# An Introduction to Recent Mobile Affective Inferring Techniques: Methods, Applications and Challenges

Ou Changkun

**Abstract**— Affective Computing has been considered as one of the essential aspects of massive human-computer interaction related projects. This paper provides a broad introduction to the recent advances in affective computing for emotion inferring based on mobile techniques. We expand the emotion inferring principles for different mobile commodity sensors, application context and their possible combinations in the recent researches. Then we compare the recent popular machine learning methods and models for these sensors, and highlight the most useful techniques and models for their performance. Our comparisons are not limited to traditional machine learning algorithm; they also include the representation learning models. In the end, we discussed few novel applications based on mobile affective computing techniques, such as how adaptive user interface and usability testing works in an emotion-aware system, as well as current limitations and open challenges of this research area.

**Index Terms**—Mobile Emotion Inferring, Convolutional Networks, Recurrent Networks, Support Vector Machine, Dialogue System, Adaptive User Interfaces

---

## 1 INTRODUCTION

Affective computing is an emerging interdisciplinary research field ranging from cognitive and social sciences to human-computer interaction (HCI) researchers with techniques like computer vision, machine learning, natural language understanding, etc. With the long-term research on emotion theory from psychology and neuro-science[37, 88], emotion has been confirmed to be a significant effect [38] on human communication, decision making and perception.

On the perspective of human-computer interaction, Picard [66] pointed out that affective computing involved projects can be used for *reducing user frustration* enabling comfortable communication of user emotion, developing infrastructure and applications to *handle affective information*, as well as building tools that *help formulate social-emotional skills*.

Recently, the ubiquitous computing[90] and wearable computing[80], which are strictly related to affective computing, have achieved the pervasive attention of scientists. Ubiquitous computing and wearable computing are the necessary products of the combination of mobile computing technology and computer individualization.

Hence, it creates new research opportunities for affective computing to inferring user emotions with the combination of smart mobile wearable devices. We simply call it *Mobile Affective Inferring*. Such devices have been widely used over the world. The key feature of smart devices is the abundant sensors that enable various affect-related signals unobtrusive monitoring. Exploring the possibility of using smart mobile devices for affective computing will benefit at least three perspectives by the potential of long-term unobtrusive monitoring users affective states: First, it influences the affect-related research literature by the wild, natural and unobtrusive study; Second, it establishes the spontaneous affect databases efficiently to evaluate new effective methods, models, systems more accurately; and third, it enhances the user-centered HCI design for the future ubiquitous computing environments.

In this paper, we present an introduction to the mobile affective computing techniques, our next section discusses the exists data source from mobile devices, and Section 3 illustrate the recent advances for each different type of data and present the state-of-the-art model and method. Next, Section 4 based on the previous information, we reasonably assume we have finished user emotion inferring stage, then

---

- *Ou Changkun is studying Human-Computer Interaction at the University of Munich, Germany, E-mail:* `hi@changkun.us`

gives two typical HCI applications in this field. At last, Section 5 and 6 discusses the current challenges of mobile affective inferring and our conclusion of this introduction. Figure 1 shows the hierarchically-structured taxonomy of this paper.
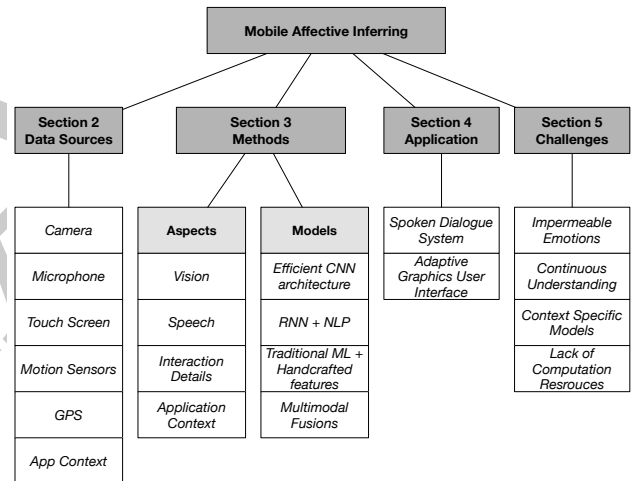


Figure 1. Hierarchically-structured taxonomy of this paper.

## 2 DATA SOURCES

Previous surveys such as [1, 93, 70, 23] put forward that different data sources should be applied to various modeling methods in multimodal affective computing. [70] also give the argument that 90% literature consider visual, audio and text information as multimodal affect analysis instead of other dimensions by their extensive literature review.

In this section, we discuss the commonly existing data sources in a typical smartphone and prepare for the modelling method in the next section. Figure 2 illustrates the data flow from sensors to user emotion state.

### 2.1 Camera

We emphasis vision sensors in the first place since face and facial expressions are undoubtedly one of the most critical nonverbal channels used by the human being to convey internal emotion [38, 70]. This part mainly discusses vision sensors, which includes RGB cameras
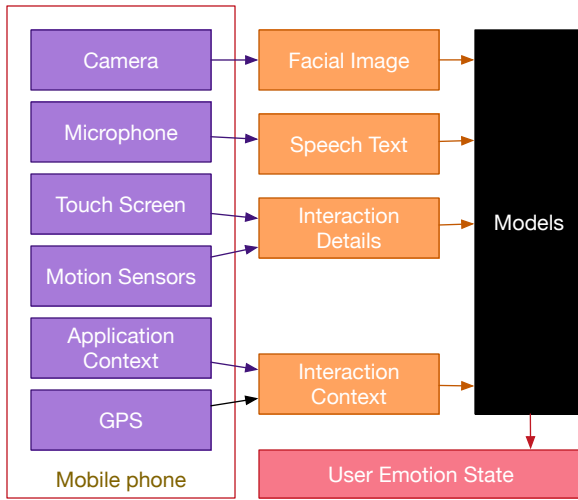
Figure 2. Data sources can be used in emotion inferring which provided by the most of commercial mobile phone devices.
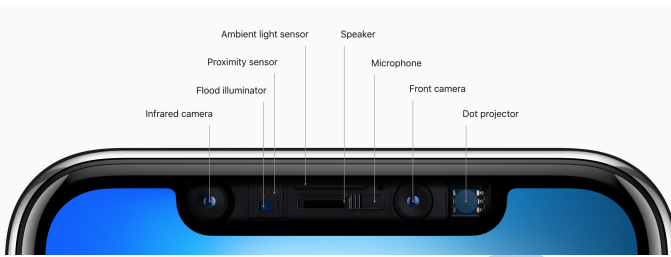


Figure 3. TrueDepth Camera in iPhone X

and depth cameras, and illustrates how vision sensor can be used for affective emotion inferring.

