

Emotional Voice Game Controller

MARJORIE ANN M. CUERDO* and REBECCA LIETZ*, University of California, Santa Cruz, USA

As technology progresses, novel forms of interaction arise. Video games in particular are highly interactive forms of media that are now going beyond traditional handheld controllers to seeing more immersive forms of player input with more embodied technologies throughout extended reality (XR). Can integrating a player's emotions enrich and/or better facilitate that experience? To explore our research question, we developed a method using machine learning techniques to detect emotion in a player's voice as video game input. In other words, we created an emotional voice game controller.

CCS Concepts: • **Computing methodologies** → **Machine learning**; • **Human-centered computing** → *Sound-based input / output*.

Additional Key Words and Phrases: emotions, neural networks, interaction, audio input, games

ACM Reference Format:

Marjorie Ann M. Cuervo and Rebecca Lietz. 2022. Emotional Voice Game Controller. In *CHI '22: ACM CHI Conference on Human Factors in Computing Systems, Apr 30–May 06, 2022, New Orleans, LA*. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Over the last few decades, technological progress has allowed for the development of highly immersive digital experiences for users. One of the most popular and effective forms of immersive technology are video games, which have been established as highly immersive forms of interaction for decades. Audiovisual stimulation in combination with an interactive narrative induces high levels of engagement and identification, which have been shown to aid in learning, emotion regulation, and even pain management (citation needed). State-of-the-art systems using extended reality (XR) approaches, such as virtual reality (VR) and augmented reality (AR), can put the user in a flow state by engaging multiple senses and providing alternative input methods that open up new possibilities for user engagement. For example, game systems may track user engagement and change the game content when the player gets bored or frustrated (citation). In recent years, a growing number of studies have explored affective gaming, or the consideration of player emotion when interacting with a game. However, most of those studies present games that are either about emotions, or that react to emotion but do not necessarily prompt the player to reflect on the emotions they are displaying. Our team has taken the concept of affective gaming and developed a game that requires the player to express certain emotions using their voice (thereby motivating reflection) and responds accordingly.

2 RELATED WORK

2.1 Role of emotions

Emotions have a large impact on various aspects of our lives. Research has linked emotional states to changes in memory (citation), attention (citation), and a range of behaviors, such as eating (citation) and driving (citation). Based

*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

on these findings, it makes sense to view emotions as highly influential in how we experience life. A long history of emotion research has led to the development of numerous emotion models, which generally fall into one of two categories: dimensional and discrete. Dimensional models, such as the circumplex model (citation), can describe a wide range of emotions along multiple axes and take into account small variations in a person's mood. While this framework gives a more realistic view of how people actually experience emotion, it can also be helpful to sort emotional states into discrete categories. The most basic model includes six core emotions: anger, sadness, disgust, happiness, surprise, and fear (citation). Discrete emotion models are particularly helpful when viewing emotion detection as a classification problem and employing machine learning techniques to solve it.

2.2 Emotion detection

Since emotions are internal states and therefore inherently difficult to measure and detect, researchers have experimented with numerous ways to identify emotions using external cues, such as facial expressions, speech, physiological data, and behaviors (citation). Behaviors such as fidgeting, mobile phone usage, or interaction on social media, can be indicative of certain emotional states. But although this type of data is relatively easy to obtain, it is often not clear whether the measured behaviors are valid indicators of the emotion they are supposed to represent. Conversely, facial expressions very obviously show a person's emotional state. In addition, the only equipment needed to collect facial data is a camera, which is already built into many personal computers. However, facial expressions are easy to fake or hide, which can make it difficult to accurately assess emotions. This method also heavily relies on the availability of suitable datasets in order to train a model reliably. Physiological responses are least likely to be intentionally modified by a person, and are therefore highly reliable measures for arousal. On the downside, measuring bodily responses usually requires additional sensors and does not account for context. Lastly, people also express emotion through speech, using different intonations, pitch, and frequency. Similar to facial expressions, speech data can be easily collected using microphones built into a user's personal device. The main disadvantage of this approach is that it, like the facial expression method, relies on the existence of suitable datasets to train a machine learning model. Additionally, recordings from built-in microphones might be noisy and hard for the system to process. Yet, systems that use speech input for emotion recognition enjoy popularity within the scientific and commercial HCI community (Garcia-Garcia citation). Khalil et al. (citation) outline the need for speech emotion recognition to employ more non-linear deep learning techniques as opposed to the more traditional techniques utilizing KNN, HMM, and SVM classifiers, which we attempt to realize with our system.

