# A Cloud-based Voting System for Emotion Recognition in Human–Computer Interaction

Ján Magyar*, Gergely Magyar† and Peter Sinčák‡
*Department of Cybernetics and Artificial Intelligence*
*Technical University of Košice*
Košice, Slovakia
Email: *jan.magyar@tuke.sk, †gergely.magyar@tuke.sk, ‡peter.sincak@tuke.sk

*Abstract*—The growing number of computer systems directly interacting with users led to a desire to equip these systems with the capability to understand human social input such as emotions. Research in affective computing defined different representations of human emotional states for machines and in recent years cloud-based emotion recognition services emerged that make the process of system development easier. In this paper, we provide an overview of the current state of emotion recognition in human–machine interaction and describe a voting-based cloud system that uses commercial services to give accurate emotion recognition results. To train the model, we used images from the Karolinska Directed Emotional Faces and Radboud Faces Database datasets. The results were processed in Microsoft Azure's Machine Learning Studio using a multiclass neural network. During the evaluation of our system, we observed a nearly 20% improvement in overall accuracy in comparison to commercial emotion recognition services.

*Index Terms*—affective computing, cloud computing, cloud service, emotion recognition, human–machine interaction, voting system

## I. INTRODUCTION

Emotions play an integral role in human communication. They induce topic changes, support response formulation and decision making [1] [2] [3] [4]. Emotion recognition is innate for humans, and as machines become more indispensable to everyday life, equipping them with similar capabilities has become increasingly desirable.

Emotion recognition depends on multiple factors and uses many modalities. Humans can perceive these together, but machines must split the task into easier subtasks, for example based on the transmission channel of the expressed emotion e.g. speech, gestures or facial expressions (which are one of the most widely used because they provide universality across cultures [5]).

In the last years, facial emotion recognition has appeared in almost every field of research and business involving human–machine interaction (HMI, often referred to as human–computer interaction) [6], which studies the process of information exchange between humans and machines [6], and its goal is to find methods to simulate the natural human perception of interactions in machines. In HMI, emotion recognition and emotion assessment are used interchangeably, however, we must emphasize that machines only estimate the emotions of the observed human. Facial emotion assessment

in machines provides additional information that can be used to enhance interaction with users.

A modern branch of computer science, affective computing, enabled the development of systems that take into consideration the emotions of their users. The study of human affect in HMI originated with Rosalind Picard's paper [7]. She defined affective systems that can recognize, interpret, process, and stimulate human affect. She also described two major areas of emotion-based human-computer interaction: detection and recognition of emotional information, and machine emotions.

In HMI, the communication channels are clearly stated and easy to control for the machine. In human interaction, two channels have been distinguished: explicit and implicit. The explicit channel is the information channel of the interaction, while the implicit channel transmits information about the speakers. Understanding the other party's emotions is one of the key tasks associated with the implicit channel [7] [8]. It often heavily influences human communication by means of emotional state changes which can result in unexpected or even inappropriate responses. Therefore, for successful human–machine interaction the implicit information should be considered, with facial emotion recognition as its fundamental part.

In this paper we present a cloud-based emotion recognition system for HMI which combines the outputs of commercially available services and creates a more accurate facial emotion assessment. Section II covers an overview of emotion models used in HMI, Section III describes the tools and voting system used in this work, Sections IV and V explain the architecture of our system and the results achieved. Section VI concludes the paper and outlines our future work.

## II. EMOTION MODELS

The word emotion generally refers to the affective aspect of consciousness, a state or feeling, or a conscious mental reaction towards an object accompanied by behavioral and physical changes [7]. The field of affective neuroscience differentiates a number of further expressions when talking about emotions [8]:

- *feeling* is a subjective representation of emotions, unique to the individual;
- *mood* is a diffuse affective state that lasts longer than emotions and is usually less intense;

- *affect* is an encompassing term, used to describe emotions, feelings, and moods together. It is commonly used interchangeably with emotions.

Of the above-mentioned three concepts, affect is the most suited for machine processing, because it has a general meaning and the best application potential in human–centered systems, such as the ones used in human–machine and human–robot interaction.

In general, emotions (or affects) consist of three main components [9]:

- *cognitive*, the inner impression of what we experience;
- *physiological*, not easily controllable bodily changes occurring in response to a person's internal state change;
- *behavioral*, external manifestations of our experience: attitude, facial expressions, human instincts.