Pure RGB cameras have been widely used in a commercial smartphone as an image sensor. For the camera with depth information on mobile (recently introduced TrueDepth Camera in iPhone X, see Figure 2.1[1]) combines infrared camera, flood illuminator, proximity sensor, ambient light sensor, front-facing camera and dot projector to provide depth images of facial information of a user.

### 2.2 Microphone

Audio sensor usually refers to built-in microphones; it collects voice information from current environments, which can infer user emotions based on their speech contents.

Before recognizing user speech, a system usually should take care and preprocess the environmental noise and detect acoustic fingerprint (i.e., voiceprint) [6] for the current user, isolate their speech from mixed audio information.

Inferring emotions from user speech can split as two part of inferring task. The first part is recognizing the speech text from the user[56, 89], then understanding or inferring from the text, namely sentiment analysis [72].

### 2.3 Touch Screen

Human emotions can be expressed in different ways, emotional communication has focused predominantly on the facial and vocal channels but has ignored the tactile channel[31], which investigated the

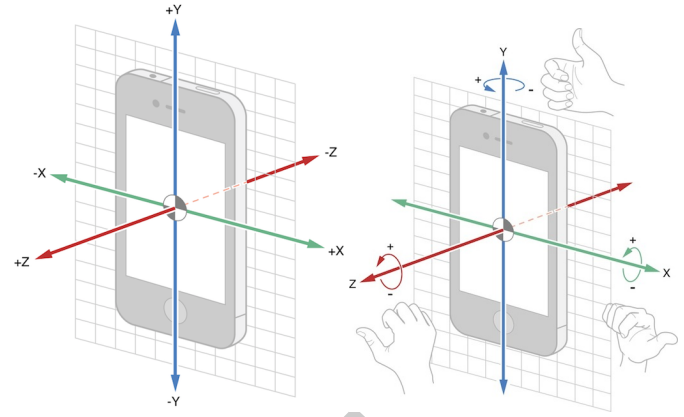[1]https://www.apple.com/iphone-x/#truedepth-camera



Figure 4. Coordinates information of Accelerometer (left) and gyroscope (right) as motion sensors in iOS CoreMotion framework.

possible expressions of user emotion in detail while they are using mobile devices with touch screen.

A capacitive touch screen provides touch position, touch pressure, touch angle through time. Among the subsequent researches[22, 77, 60, 3], researchers explored yield results that human emotions can be inferred by capacitive touch channel in a specific application context based on these features, whose are the existed typical research on emotion inferring only with touch screen interface.

Interestingly, 3D touch screen was introduced in commercial devices a few years ago, some of the researches investigated the possibility of haptic based application [17]. [52, 48] shows a system with haptic touch response essentially can express and influence user emotions. and [4] concludes that haptic-based affect detection remains an understudied topic.

### 2.4 Motion Sensors

Motion sensors typically combine gyroscope and accelerometer, which are yet another interaction detail information[93]. An accelerometer measures proper acceleration (acceleration it experiences relative to free fall), felt by people or objects. Most smartphone accelerometers trade large value range for high precision. The gyroscope can be a handy tool because it moves in peculiar ways and defies gravity. Gyroscopes have been around for a century, and they are now used everywhere from airplanes, toy helicopters to smartphones. A gyroscope allows a smartphone to measure and maintain orientation. Gyroscopic sensors can monitor and control device positions, orientation, direction, angular motion, and rotation. Figure 4 shows the coordinates information of accelerometer and gyroscope sensors.

With these motion sensors, interaction details information such as device holding posture, device moving trajectory can be inferred from these sensor data [60, 70].

### 2.5 GPS

GPS sensors provide geographical information of a user, and it detects the location of the smartphone using 1) GPS[85]; 2) Laceration/Triangulation of cell towers or wifi networks (with a database of known locations for towers and networks)[73]; 3) Location of an associated cell tower or WiFi networks[69].

However, GPS will not work indoors and can quickly kill the battery. Smartphones can try to automatically select best-suited alternative location provider (GPS, cell towers, WiFi), mostly based on desired precision. With the location, we can study the relationship between life patterns and affective states. For example, most people in playground feel happy while most feel sad in a cemetery.

The location provides additional information to verify the subjective report from participants of affective studies. It may also help to build a confidence mechanism [85] for subjective reports. Attaching the location to the subjective report would produce confident weights
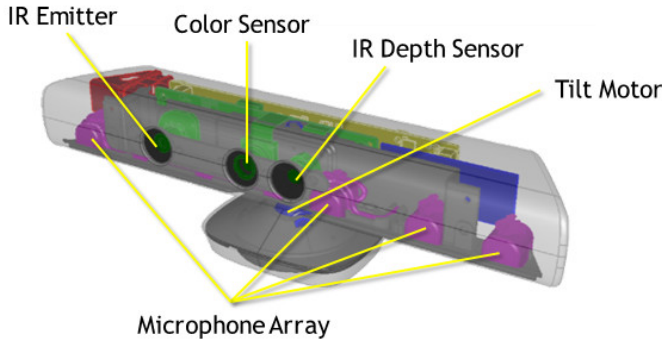
Figure 5. Principle of Microsoft Kinect.

Table 1. Complexity comparison of CNN models, smaller number refers good performance

| Model | Cls. Error (%) | Complexity (MFLOPs) |
|-------|----------------|---------------------|
| AlexNet[46] | 42.8 | 720 |
| VGG-16[78] | 28.5 | 15300 |
| SqueezeNet[36] | 42.5 | 833 |
| MobileNet[34] | 31.6 | 325 |
| ShuffleNet[94] | 31.0 | 292 |
| Xception[10] | 21.0 | 486 |

to measure the significance of collected subjective reports. For example, a participant reported that he was happy in a cemetery. But, in common sense, people in the cemetery would be sad. Thus, we could set a low weight as a confidence value to the report.

### 2.6 Application Context

As we slightly mentioned in the previous section, most of the mobile affective inferring techniques based on a touch screen and motion sensors[22, 77, 60, 3] are based on specific application context, for example, application user interface. It can be another confidence mechanism for inferring system. Vice versa, the inferring systems are only modeling for a specific context. However, this could be a drawback of multimodal emotion inferring since a complete system will integrate more models.

