders*.

The `Speech Recognition Module` uses two main methods, Speech Enrollment and Speech Verification, whose functions rely on Verification Profile - Create Enrollment and Speaker Recognition - Verification APIs, respectively. The enrollment supported by Speaker Recognition - Create Enrollment API is text-dependent: the speaker needs to choose a specific phrase to use in both enrollment and verification. The list of possible phrases is found by making another API call to the function Verification Phrase - List All Supported Verification Phrases. The service requires at least three enrollment captures for each speaker before the profile can be used in verification scenarios. It is also recommended to use the same device in both enrollment and verification. The audio length should be at least 1 s and no longer than 15 s.

The Speaker Recognition - Verification has the same technical requirements for the audio file format and length. During verification, a *confidence* level is returned associated with the verification result. Usually, a *Low* or *Normal* confidence level could mean that the person is not speaking with the same timbre of voice or the device he/she is using to authenticate is not working properly.

In *Smart Peephole* application, the `Emotion Recognition Module` processes the face images resulting from the previous detection step. According to the response, the execution flow of the system could vary producing a different action for each different response code. At present, playing of a different song while entering the door is associated to each emotion. The main problem encountered during the implementation is the ambiguous way the humans show their emotion, according to the level of expression. For example, fear can be misunderstood with sadness, happiness with neutrality, anger with disgust and so on. Hence, sometimes the module might produce false positives. A related problem is that users may have to accentuate their facial expressions to allow recognizing the emotion flawlessly an unambiguously.

# 4   Experimental Results

This section presents experiments for two out of three modules: Face and Emotion Recognition.

## 4.1   Experiments on Face Recognition

For the face recognition module two different experiments were carried out in order to confront them with each other. The first one was based on a subset of the LFW (Labeled Faces in the Wild)[2] dataset [5] consisting of 55 persons (where 5 persons are genuine and the rest are considered as impostors) each of them with a single photo (of course the test photo of the genuine persons is different from the photos used to register them in the system). The second experiment dealt with testing the aforementioned module with high quality pictures coming from the same device. This last dataset, referred to as HD-MCS[3], includes 3 pictures of 15 different persons with 3 different head positions (straight, half-left and half-right). Two of these persons are considered as genuine users and, thus, are registered in the system. The others are used to estimate the expected threshold using False Acceptance (FAR) and False Rejection (FRR) curves.

Figure 2 shows The FRR/FAR plot for the subset of LFW. Equal Error Rate (EER) of 0.02 is achieved with a threshold of 0.61. It is worth noticing that the reduced LFW set includes problematic pairs of subjects, e.g., pictures of the Olsen twins (Mary Kate was considered as a genuine user and Ashley as an impostor).

Figure 3 shows an optimal EER of 0.0 achieved with a threshold of 0.5. It is possible to observe that the photo quality does not impact enormously on the setting of the recognition threshold (0.61 in the first experiment and 0.5 in this experiment) but the achieved accuracy does.
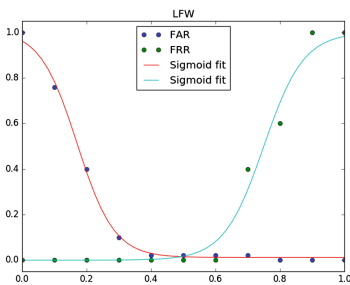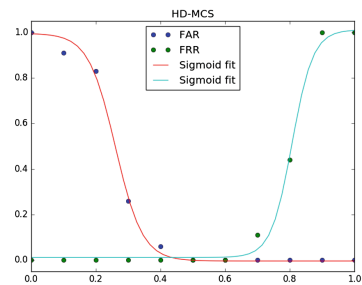


**Fig. 2.** Result for LFW subset.



**Fig. 3.** Results for HD-MCS.

---

## 4.2    Experiments on Emotion Recognition

The Karolinska Directed Emotional Faces (KDEF)[4] [8] is a set of 4900 pictures
of human facial expressions of emotion. The dataset was collected in 1998 at
Karolinska Institutet, Department of Clinical Neuroscience, Section of Psychol-
ogy, Stockholm, Sweden. Participants were 70 amateur actors, 35 females and
35 males, between 20 and 30 years of age. No beards, mustaches, earrings or
eyeglasses, and preferably no visible makeup was allowed during photo session.
Each subject reproduced 7 natural but strong and clear expressions: neutral,
happy, angry, afraid, disgusted, sad, surprised. Face images were captured at
5 different angles: $-90, -45, 0, +45, +90$ degrees, i.e., full left profile, half-left
profile, straight, half-right profile, full right profile, in two sessions.