In our work, we focus on the behavioral component of emotions, more specifically on facial expression. The representation of recognized affects has been a subject of research for a long time, with a number of different emotion models emerging to enable affect representation in machines.

### A. Emotion models

Humans can clearly recognize the emotions of others. This apparent ease of recognition has led to the identification of a number of basic emotions, which are thought to be universal among all cultures. Some experts have questioned this understanding of emotions, and over the years, new aspects of affects have been defined [10].

We can divide emotional models into three groups:

1) discrete emotion models
2) dimensional emotion models
3) hybrid emotion models

### B. Discrete emotion models

In discrete emotion theory, all humans are thought to have an innate set of basic emotions that are cross-culturally recognizable. These basic emotions are described as "discrete" because they are believed to be distinguishable by an individual's facial expression and biological processes. This means that we should be able to tell what emotion a person is feeling by looking at his or her brain activity and/or physiology, without knowledge of the larger context of the eliciting event [11].

A popular example of a discrete emotion model is Paul Ekman and his colleagues' cross-cultural study [11], in which they concluded that the six basic emotions are anger, disgust, fear, happiness, sadness, and surprise. Ekman argued that there are particular characteristics attached to each of these emotions, allowing them to be expressed in varying degrees [12]. Each emotion acts as a discrete category rather than an individual emotional state. Moreover, with additional ontologies, it is possible to extend the existing model with simple relations to further refine the available information (see Fig. 1). The emotion recognition services used in this study all use this model to describe the emotion recognized on a human face.
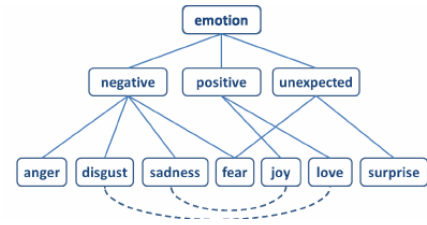


Fig. 1. Emotion ontology for the six Ekman emotions (plus 'Love'). Solid lines indicate inheritance, dashed lines indicate opposites [13]
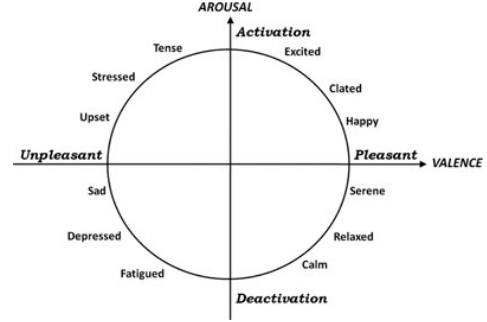


Fig. 2. A graphical representation of the 2D circumplex model of affect with the horizontal axis representing the valence dimension and the vertical axis representing the arousal or activation dimension [20]

### C. Dimensional emotion models

The discrete emotion theory was heavily criticized throughout the years because of the lack of correspondence between emotions and brain activity, variability in facial expressions and behavior [14]. Newer theories suggest that emotions are highly intercorrelated [15] [16]. An extensive and detailed study of the intercorrelations among emotional experiences has yielded two-dimensional models of affective experience [17]. One of the more widely used representations of such 2D models is the circumplex model of affect (see Fig. 2) which suggests that all affective states arise from two independent neurophysiological systems: valence and arousal [18]. They represent the model's axes and refer to connotation (negative or positive) and intensity (low or high) of the felt emotion [19]. According to this model, every affective experience is the consequence of a linear combination of these two independent systems.

### D. Hybrid emotion models

Other theories combine discrete and dimensional models, in their large variety, it is worth mentioning two of them [9]. The Diener-Smith-Fujite model distinguishes among discrete emotions after having discerned between two groups characterized by positive and negative valence [21]. Another broadly used example is Plutchik's model, a three-dimensional model that is a hybrid of both basic-complex categories and dimensional theories [22]. It arranges emotions in concentric circles where inner circles contain basic emotions and other circles contain more complex emotions. These complex emotions can be acquired by a mixing of emotions inspired by the theory of colors as explained in [23].