### 3 METHODS

Emotion inferring problems are mostly considered as a classification problem, which classifies three states of user emotion: Happy, Unhappy, Neutral. The reason for this consideration is argued by in technical constraints: more type of emotion we need to classify with more data we need to prepare. In this section, we will see the most accurate specialized method, models and their datasets in different data type aspects mainly all based on machine learning methods.

### 3.1 Vision Aspect

The standard RGB camera brings us focusing on how conduct emotion recognition with RGB images. Through depth camera was recently introduced on a commercial mobile phone, its principle is as same as Microsoft Kinect (see Figure 5). Considering these two different sensor aspects, we dive into two distinct research area on vision sensors.

#### 3.1.1 Plain recognition

Recently, convolutional neural networks (CNN) method has successfully make break-through contributions to computer vision as well as its application to emotion inferring. AlexNet [46] popularized deep convolutional neural networks by winning the ImageNet Challenge (A large-scale image classification challenge). Subsequently, other powerful CNN architecture was proposed such as VGG[78], Inception series[83, 84, 82], ResNet[30], DenseNet[35] and CapsNet[75]. However, the most accurate CNN's usually have hundreds of layers and thousands of channels whereas it is entirely not possible to deploy them to the mobile system. The increasing of mobile emotion inferring needs of running high quality deep neural networks on embedded devices encourage the study on efficient model designs[29].

SqueezeNet [36] is the first model that reduces parameters and computation significantly while maintaining accuracy. MobileNet [34], ShuffleNet [94] and Xception [10] utilizes the depthwise separable convolutions among lightweight models. Table 1 shows the complexity comparison of these CNNs.

The classification of RGB images is just the final step of feedforward propagation. To determine the facial information inside an image
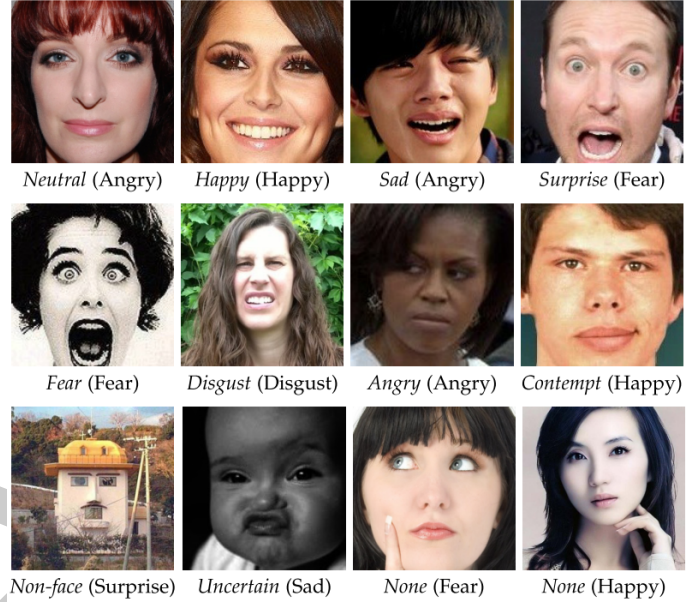


Figure 6. Samples of AffectNet database and classification results. The emotion expression labels is written in parentheses.

essentially become more difficult than only a classification, then the problem refers to Landmark detection. The recent break-through contributions in this area are the Region-based CNN (R-CNN) approach [25] that to bounding-box object detection is to attend to a manageable number of candidate object regions. The state-of-the-art contribution is Mask-R-CNN [28], which serve a conceptually simple, flexible, and general framework for object instance segmentation.

With all these computer vision methods, the ICML 2013 Challenges in Representation Learning introduced the Facial Expression Recognition 2013 (FER-2013) database [26]. Fortunately, human actors/subjects databases portray the basic emotions of external human emotion has been created, which solves the problem of training data.

Benitez-Quiroz et al. [19] proposed EmotioNet database that extracts features by using Gabor filters. Their database is subject-independent and cross-database experiments. [58] uses FER-Wild database, and trained them on AlexNet with noise estimation methods and archived 82.12% accuracy on FER-Wild. AffectNet [57] is the state-of-the-art database that proposed gives the most extensive database of facial expression, valence, and arousal in the wild (see Figure 6). In their paper, various evaluation metrics shows that their deep neural network gives the state-of-the-art performance in facial emotion recognition.

#### 3.1.2 Depth recognition

As we discussed in the previous section, depth camera principle is as same as Microsoft Kinect, the main difference between depth camera and an RGB camera is it provides 3D facial information, which leads the model difference in this field. Unlike Kinect, depth camera in most

cases can only offer facial reconstruction model information instead of body gesture. Thus, depth recognition mainly focuses on modeling 3D facial points.

Paper [9] was the recent research that considering 3D modelling. They propose a real-time 3D model-based method that continuously recognizes dimensional emotions from facial expressions in natural communications. The most challenge parts of their research cover the 3D facial information reconstructed from 2D images. Zhang et al. [95] proposed their exploration on 3D facial points modeling of emotion recognition that directly gets depth information from Kinect. However, their recognition only gives three different state recognition.

Despite the already existing depth cameras in mobile phones, researchers in this area are rare and not in popular demand. We believe that primary reason is 3D modeling requires extensive computation which is not possible from mobile devices at the moment. We will discuss this issue detailly in the section of challenges.

## 3.2   Voice Aspect

Voice aspect as we discussed in the previous section, emotion inferring from user speech is primarily processing user speech. The first part is to deducing user speech text from their voice, and the second part of recognition is calculated sentiment from these documents.

### 3.2.1   Speech recognition

Speech recognition is a board research area, and there exist broad approaches to archive this goal. Previous years commercial systems modeling speech recognition by using Hiden Markov model which archived good performance, however with the raising of deep learning methods, recurrent neural network [56], Long-short term memory cell [33] and attention mechanism [89]. It is laborious to say which is the state-of-the-art model since speech recognition is much more complicated than a typical vision task. Consider there exists very successful commercial system such as Google Speech API [2] can performs stream speech recognition with returning the speech text information. We don't consider this area in detail for the primary goal of mobile affective computing.