The original KDEF dataset contains also images of full right and full left
faces, however, because the proposed system only recognizes straight, half-left
and half-right face positions, the original dataset was reduced to contain only
straight, half-left/ right pictures. In this manner, the test dataset contains 140
photos for each of the 7 emotions (leaving out the contempt emotion). Contempt
was not considered because it is not included in the set of the six basic emotions
of anger, disgust, fear, sadness, happiness and surprise, but is rather considered
a mix of anger and disgust [12]. The following Tables 1, 2, and 3 summarize
the results of each orientation of face in the pictures, half-left, half-right and
frontal, respectively, for the 7 emotions. The results reported on the tables are
represented in percentages, where the columns represent the emotion detected by
the system and the rows represent the true emotion of the face image. The last
column of the table represents the percentage of photos rejected by the system.
Tables 1 and 2 show that the number of possibly unclassified expressions is about
the same in the two half profile poses. All tables show that Anger, Disgust and
Fear are the most problematic expressions. Table 3 (frontal pose) also includes
the results corresponding to the emotion of contempt. It was introduced because
the system sometimes fails in recognizing contempt instead of anger or fear.

**Table 1.** Confusion matrix for emotion recognition in images were face is half-left.

|  | Anger | Disgust | Fear | Sadness | Happiness | Surprise | Neutral | Error |
|---|---|---|---|---|---|---|---|---|
| Anger | 10.7143 | 4.2857 | 0 | 3.5714 | 0 | 0 | 81.4286 | **0** |
| Disgust | 7.1426 | 57.8571 | 0 | 16.4286 | 2.8571 | 0 | 14.2857 | **1.4286** |
| Fear | 1.4287 | 2.1429 | 3.5714 | 15 | 7.1429 | 36.4286 | 34.2857 | **0** |
| Sadness | 0 | 0 | 0 | 63.5714 | 1.4286 | 0 | 34.4286 | **0** |
| Happiness | 0 | 0 | 0 | 0 | 100 | 0 | 0 | **0** |
| Surprise | 0 | 0 | 0 | 0 | 0.7143 | 85 | 14.2857 | **0** |
| Neutral | 0 | 0 | 0 | 0 | 0 | 0 | 99.2857 | **0.7143** |

---

[4] http://www.emotionlab.se/resources/kdef.

**Table 2.** Confusion matrix for emotion recognition in images were face is half-right.

|           | Anger   | Disgust | Fear   | Sadness | Happiness | Surprise | Neutral | Error |
|-----------|---------|---------|--------|---------|-----------|----------|---------|---------|
| Anger     | 17.1429 | 2.1429  | 0      | 0.7143  | 0.7143    | 0.7143   | 78.5714 | **0** |
| Disgust   | 13.5714 | 53.5714 | 0      | 13.5714 | 2.1429    | 0        | 17.1429 | **0** |
| Fear      | 1.4286  | 2.1429  | 2.8571 | 13.5714 | 4.2857    | 36.4286  | 39.2857 | **0** |
| Sadness   | 0       | 0       | 0      | 51.4285 | 0.7142    | 0        | 47.8571 | **0** |
| Happiness | 0       | 0       | 0      | 0       | 99.2857   | 0        | 0.7143  | **0** |
| Surprise  | 0       | 0       | 0      | 0       | 0.7143    | 78.5714  | 20      | **0** |
| Neutral   | 0       | 0       | 0      | 0       | 0         | 0        | 99.2857 | **0.71423** |