## III. Used Tools and Methods

Our study's aim was to create a cloud service providing more accurate emotion recognition than free commercial solutions. The service is based on a voting system that uses the results of commercial emotion recognition APIs to calculate a weighted probability for the presence of various emotions. We focused on the recognition of anger, contempt, disgust, fear, happiness, neutral, sadness, and surprise. To train the voting system model, a database of 2,581 images was created. On this sample, we evaluated the performance of some emotion recognition APIs, and selected four to use in our service. The voting system was implemented as a multiclass neural network model in Microsoft Azure Machine Learning Studio (ML Studio) [24]. In this section, we describe the image database that was used, provide a short overview of the emotion recognition APIs we evaluated, and explain the concept of a voting system.

### A. Image database

For evaluating commercial emotion recognition APIs and to train our model, we selected images from two freely available datasets: the Karolinska Directed Emotional Faces (KDEF) [25] and the Radboud Faces Database (RaFD) [26].

The KDEF was developed in 1998 at Karolinska Institutet and it contains 4,900 images of 70 individuals displaying 7 emotional expressions (anger, disgust, fear, happiness, neutral, sadness, surprise), each expression photographed twice from 5 different angles. We selected only frontal images, two per each emotion per individual, except for one individual for whom only one set of images was used.

For the development of RaFD, pictures of 67 models were taken. The models expressed eight emotions, contempt in addition to the ones present in KDEF. Photographs were taken simultaneously from five different angles for three different gaze directions of the model. We selected frontal images for each emotion expression of every model for all gaze directions. This resulted in a dataset of a total 2,581 images: 973 from KDEF and 1,608 from RaFD.

### B. Emotion recognition APIs

In recent years, a number of emotion recognition services were made publicly available. Due to the complexity of emotion recognition, these services make the development of systems incorporating emotion recognition easier, and are capable of exploring additional features of faces, not only emotion. After the user uploads an image or provides an URL address to it, the services return information about the detected features; the results of emotion recognition are usually provided in the form of number values which represent the probability of the given emotion's presence in the image. For this study, we selected emotion recognition APIs based on two criteria: 1) the API must be freely available (with a limited number of calls or for a limited period of time), and 2) the API must use a discrete emotion model to assess emotions. We selected eight services meeting these requirements for testing.

The **Affectiva SDK** [27] and **Affectiva Emotion as a Service** [28] detect seven emotions: anger, contempt, disgust, fear, joy, sadness, and surprise. The API is not available for free, an academic license can be acquired for research purposes, while the SDK can be used for free for a limited period of time. For evaluating the performance of Affectiva's emotion recognition over our image database, we used the Affectiva SDK.

Amazon's **Rekognition** [29] is capable of detecting seven emotions (anger, calmness, confusion, disgust, happiness, sadness, and surprise). The category Unknown is returned if the API fails to detect any of the emotions listed. Unlike other APIs, Rekognition provides results only for the emotions with the three highest confidence values.

**Face++ Cognitive Services' Emotion Recognition** [30] is currently in beta with support for the detection of anger, disgust, fear, happiness, neutral, sadness, and surprise. The free API key can be used forever with a limitation of one call per second.

**Google Vision** [31] can detect four emotional states on the face, namely anger, jow, sorrow, and surprise. The probability values for each emotion aren't from the range 0–1, they are the discrete values unknown, very unlikely, unlikely, possible, likely, very likely. In our tests, we converted these values to numeric values in equal intervals with very unlikely being 0 and very likely 1 (the value unknown was represented with NULL). The service is free of charge for up to 1,000 calls per month.

The **Kairos API** [32] detects anger, disgust, fear, joy, sadness, and surprise. The service is available for free with limited usage: 25 API transactions/minute are possible with a daily cap of 1,500.

The **Microsoft Face API** [33] is available as part of the company's cognitive services. It can detect eight emotions: anger, contempt, disgust, fear, happiness, neutral, sadness, and surprise. In its free version, the API key can be used for 30,000 transactions with 20 per minute.

The **F.A.C.E. API** [34] by Sightcorp detects anger, disgust, fear, happiness, sadness, and surprise. The API is available for a two-week long free trial with a maximum of 5,000 calls and a 1 call/second cap.

The **Sighthound Cloud API** [35] detects anger, disgust, fear, happiness, neutral, sadness, and surprise. The API can be used for a trial period with 5,000 calls.