### 3.2.2   Sentiment Analysis

Sentiment analysis requires text understanding, and it is not an easy problem to solve. Some machine learning techniques, including various supervised and unsupervised algorithms, are being utilized. Some algorithms rank the importance of sentences within the text and then construct a summary out of essential sentences, others are end-to-end generative models. After we have the speech text from the user, sentiment analysis can be performed to evaluate user emotions for each speech sentence or a chunk of speech contents during a period.

[72] provides board and comprehensively survey on sentiment analysis. We conclude here for the general steps of sentiment analysis:

The approach to extract sentiment from speech text is as follows:

- Tokenize each word via a public sentiment calculation dataset;

- For each word, compare it with positive sentiments and negative sentiments word in the dictionary. Then increment positive count or negative count.

- Finally, based on the positive count and negative count, once can get result percentage of sentiment to decide the polarity.

Sentiment calculation for text essentially clear defined engineering, and the final re sentiment value suits of user speech can be a training feature to different modalities.

## 3.3   Touch Aspect

Touch interaction modality in previous research all considered using a handcrafted feature for touch behavior and using kernel SVM to train linear models for classification. [22] is the first application specific in game, which the recognition rates are very robust even in naturalistic

Table 2. Comparison of classifiers.

| Feature | Unit (%) | |
|---|---|---|
| Deviation in number of strikes | | |
| Deviation in number of taps | | |
| Mode of strike length | Millimeter | |
| Average strike length | Millimeter | |
| Mode of strike speed | Meter/second | |
| Average strike speed | Meter/second Mode of delay | Millisecond |
| Average delay | Millisecond | |
| Total delay | Second | |
| Turnaround time | Second | |

Table 3. Comparison of classifiers.

| Classifier | Accuracy (%) |
|---|---|
| Nave Bayes | 86.13% |
| kNN | 92.82% |
| SVM | 96.75% |

settings in the context of smartphone-based computer games. [77] proposed a reasonable handcrafted features, for three classes (happy, unhappy, neutral); The recent studies in [3] has 7 proposed features, for four classes (Excited, Relaxed, Frustrated, Bored) classification; and [87] compares four discriminative models, namely the Nave Bayes, K-Nearest Neighbor (KNN), Decision Tree and Support Vector Machine (SVM) were explored, with SVM giving the highest accuracy of 96.75%. Table 2 shows a feature set of touch interaction information and Table 3 shows the performance comparison for different classifiers. However readers should be aware that these provide solution doesn't provide any stability analysis of their classification model, it is possible can be counted as overfitting if readers cannot recur the accuracy.

Touch interaction is typical interaction information on mobile which outputs from users. One can conclude here is SVM as the most accurate model for a handcrafted model in this research direction.

## 3.4   Sensors Fusion

Multimodal fusions of previous sensors is necessarily a complex challenge. Much many research has focused on this analysis, in particular, [92] gave bimodal affective recognition via facial emotion recognition and combined audio with visual modalities so that the final affect recognition accuracy is greatly improved to almost 90%. [61] proposed a multimodal system with real-time feedback for public speaking. The system has been developed within the paradigm of positive computing which focuses on designing for user wellbeing.

As multimodal sentiment analysis and emotion recognition research continue to gain popularity among the AI and NLP research communities, there needs for a timely, thorough literature review to define future directions, which can, in particular, further the progress of early-stage researchers interested in this multidisciplinary field (see Figure 2).

The recent surveys of multimodal affect analysis on mobile [16, 71, 70], focuses mainly on state of the art in collecting sample data, and reports performance comparison of selected multimodal and unimodal systems, as opposed to comprehensively reviewing key individual systems and approaches, from the growing literature in the field.

Table 4 from [16] shows the most commonly used experiments result of fusion modalities.

On the same training and test sets, the classification experiment using SVM, NN (Neural network) and ELM (Extreme learning machine). ELM outperformed NN by 12% in terms of accuracy (see Table 5). Regarding training time, the ELM outperformed SVM and NN by a considerable margin (Table 7). A real-time multimodal sentiment

Table 4. Results of feature-level fusion

| Combination of modalities | Precision |
|---|---|
| Accuracy of the experiment carried out on Textual Modality | 0.619 |
| Accuracy of the experiment carried out on Audio Modality | 0.652 |
| Accuracy of the experiment carried out on Vision Modality | 0.681 |
| Experiment using only visual and text-based features | 0.7245 |
| Result obtained using visual and audio-based features | 0.7321 |
| Result obtained using audio and text-based features | 0.7115 |
| Accuracy of feature-level fusion of three modalities | 0.782 |

Table 5. Comparison of classifiers.

| Classifier | Recall (%) | Training time |
|---|---|---|
| SVM | 77.03 | 2.7 min |
| ELM | 77.10 | 25 s |
| NN | 57.81 | 2.9 min |

analysis engine, the NN as a classifier which provided the best performance regarding both accuracy.

In conclusion, we discussed the technique method in a different method. We conclude here in this section for each type of data:

- *Vision data*: CNN models perform the state-of-the-art performance of facial analysis;

- *Speech data*: RNN models performs the state-of-the-art performance of speech recognition and sentiment analysis;

- *Interaction data*: handcrafted features are the most commonly used feature and neural networks for this kind of unstructured data are unminded;

- *Context data*: application context and geographic location are normally used as a validation dimension for emotion inferring, researchers don't consider them how to use as a training feature properly.

## 4 APPLICATIONS CASE STUDY

Assuming that user emotions has inferred by models or systems, [11] address a variety of issues related to the development of affective loop as well as synthesis of the appropriate affective expressions. In subsequent researches emotional systems was sugussted in [13] since emotional-sensitive system make our interactions with machines like human to human communications[66]. expressions. In this section, we illustrate the most popular applications of emotion-sensitive HCI system in the mobile affective research area.

### 4.1 Case 1: Spoken Dialogue Systems

A voice interaction system capable of sending the users emotional messages is able to improving the intelligence and interaction experience of a voice user interface [81].

In the early research stage, [45] proposed a personalized voice-emotion user interface in desktop system regardless of speaker's age, sex or language is presented. They experimented with participants and the results showed that voice emotion sensitive agents are feasible.

The most recent papers in emotion-sensitive voice user interface considers emotional voice tones to caring users [8], as well as the voice control system in an in-vehicle infotainment system [42].

Due to the insivible property of voice interaction system, it is almost endless of how we integrate user emotions to a spoken dialogue system, [54] addressed emotional spoken system by a markup language, which means developers can easily integrate user emotion state to adjust the system voice.