**Table 3.** Confusion matrix for emotion recognition in images were face is frontal.

|           | Anger  | Disgust | Contempt | Fear    | Sadness | Happiness | Surprise | Neutral | Error |
|-----------|--------|---------|----------|---------|---------|-----------|----------|---------|-------|
| Anger     | 55     | 2.1429  | 2.8571   | 0.7143  | 2.8571  | 0         | 0.7143   | 35.7143 | **0** |
| Disgust   | 7.8571 | 71.4286 | 0        | 0       | 14.2857 | 2.1429    | 0        | 4.2857  | **0** |
| Fear      | 1.4286 | 0.7143  | 2.1429   | 17.142  | 20.7143 | 4.2857    | 42.8571  | 10.7143 | **0** |
| Sadness   | 0      | 0       | 0        | 0       | 86.4286 | 0.7143    | 0.7143   | 12.1429 | **0** |
| Happiness | 0      | 0       | 0        | 0       | 0       | 100       | 0        | 0       | **0** |
| Surprise  | 0      | 0       | 0        | 0       | 0       | 0.7143    | 95.7143  | 3.5714  | **0** |
| Neutral   | 0      | 0       | 0        | 0       | 0       | 0         | 0        | 100     | **0** |

The fact that the emotion of anger is swapped by the contempt one is not surprising. On the other hand, confusing fear with contempt implies a critical error in the emotion detection system.

## 5    Conclusions

This paper has presented *Smart Peephole*, a biometrics security system based on cloud services (biometrics as a service), that has the capability to observe incoming subjects at the house doorstep. It plays the role of a classic door peephole taking decisions to open or not the door for a certain person. The system is composed of four main modules: Peephole module, Face detection module, Speech verification module, and Emotion detection module, that rely on APIs provided by Microsoft Cognitive Services. The tests carried out for this project highlight that the system has optimal capability of distinguishing enrolled persons from the so called intruders. Accuracy is optimal upon recognition of a frontal or near-frontal person picture, which is a normal situation given the kind of application. Regarding the future works, we would like to integrate *Smart Peephole* into an overall smart ambient project, to construct a smart environment system which increases the house awareness with respect to the habits of its inhabitants. This would imply the implementation of a smart algorithm that, upon recognizing a certain person, remembers the recognized person' s habits

and manners around the house and tries, with respect to these habits, to make him/her feel as comfortable as possible, trying at the same time to diminish energy consumption.

# References

1. Ahonen, T., Hadid, A., Pietikainen, M.: Face description with local binary patterns: application to face recognition. IEEE Trans. Pattern Anal. Mach. Intell. **28**(12), 2037–2041 (2006)
2. Jain, A.K., Ross, A., Prabhakar, S.: An introduction to biometric recognition. IEEE Trans. Circ. Syst. Video Technol. **14**(1), 4–20 (2004)
3. Clarke, R.: Human identification in information systems: management challenges and public policy issues. Inf. Technol. People **7**(4), 6–37 (1994)
4. Farnebäck, G.: Two-frame motion estimation based on polynomial expansion. In: Proceedings of the 13th Scandinavian Conference on Image Analysis (2003)
5. Huang, G.B., Ramesh, M., Berg, T., Learned-Miller, E.: Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07–49, University of Massachusetts, Amherst (2007)
6. Kaehler, A., Bradski, G.: Learning OpenCV 3: Computer Vision in C++ with the OpenCV Library. O'Reilly Media Inc., Sebastopol (2016)
7. Muda, L., Begam, M., Elamvazuthi, I.: Voice recognition algorithms using Mel Frequency Cepstral Coefficient (MFCC) and Dynamic Time Warping (DTW) techniques. J. Comput. **2**(3) (2010)
8. Lundqvist, D., Flykt, A.,Ohman, A.: The karolinska directed emotional faces (KDEF). CD ROM from Department of Clinical Neuroscience, Psychology section, Karolinska Institutet (1998)
9. Microsoft: Truly consistent hybrid cloud with microsoft azure. Website (2017), https://azure.microsoft.com/mediahandler/les/resourceles/bf2fe090-ec7c-4463-92e7-92501d86dd28/Truly%20Consistent%20Hybrid%20Cloud%20with%20Microsoft%20Azure.pdf
10. Nickel, J.: Mastering Identity and Access Management with Microsoft Azure. Packt Publishing Ltd., Birmingham (2016)
11. Ojala, T., Pietikäinen, M., Mäenpää, T.: Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. IEEE Trans. Pattern Anal. Mach. Intell. **24**(7), 971–987 (2002)
12. TenHouten, W.D.: A General Theory of Emotions and Social Life. Routledge, London (2006)
13. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. In: Proceedings of IEEE Conference on Computer Vision and Pattern Recognition - CVPR 2001, pp. 511–518 (2001)
14. Stallings, W., Brown, L.: Computer Security: Principles and Practice, 3rd edn. Pearson Education, London (2015)