### C. Voting system

The above-mentioned emotion recognition APIs show a reasonably good accuracy for the recognition of common emotions, but their performance drops for more complex emotions. This might be due to the lower level of expressiveness of some emotions, or the nature of the dataset that was used to train the API's model which can show a bias towards common emotions (e.g. more images depict happiness than contempt). If an application is expected to perform at the same level of accuracy across all emotions, alternative solutions must be used.
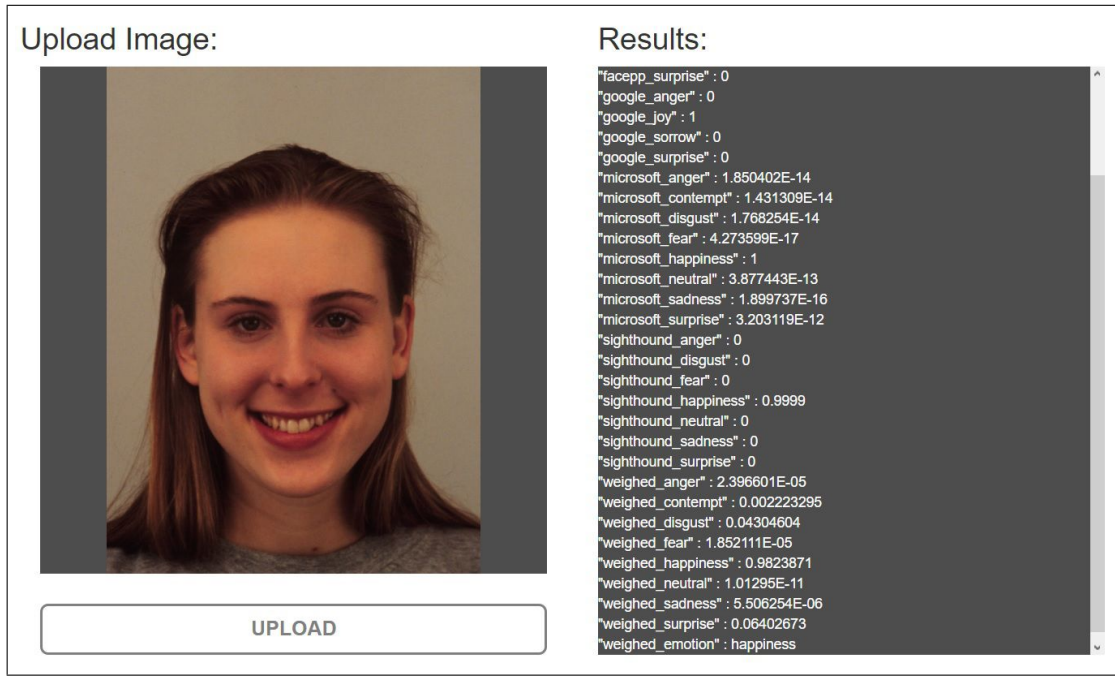
Fig. 3. User interface of web application with uploaded image and results displayed

To tackle this problem, multiple emotion recognition APIs can be consulted. This raises a further problem when there is a discrepancy between the emotion recognition results. For example, if one API detects sadness and the other one contempt, which result should be considered to be correct? A voting system solves this problem by taking the output of APIs and assigning weights to them that reflect the level of expertise of APIs in the recognition of the given emotion. A weighted sum of the probabilities of emotions' presence can provide a higher accuracy.

## IV. CLOUD-BASED EMOTION RECOGNITION VOTING SYSTEM

We designed and implemented a web application hosted on Microsoft Azure. The web application has a simple user interface offering file upload (see Fig. 3). After the image is uploaded, the application sends it to emotion recognition cloud services and gathers their output. This output is then sent to a voting system cloud service deployed through Microsoft Azure Machine Learning Studio that evaluates these results and provides a weighted output for them. The results of the voting system, along with the partial emotion recognition results from the used APIs are then displayed to the user in a key–value pair format. The results are saved in a database and can be used for further training of our emotion recognition voting system. The web application's system architecture is shown in Fig. 4, its workflow is depicted in Fig. 5.