With the raising of voice assistent, there already exists successful commercial voice system such as Apple Siri[3], Google Assistant[4], Amazon Alexa[5] and Microsoft Cortana[6]. Voice user interfaces design become on board of user experience design. Even thgough they provides this kind of markup language, there still a huge unmined research to an open problem of how to evaluate this kind of emotional voice system.

### 4.2 Case 2: Adaptive Graphics User Interfaces

Adaptive user interfaces has been researched for years [76, 47] and addressed in many ways. User emotion is one of the aspect of adaptive user interface.

Dalvand et al.[14] introduced an adaptive user interface, the colors of user interface change according to the emotional states and mood state of users. Emotional states of a user are specified according to his/her interactions with the keyboard. After detection of emotions, the user mood reflects appropriate colors, and [39] as another example of this adaption. However they didn't have a good evaluation of such kind of system. A recent paper [20] also addresses UI adaptation by user emotions (positive, negative and neutral) at run-time. Their prototype was tested successfully of how it reacts to emotions (negative).

In conclusion, the adaptive graphics user interfaces system suguessts covered components with inferring engine (integratable with techniques we discussed in previous section), the adaptation engine (with typical personalized system rules) and the interactive system (for normal graphics user interfaces). It also evidenced that GUI changes denotes emotion changes in run-time, which was particularly beneficial for most users.

## 5 CHALLENGES

Politou et al. [69] discussed most of the challenges in smartphone affective research which covers *privacy*, *data misuse*, *trust and engagement*, *multimodal fusion*, *resource constraints*, *affect modelling and representation*, *cultural differences* and *system building costs*. From the technique prespective in this paper mainly focused, challenges and limitations can be shrink to the following sections.

### 5.1 Impermeable Emotions

The primary limitation of traditional affective computing research refers to as impermeable emotions[65]. Impermeable emotions is broad, with many of these modalities being inaccessible (e.g., blood chemistry, brain activity, neurotransmitters), and many others being too non-differentiated. This makes it unlikely that collecting the necessary data will be possible or feasible in the near.

There is a time to express emotion, and a time to forbear; a time to sense what others are feeling and a time to ignore feelings. In every time, human need a balance when they express their emotions, and this balance is missing in computing. Figure 7 illustrates a map of human emotions.

In most cases, researches feel positive and argue that impermeable emotions can be expressed explicitly in other expressed Emotions[63] and implicity emotions are trivial and not important for most of the case of affective computing application when we need a emotion state [93]. However, impermeable emotions still a open problem and challenges in affective computing research area.

### 5.2 Continious Understanding

Emotions are instantaneous. Continuous emotion state understanding is much more challenge since its not related to external emotion but also related to internal emotions, various emotions can be expressed as a map *(see Figure 7)*.

Most of the researches we introduced in privous section are trated emotion inferring problem as a classification problem, whereas the

---

[3] https://www.apple.com/ios/siri/
[4] https://assistant.google.com/
[5] https://developer.amazon.com/alexa
[6] https://www.microsoft.com/en-us/windows/cortana

Figure 7. A map of human emotions. Image from[79].

human emotions are always passive and instantaneous. People's expression of emotion is so idiosyncratic and variable, that there is little hope of accurately recognizing an individuals emotional state from the available data sometime.

## 5.3 Context Specific Models

A potential ethical limitation in research studies comes from the fact that perceptions of emotions and personality are not universal but they are highly dependent on the mobile application context [22, 77, 3, 87] as well as the cultural differences between humans [55, 51, 24]. According to these perspectives, a reasonable challenge is how an affective recognition model should be designed and built in a application and cultural transparent way.

## 5.4 Lack of Computation Resources

Affective information, like emotions and personality traits, can be inferred from various communication channels such as facial and body expressions, speech, text and embedded smartphone sensors or biosensors. However, a challenging issue that impacts the effectiveness of affective technology is the fusion of these modalities for building competent multimodal affective recognition systems. A multimodal affective recognition system is a system that, via multiple inputs, retrieves affective information from various types of sources and associates input data with a finite set of affective states[21].

The problems in regard to the implementation of an effective system with multimodal fusion are plenty and they mainly concern the great variations of data in terms of structure and content, the varied velocity of data reception, the different sampling rate and quality of received data and the continuous growing of data size. Existing complex technical solutions should be implemented in an energy and power saving approach by keeping in mind that mobile phones have limited operation time and processing power.

Additionally, mobile affective applications should be implemented with a higher portability to address interoperability issues, since mobile devices may not only be equipped with different mobile platforms and operation environments [41], but with different sensor capabilities as well. Consequently, given the diversity of mobile devices in availability and capabilities, modelling and predicting the energy and processing requirements to accomplish even a particular emotion inferring task remains a complex issue.

## 6 CONCLUSIONS

In this paper, we investigated the recent advances in mobile affective computing related to human-computer interaction projects and inferring techniques.

Section 2 addresses different data sources in various mobile commodity sensors for emotion inferring in previous studies. These in-

clude camera, touch screen, motion sensors, microphone, GPS, and application context.

Next, in the Section 3, we first carried out the review of emotion inferring methods based on different type of data source, and compared the tested methods and inferring models from previous researchers. In these comparisons, we first reviewed various models for user emotion inferring, researchers usually transfers emotion inferring problem into a classification problem. As a classification problem, most researchers consider user emotions can be inferred to three different states (Happy, Unhappy, Neutral). In each subsection, we highlighted the most useful methods for different type of emotion inferring that concluded by the most recent research papers, such as CNN as the best way of vision type, RNN as the best way of auditory type, and handcrafted feature with SVM as the best way of interaction details. At the end of this section, we considered the combinations of these types of data. According to our investigation, the most commonly used data type are suggested to *vision and audio* data in mobile affective computing; However, the combination of *vision, audio, interaction details, and context four aspects fusion are unmined open topics*.

In Section 4, we survey two novel applications in human-computer interaction related projects driven by emotion inferring. Voice user interfaces consider user emotion as inputs and suggest to please users by adjusting system voice tone; Adaptive graphics user interfaces then considers emotion state as the context of UI theme.

Despite we researched the scientific approaches of mobile emotion inferring in human-computer interaction related topics, there still apparent challenges in this area. Section 5 pointed out the current problems of this research area. The main challenges of this area are *impermeable emotions* and *continuous understanding*. Moreover, the generalisability of mobile affective computing applications are subject to certain limitations. For instance, most of *inferring models are context specific* and multimodal inferring then *requires large computation resources*.

Nowadays, new technologies and methods provide us new opportunities of affect emotion inferring in an unobtrusive mobile device. Since the complexity of the interpretation of human behavior at a very deep level is tremendous and requires a highly interdisciplinary collaboration, we believe the true break-throughs application in this field can be established by precisely modeling and new sensing technologies in the future.

## REFERENCES

[1] "Affective computing: A review". In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 3784 LNCS (2005), pp. 981–995.

[2] Jeremy N Bailenson, Nick Yee, Scott Brave, Dan Merget, and David Koslow. "Virtual interpersonal touch: expressing and recognizing emotions through haptic devices". In: *Human-Computer Interact.* 22.November (2007), pp. 325–353.

[3] S Bhattacharya. "A predictive linear regression model for affective state detection of mobile touch screen users". In: *International Journal of Mobile Human Computer Interaction* 9.1 (2017), pp. 30–44.

[4] Samit Bhattacharya. "Model for Affective State Detection of Mobile Touch Screen Users". In: 9.1 (2017), pp. 10–13.

[5] Jeffrey P Bigham and Michael S Bernstein. "Human-Computer Interaction and Collective Intelligence". In: *The handbook of collective intelligence* (2014).

[6] Andrew Boles and Paul Rad. "Voice biometrics: Deep learning-based voiceprint authentication system". In: *System of Systems Engineering Conference (SoSE), 2017 12th*. IEEE. 2017, pp. 1–6.

[7] Jessica R. Cauchard, Kevin Y. Zhai, Marco Spadafora, and James A. Landay. "Emotion encoding in human-drone interaction". In: *ACM/IEEE International Conference on Human-Robot Interaction* 2016-April (2016), pp. 263–270.

[8] Woosuk Chang, Angel Camille Lang, Miki Nobumori, Luca Rigazio, Gregory Senay, and Akihiko Sugiura. *Intelligent caring user interface*. US Patent App. 15/085,761. 2016.

[9] Hui Chen, Jiangdong Li, Fengjun Zhang, Yang Li, and Hongan Wang. "3D model-based continuous emotion recognition". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015, pp. 1836–1845.

[10] François Chollet. "Xception: Deep Learning with Depthwise Separable Convolutions". In: *arXiv preprint arXiv:1610.02357* (2016).

[11] Cristina Conati, Stacy Marsella, and Ana Paiva. "Affective interactions". In: *Proc. 10th Int. Conf. Intell. user interfaces - IUI '05* (2005), p. 7.

[12] Roddy Cowie. "Ethical issues in affective computing". In: *The Oxford Handbook of Affective Computing* (2015), p. 334.

[13] Elizabeth a Crane, N Sadat Shami, and Christian Peter. "Let's get emotional: emotion research in human computer interaction". In: *Proceedings of ACM CHI 2007 Conference on Human Factors in Computing Systems* 2 (2007), pp. 2101–2104.

[14] Kasra Dalvand. "An Adaptive User-Interface Based on User's Emotion". In: ().

[15] Charles Darwin and Phillip Prodger. *The expression of the emotions in man and animals*. Oxford University Press, USA, 1998.

[16] Sidney K D'mello and Jacqueline Kory. "A review and meta-analysis of multimodal affect detection systems". In: *ACM Computing Surveys (CSUR)* 47.3 (2015), p. 43.

[17] Mohamad A Eid, Senior Member, and Hussein A L Osman. "Affective Haptics : Current Research and Future Directions". In: 4 (2016).

[18] "Emotion recognition in human-computer interaction". In: *Neural Networks* 18.4 (2005), pp. 389–405.

[19] C Fabian Benitez-Quiroz, Ramprakash Srinivasan, and Aleix M Martinez. "Emotionet: An accurate, real-time algorithm for the automatic annotation of a million facial expressions in the wild". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, pp. 5562–5570.

[20] Julián Andrés Galindo, Sophie Dupuy-chessa, Éric Céret, Université Grenoble Alpes, Bâtiment Imag, Domaine Universitaire, and F Grenoble. "Toward a User Interface Adaptation Approach Driven by User Emotions". In: c (2017), pp. 12–17.

[21] Raghu K Ganti, Fan Ye, and Hui Lei. "Mobile crowdsensing: current state and future challenges". In: *IEEE Communications Magazine* 49.11 (2011).

[22] Yuan Gao, Nadia Bianchi-Berthouze, and Hongying Meng. "What Does Touch Tell Us about Emotions in Touchscreen-Based Gameplay?" In: *ACM Trans. Comput. Interact.* 19.4 (2012), pp. 1–30.

[23] Jose Maria Garcia-Garcia, Victor M R Penichet, and Maria D Lozano. "Emotion Detection: A Technology review". In: *Interacción 2017* (2017).

[24] Maria Gendron, Debi Roberson, Jacoba Marietta van der Vyver, and Lisa Feldman Barrett. "Perceptions of emotion from facial expressions are not culturally universal: evidence from a remote culture." In: *Emotion* 14.2 (2014), p. 251.

[25] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. "Rich feature hierarchies for accurate object detection and semantic segmentation". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014, pp. 580–587.

[26] Ian J Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, et al. "Challenges in representation learning: A report on three machine learning contests". In: *International Conference on Neural Information Processing*. Springer. 2013, pp. 117–124.

[27] Mariam Hassib, Daniel Buschek, PawelW W Wozniak, and Florian Alt. "HeartChat: Heart Rate Augmented Mobile Chat to Support Empathy and Awareness". In: *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (2017), pp. 2239–2251.

[28] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. "Mask r-cnn". In: *arXiv preprint arXiv:1703.06870* (2017).

[29] Kaiming He and Jian Sun. "Convolutional neural networks at constrained time cost". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015, pp. 5353–5360.

[30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conf. Comput. Vis. Pattern Recognit.* (2016), pp. 770–778.

[31] Matthew J Hertenstein, Rachel Holmes, Margaret McCullough, and Dacher Keltner. "The communication of emotion via touch." In: *Emotion* 9.4 (2009), p. 566.

[32] Heinke Hihn, Sascha Meudt, and Friedhelm Schwenker. "Inferring mental overload based on postural behavior and gestures". In: *Proceedings of the 2nd workshop on Emotion Representations and Modelling for Companion Systems - ERM4CT '16* (2016), pp. 1–4.