From the APIs mentioned in section III-B, we used Face++ Cognitive Services' Emotion Recognition, Google Vision, Microsoft Face, and Sighthound Cloud API. Affectiva's Emotion as a Service was omitted due to a lack of free subscription, while the F.A.C.E. API's free trial version is available only for
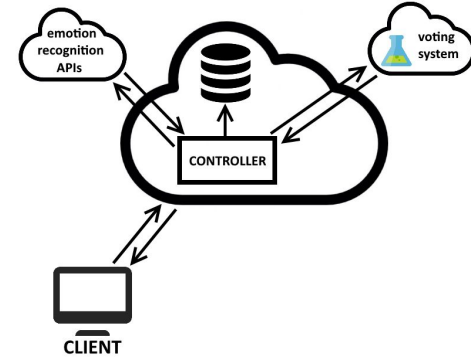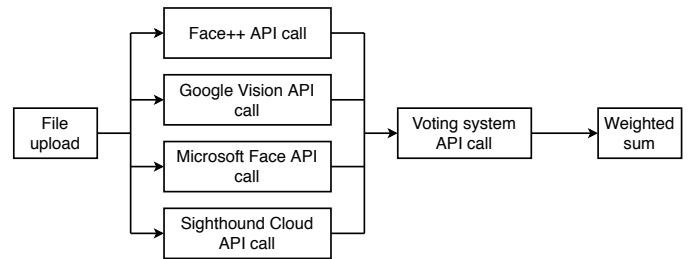


Fig. 4. System architecture



Fig. 5. Workflow of emotion recognition

two weeks. Amazon Rekognition returned probability values only for the three most probable emotions, which made it unsuited for use in a voting system, while the Kairos API not always detected faces on the images from the KDEF and RaFD datasets, and its accuracy was lower compared to other APIs (see Table II).

TABLE I
CONFUSION MATRIX FOR THE PERFORMANCE OF THE VOTING SYSTEM MODEL (A – ANGER, C – CONTEMPT, D – DIGUST, F – FEAR, H – HAPPINESS, N – NEUTRAL, SA – SADNESS, SU – SURPRISE)

| | | Actual class | | | | | | | | Precision (%) |
|---|---|---|---|---|---|---|---|---|---|---|
| | | A | C | D | F | H | N | Sa | Su | |
| Predicted class | A | 303 | 0 | 3 | 3 | 0 | 0 | 6 | 0 | 96.19 |
| | C | 3 | 188 | 0 | 0 | 0 | 0 | 1 | 0 | 97.92 |
| | D | 7 | 0 | 327 | 5 | 0 | 0 | 3 | 0 | 95.61 |
| | F | 4 | 0 | 2 | 298 | 0 | 0 | 7 | 6 | 94.01 |
| | H | 0 | 0 | 0 | 1 | 340 | 0 | 0 | 0 | 99.71 |
| | N | 10 | 13 | 0 | 1 | 0 | 340 | 7 | 2 | 91.15 |
| | Sa | 13 | 0 | 8 | 12 | 0 | 0 | 316 | 0 | 90.54 |
| | Su | 0 | 0 | 0 | 20 | 0 | 0 | 0 | 332 | 94.32 |
| Recall (%) | | 89.12 | 93.53 | 96.18 | 87.65 | 100 | 100 | 92.94 | 97.65 | |

TABLE II
ACCURACY OF EMOTION RECOGNITION APIS (VALUES IN %)

| API | KDEF | RaFD | Overall |
|---|---|---|---|
| Affectiva | 54.57 | 64.68 | 60.87 |
| F.A.C.E. | 63.21 | 61.32 | 62.03 |
| Kairos | 45.74 | 26.99 | 34.06 |
| Rekognition | 52.62 | 39.74 | 44.6 |
| Face++ | 77.08 | 71.33 | 73.5 |
| Google | 43.47 | 36.63 | 39.21 |
| MS Face | 75.33 | 76.24 | 75.9 |
| Sighthound | 62.18 | 72.33 | 68.5 |
| Voting system | 90.75 | 97.08 | 94.69 |

The voting system model was implemented in Microsoft Azure Machine Learning Studio (ML Studio) which is a free tool for creating, training and deploying models as APIs. We created a file containing results from all four used APIs for all the images in our dataset and the emotion present in the images. This file was uploaded to ML Studio and was used to train a multiclass neural network with 100 hidden nodes. For possible learning rates, we set 0.01, 0.02, and 0.04, from which ML Studio selected the most appropriate during training the model. The number of iterations was set similarly to be from the interval 164–500 with number of points of 3. The initial learning weight was 0.1, the momentum 0, and the Min-Max normalizer was used. The neural network had 26 input nodes for the probability of each recognized emotion from four APIs, and had eight output nodes for the weighted sum of probabilities for the detected emotions (anger, contempt, disgust, fear, happiness, neutral, sadness, surprise). The data about the images was split randomly into a 75% training set and a 25% testing set for evaluation. The dataset splitting was not stratified, not enforcing a proportional representation of emotions in the training and testing sets. The performance of the voting-system model over the combined dataset is described in the next section.