[33] Sepp Hochreiter and Jürgen Schmidhuber. "Long short-term memory". In: *Neural computation* 9.8 (1997), pp. 1735–1780.

[34] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. "Mobilenets: Efficient convolutional neural networks for mobile vision applications". In: *arXiv preprint arXiv:1704.04861* (2017).

[35] Forrest Iandola, Matt Moskewicz, Sergey Karayev, Ross Girshick, Trevor Darrell, and Kurt Keutzer. "Densenet: Implementing efficient convnet descriptor pyramids". In: *arXiv preprint arXiv:1404.1869* (2014).

[36] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and¡ 0.5 MB model size". In: *arXiv preprint arXiv:1602.07360* (2016).

[37] William James. "What is an emotion?" In: *Mind* 9.34 (1884), pp. 188–205.

[38] William James. *What is an Emotion?* Simon and Schuster, 2013.

[39] Robin Kaiser and Karina Oertel. "Emotions in HCI  An Affective E-Learning System". In: *Computer (Long. Beach. Calif).* (2006), pp. 105–106.

[40] Herman Kamper, Aren Jansen, and Sharon Goldwater. "A segmental framework for fully-unsupervised large-vocabulary speech recognition". In: *Computer Speech & Language* 46 (2017), pp. 154–174.

[41] Wazir Zada Khan, Yang Xiang, Mohammed Y Aalsalem, and Quratulain Arshad. "Mobile phone sensing systems: A survey". In: *IEEE Communications Surveys & Tutorials* 15.1 (2013), pp. 402–427.

[42] Dong-hyu Kim and Heejin Lee. "Effects of user experience on user resistance to change to the voice user interface of an in-vehicle infotainment system: Implications for platform and standards competition". In: *International Journal of Information Management* 36.4 (2016), pp. 653–667.

[43] Hyun-Jun Kim and Young Sang Choi. "Exploring emotional preference for smartphone applications". In: *2012 IEEE Consumer Communications and Networking Conference (CCNC)*. 2012, pp. 245–249. DOI: 10.1109/CCNC.2012.6181095.

[44] Joohee Kim, Na Hyeon Lee, Byung-Chull Bae, and Jun Dong Cho. "A Feedback System for the Prevention of Forward Head Posture in Sedentary Work Environments". In: *Proceedings of the 2016 ACM Conference Companion Publication on Designing Interactive Systems - DIS '16 Companion* (2016), pp. 161–164.

[45] V Kostov and S Fukuda. "Emotion in user interface, voice interaction system". In: *Systems, Man, and Cybernetics, 2000 IEEE International Conference on*. Vol. 2. IEEE. 2000, pp. 798–803.

[46] Alex Krizhevsky and Geoffrey E Hinton. "ImageNet Classification with Deep Convolutional Neural Networks". In: (2012), pp. 1–9.

[47] Pat Langley. "Machine learning for adaptive user interfaces". In: *KI-97: Advances in artificial intelligence*. Springer. 1997, pp. 53–62.

[48] Rodrigo C Lentini, Beatrice Ionascu, Friederike A Eyssel, Scandar Copti, and Mohamad Eid. "Authoring Tactile Gestures : Case Study for Emotion Stimulation". In: (2017).

[49] Rodrigo C Lentini, Beatrice Ionascu, Friederike A Eyssel, Scandar Copti, and Mohamad Eid. "Authoring Tactile Gestures : Case Study for Emotion Stimulation". In: (2017).

[50] Daniel Lopez-Martinez and Rosalind Picard. "Multi-task Neural Networks for Personalized Pain Recognition from Physiological Signals". In: (2017), pp. 3–6.

[51] Takahiko Masuda, Phoebe C Ellsworth, Batja Mesquita, Janxin Leu, Shigehito Tanida, and Ellen Van de Veerdonk. "Placing the face in context: cultural differences in the perception of facial emotion." In: *Journal of personality and social psychology* 94.3 (2008), p. 365.

[52] Antonella Mazzoni and Nick Bryan-kinns. "Mood Glove : A haptic wearable prototype system to enhance mood music in film". In: *Entertainment Computing* 17 (2016), pp. 9–17.

[53] Daniel McDuff, Abdelrahman Mahmoud, Mohammad Mavadati, May Amr, Jay Turcot, and Rana el Kaliouby. "AFFDEX SDK: A Cross-Platform Real-Time Multi-Face Expression Recognition Toolkit". In: *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems - CHI EA '16* (2016), pp. 3723–3726.

[54] Michael F McTear. "The Rise of the Conversational Interface: A New Kid on the Block?" In: *International Workshop on Future and Emerging Trends in Language Technology*. Springer. 2016, pp. 38–49.

[55] Batja Mesquita and Nico H Frijda. "Cultural variations in emotions: a review." In: *Psychological bulletin* 112.2 (1992), p. 179.

[56] Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernockỳ, and Sanjeev Khudanpur. "Recurrent neural network based language model." In: *Interspeech*. Vol. 2. 2010, p. 3.

[57] Ali Mollahosseini, Behzad Hasani, and Mohammad H. Mahoor. *AffectNet: A Database for Facial Expression, Valence, and Arousal Computing in the Wild*. 2017. arXiv: 1708.03985.

[58] Ali Mollahosseini, Behzad Hasani, Michelle J Salvador, Hojjat Abdollahi, David Chan, and Mohammad H Mahoor. "Facial expression recognition from world wild web". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. 2016, pp. 58–65.

[59] "Motion, Emotion, and Form: Exploring Affective Dimensions of Shape". In: *Late-Breaking Work: Designing Interactive Systems* (2016), pp. 1430–1437.

[60] Aske Mottelson and Kasper Hornbæk. "An affect detection technique using mobile commodity sensors in the wild". In: *ACM Int. Jt. Conf. Pervasive Ubiquitous Comput.* (2016), pp. 781–792.

[61] "Multimodal positive computing system for public speaking with real-time feedback". In: *Proceedings of the 18th ACM International Conference on Multimodal Interaction - ICMI 2016* (2016), pp. 541–545.

[62] Ganapreeta R. Naidu. "A Computational Model of Culture-Specific Emotion Detection for Artificial Agents in the Learning Domain". In: *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction - ICMI '15* (2015), pp. 635–639.

[63] Brian Parkinson. *Ideas and realities of emotion*. Psychology Press, 1995.