## V. EVALUATION AND RESULTS

To evaluate the performance of the APIs mentioned in section III-B as well as our voting-system based model, we tested them on the dataset containing images from KDEF and RaFD. The performance of the classification system is presented in a confusion matrix in Table I. Precision describes the predictive

performance of the model, showing the percentage of correct classification within a predicted class. Recall describes the classification performance and represents the percentage of correctly classified data for any class.

The voting-system based model showed the lowest precision in recognizing the emotions sadness, neutral, fear, and surprise. This is due both to the low level of expressiveness of these emotions and also their similarity to each other (neutral–sadness, fear–surprise). It is important to note that the voting system's results were still considerably better for each emotion than any individual API's, a more detailed emotion-wise comparison of API performance will be the subject of a follow-up paper.

The overall accuracy (number of correct classifications over the size of the dataset) is presented in percentages in Table II for the two individual image datasets and for the combined dataset. We first include the APIs that were not used for our model, then the four APIs used, and finally our model aggregating their results.

Table II shows that our voting system aggregating the results from four APIs performs significantly better than each individual API and the APIs that were not included in the aggregagtion process. The low accuracy of the Kairos API is due to the API not detecting a face in 348 images, while the Google API performed badly because of its use of fewer emotion categories. Compared to the best-performing single emotion recognition API – Microsoft Face API – our voting system showed an almost 20% overall improvement.

## VI. CONCLUSION

This paper describes a cloud-based facial emotion assessment system which uses the outputs of different publicly available emotion recognition services. These data are processed using the Machine Learning Studio developed by Microsoft. The system was tested using a combined dataset of the KDEF and RaFD facial expression databases. The results showed that our system outperformed the existing recognition services. The main advantage of our system is that it can be deployed and used for free in a small-scale scenario, since all the used APIs have a free-tier usage. To use it in a larger scale and send several images in a small amount of time, users must subscribe to the given services.

The next steps in the development of the system involve creating a public API so other researchers could call our service and get more accurate emotion recognition results than the ones from the tested services. This way the resulting system could be used in various applications of affective computing, such as human–robot interaction. By labeling the uploaded images we can create training datasets to improve the accuracy of our model. The voting system based service described in this paper also needs to be tested in real human–machine interaction, which will be the subject of further research.

## References

[1] S. R. Fussell, *The verbal communication of emotions: interdisciplinary perspectives*. Psychology Press, 2002.

[2] K. Oatley, "Two movements in emotions: Communication and reflection," *Emotion Review*, vol. 2, no. 1, pp. 29–35, 2010.

[3] E. Altenmüller, S. Schmidt, and E. Zimmermann, *The evolution of emotional communication: From sounds in nonhuman mammals to speech and music in man*. OUP Oxford, 2013.

[4] T. Lanciano and A. Curci, "Does emotions communication ability affect psychological well-being? a study with the mayer–salovey–caruso emotional intelligence test (msceit) v2. 0," *Health communication*, vol. 30, no. 11, pp. 1112–1121, 2015.

[5] L. De Sonneville, C. Verschoor, C. Njiokiktjien, V. Op het Veld, N. Toorenaar, and M. Vranken, "Facial identity and facial emotions: speed, accuracy, and processing strategies in children and adults," *Journal of Clinical and experimental neuropsychology*, vol. 24, no. 2, pp. 200–213, 2002.

[6] G. Johannsen, "Human-machine interaction," *Control Systems, Robotics and Automation*, vol. 21, pp. 132–62, 2009.

[7] R. W. Picard *et al.*, "Affective computing," 1995.

[8] E. Fox, *Emotion science: An integration of cognitive and neuroscience approaches*. Palgrave Macmillan, 2008.

[9] K. Lochner, "Affect, mood, and emotions," pp. 43–67, 2016.

[10] M. Gendron and L. Feldman Barrett, "Reconstructing the past: A century of ideas about emotion in psychology," *Emotion review*, vol. 1, no. 4, pp. 316–339, 2009.

[11] P. Ekman, "An argument for basic emotions," *Cognition & emotion*, vol. 6, no. 3-4, pp. 169–200, 1992.

[12] P. Ekman and W. V. Friesen, *Facial Action Coding System: Investigatoris Guide*. Consulting Psychologists Press, 1978.

[13] K. Roberts, M. A. Roach, J. Johnson, J. Guthrie, and S. M. Harabagiu, "Empatweet: Annotating and detecting emotions on twitter." in *LREC*, vol. 12. Citeseer, 2012, pp. 3806–3813.

[14] L. F. Barrett, "Variety is the spice of life: A psychological construction approach to understanding variability in emotion," *Cognition and Emotion*, vol. 23, no. 7, pp. 1284–1306, 2009.

[15] J. A. Russell and J. M. Carroll, "On the bipolarity of positive and negative affect." *Psychological bulletin*, vol. 125, no. 1, p. 3, 1999.

[16] D. Watson, D. Wiese, J. Vaidya, and A. Tellegen, "The two general activation systems of affect: Structural findings, evolutionary considerations, and psychobiological evidence." *Journal of personality and social psychology*, vol. 76, no. 5, p. 820, 1999.

[17] R. J. Larsen and E. Diener, "Promises and problems with the circumplex model of emotion." 1992.

[18] J. Posner, J. A. Russell, and B. S. Peterson, "The circumplex model of affect: An integrative approach to affective neuroscience, cognitive development, and psychopathology," *Development and psychopathology*, vol. 17, no. 3, pp. 715–734, 2005.

[19] J. A. Russell, "A circumplex model of affect." *Journal of personality and social psychology*, vol. 39, no. 6, p. 1161, 1980.

[20] P. A. Abhang and B. W. Gawali, "Correlation of eeg images and speech signals for emotion analysis," *British Journal of Applied Science & Technology*, vol. 10, no. 5, pp. 1–13, 2015.

[21] E. Diener, H. Smith, and F. Fujita, "The personality structure of affect." *Journal of personality and social psychology*, vol. 69, no. 1, p. 130, 1995.

[22] R. Pluchik, "A general psychoevolutionary theory of emotions," *Emotion: Theory, research and experience*, vol. 1, pp. 3–33, 1980.

[23] R. Plutchik and H. Kellerman, *Emotion, psychopathology, and psychotherapy*. Academic press, 2013, vol. 5.

[24] "Microsoft Azure Machine Learning Studio," https://studio.azureml.net, accessed: 2018-04-22.

[25] D. Lundqvist, A. Flykt, and A. Öhman, "The Karolinska directed emotional faces (KDEF)," *CD ROM from Department of Clinical Neuroscience, Psychology section, Karolinska Institutet*, no. 1998, 1998.

[26] O. Langner, R. Dotsch, G. Bijlstra, D. H. Wigboldus, S. T. Hawk, and A. Van Knippenberg, "Presentation and validation of the Radboud Faces Database," *Cognition and emotion*, vol. 24, no. 8, pp. 1377–1388, 2010.

[27] "Affectiva SDK," https://www.affectiva.com/product/emotion-sdk/, accessed: 2018-04-22.

[28] "Affectiva Emotion as a Service," https://www.affectiva.com/product/emotion-as-a-service/, accessed: 2018-04-22.

[29] "Amazon Rekognition," https://aws.amazon.com/rekognition/, accessed: 2018-04-22.

[30] "Face++ Cognitive Services - Emotion Recognition," https://www.faceplusplus.com/emotion-recognition/, accessed: 2018-04-22.

[31] "Google Vision API," https://cloud.google.com/vision/, accessed: 2018-04-22.

[32] "Kairos API," https://www.kairos.com/docs/api/, accessed: 2018-04-22.

[33] "Microsoft Face API," https://azure.microsoft.com/en-us/services/cognitive-services/face/, accessed: 2018-04-22.

[34] "F.A.C.E. API by Sightcorp," https://face-api.sightcorp.com, accessed: 2018-04-22.

[35] "Sighthound Cloud API," https://www.sighthound.com/products/cloud, accessed: 2018-04-22.