[64] Phuong Pham and Jingtao Wang. "Understanding Emotional Responses to Mobile Video Advertisements via Physiological Signal Sensing and Facial Expression Analysis". In: *Proceedings of the 22nd International Conference on Intelligent User Interfaces - IUI '17* (2017), pp. 67–78.

[65] Rosalind W Picard. "Affective computing: challenges". In: *International Journal of Human-Computer Studies* 59.1 (2003), pp. 55–64.

[66] Rosalind W Picard. "Affective Computing for HCI". In: *Human Computer Interaction: Ergonomics and User Interfaces. Proceedings of HCI International '99 8th International Conference on Human Computer Interaction* (1999).

[67] Rosalind W Picard. "Affective computing: from laughter to IEEE". In: *IEEE Transactions on Affective Computing* 1.1 (2010), pp. 11–17.

[68] Matthew Pike, Richard Ramchurn, and Max L Wilson. "Two-way Affect Loops in Multimedia Experiences". In: *Proceedings of the 2015 British HCI Conference* (2015), pp. 117–118.

[69] Eugenia Politou, Efthimios Alepis, and Constantinos Patsakis. "A survey on mobile affective computing". In: *Comput. Sci. Rev.* (2017).

[70] Soujanya Poria, Erik Cambria, Rajiv Bajpai, and Amir Hussain. "A review of affective computing : From unimodal analysis to multimodal fusion". In: *Information Fusion* 37 (2017), pp. 98–125.

[71] Soujanya Poria, Erik Cambria, Newton Howard, Guang-Bin Huang, and Amir Hussain. "Fusing audio, visual and textual clues for sentiment analysis from multimodal content". In: *Neurocomputing* 174 (2016), pp. 50–59.

[72] Srinivasan Rajalakshmi, Sarah John Asha, and N. Pazhaniraja. "A comprehensive survey on sentiment analysis". In: *2017 Fourth International Conference on Signal Processing, Communication and Networking (ICSCN)* (2017), pp. 1–5.

[73] Rajib Rana, Margee Hume, John Reilly, Raja Jurdak, and Jeffrey Soar. "Opportunistic and context-aware affect sensing on smartphones: the concept, challenges and opportunities". In: *arXiv preprint arXiv:1502.02796* (2015).

[74] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. "You only look once: Unified, real-time object detection". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, pp. 779–788.

[75] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. "Dynamic Routing Between Capsules". In: *Advances in Neural Information Processing Systems*. 2017, pp. 3857–3867.

[76] Matthias Schneider-Hufschmidt, Uwe Malinowski, and Thomas Kuhme. *Adaptive user interfaces: Principles and practice*. Elsevier Science Inc., 1993.

[77] Sachin Shah, J. Narasimha Teja, and Samit Bhattacharya. "Towards affective touch interaction: predicting mobile user emotion from finger strokes". In: *J. Interact. Sci.* 3.1 (2015), p. 6.

[78] Karen Simonyan and Andrew Zisserman. "VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE-SCALE IMAGE RECOGNITION". In: (2015), pp. 1–14. arXiv: `arXiv:1409.1556v6`.

[79] stanchew. *A map of human emotions*. `https://stanchew.wordpress.com/2012/04/23/a-map-of-human-emotions/`. [Online; accessed 20-November-2017]. 2012.

[80] Thad Starner. "Human-powered wearable computing". In: *IBM systems Journal* 35.3.4 (1996), pp. 618–629.

[81] Kevin J Surace, George M White, Byron B Reeves, Clifford I Nass, Mark D Campbell, Roy D Albert, and James P Giangola. *Voice user interface with personality*. US Patent 6,334,103. 2001.

[82] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alex Alemi. "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning". In: (2016), pp. 4278–4284. arXiv: `1602.07261`.

[83] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. "Going Deeper with Convolutions". In: (2014).

[84] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. "Rethinking the Inception Architecture for Computer Vision". In: (2015). arXiv: `1512.00567`.

[85] Guang Tan, Hongbo Jiang, Shengkai Zhang, Zhimeng Yin, and Anne-Marie Kermarrec. "Connectivity-based and anchor-free localization in large-scale 2d/3d sensor networks". In: *ACM Transactions on Sensor Networks (TOSN)* 10.1 (2013), p. 6.

[86] Jordan Tewell, Jon Bird, and George R. Buchanan. "The Heat is On: A Temperature Display for Conveying Affective Feedback". In: *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems - CHI '17* (2017), pp. 1756–1767.

[87] Subrata Tikadar, Sharath Kazipeta, Chandrakanth Ganji, and Samit Bhattacharya. "A Minimalist Approach for Identifying Affective States for Mobile Interaction Design". In: *Human-Computer Interact. - INTERACT 2017 16th IFIP TC 13 Int. Conf. Mumbai, India, Sept. 25–29, 2017, Proceedings, Part I*. Ed. by Regina Bernhaupt, Girish Dalvi, Anirudha Joshi, Devanuj K. Balkrishan, Jacki O'Neill, and Marco Winckler. Cham: Springer International Publishing, 2017, pp. 3–12.

[88] Sherry Turkle. *The second self: Computers and the human spirit*. Mit Press, 2005.

[89] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. "Attention is All You Need". In: 2017.

[90] Mark Weiser. "The computer for the 21st century". In: *Scientific american* 265.3 (1991), pp. 94–104.

[91] Novita Belinda Wunarso and Yustinus Eko Soelistio. "Towards Indonesian Speech-Emotion Automatic Recognition ( I-SpEAR )". In: 2017 (2017), pp. 8–11.

[92] Zhihong Zeng, Jilin Tu, Ming Liu, Tong Zhang, Nicholas Rizzolo, Zhenqiu Zhang, Thomas S. Huang, Dan Roth, and Stephen Levinson. "Bimodal HCI-related affect recognition". In: *Proceedings of the 6th international conference on Multimodal interfaces - ICMI '04* (2004), p. 137.

[93] Shengkai Zhang and Pan Hui. "A Survey on Mobile Affective Computing". In: 1 (2014). arXiv: `1410.1648`.

[94] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. "Shufflenet: An extremely efficient convolutional neural network for mobile devices". In: *arXiv preprint arXiv:1707.01083* (2017).

[95] Zhan Zhang, Liqing Cui, Xiaoqian Liu, and Tingshao Zhu. "Emotion detection using Kinect 3D facial points". In: *Web Intelligence (WI), 2016 IEEE/WIC/ACM International Conference on*. IEEE. 2016, pp. 407–410.