2.3 Emotion in video games

As highly interactive media, games have the ability to evoke a diverse range of specific emotional experiences for players. Game designers manipulate Mechanics-Dynamics-Aesthetics (MDA) [1] to elicit particular types of player experiences according to their vision for their game. Regardless, many methods of game performance evaluation simply rely on achieving a general state of flow [3]. Games are useful to study emotions, as they can be broken down into components. Players' emotions can be influenced through anything from the game's narrative, aesthetic representation (e.g. colors, imagery, sound, etc.), presented challenges, to even sociocultural context.

In addition to inducing flow, people often express that games should be fun. Lazzaro claimed that there are different types of "fun" that should be differentiated: Easy Fun, Serious Fun, Hard Fun, and People Fun [4]. These types of fun correspond to specific types of human emotions. Easy Fun is more relaxed and focused on feelings of curiosity, wonder, and awe. Serious Fun is a bit more repetitively involved and evokes excitement and zen focus. Hard Fun involves intense

fiero – the feeling of triumphing over hardship or experiencing relief after frustration. People Fun occurs with other people, experiencing amusement and admiration.

While categorization can be useful for design, it is becoming more apparent over time that emotions in games are more complex than initially assumed. Bopp et al. [5] found that when reflecting on emotionally-moving game experiences, players most enjoyed and appreciated negatively valenced emotions, counter to what one may initially expect. It can no longer be claimed that a positive player experience is solely depending on experiencing purely positive emotions. Serious games for learning are found to be effective when involving emotionally challenging experiences [], so the need to further understand emotions in games goes beyond entertainment purposes.

2.4 Affect-adaptive gaming

Sundstrom [2] describes an affective loop as "an interaction process where:

- (1) the user first expresses her emotions through some physical interaction involving the body, for example, through gestures or manipulations of an artifact
- (2) the system (or another user through the system) then responds through generating affective expression, using for example, colors, animations, and haptics
- (3) this in turn affects the user (both mind and body) making the user respond and step-by-step feel more and more involved with the system"

Keeping those requirements in mind, it is easy to see how games are a natural way for an affective loop to be realized. When people are playing a game, they are required to constantly provide input and receive output (visual, auditory, and/or haptic feedback) from the system. Input to this loop can vary from behavioral data (game metrics) to objective data (bodily responses).

3 METHODS

To accomplish creating our emotional voice game controller, our process involved the following steps:

- (1) Pre-processing the audio data
- (2) Training the neural network
- (3) Creating the API
- (4) Developing the game and connecting to API

3.1 Pre-processing the audio data

We used the Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS). We used the speech audio files from 24 actors either saying "Kids are talking by the door" or "Dogs are sitting by the door" with any of the following eight emotions: neutral, calm, happy, sad, angry, fearful, disgust, surprised.

Due to lack of appropriate hardware, we turned to Google Colaboratory to access GPU power. To be used in a neural network, the voice audio files were converted into spectrogram images using the librosa library. After running the code in this [Colab notebook](#), the spectrogram images were sorted into the respectively labeled folders. We then downloaded the sorted data folder locally and uploaded it to a Google Drive folder for Colab to access later.

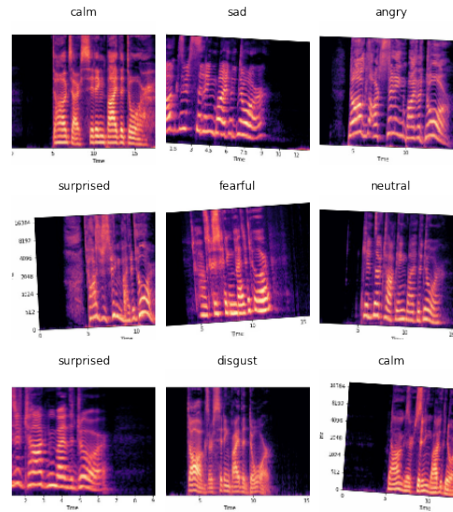


Fig. 1. Examples of categorized spectrograms according to emotions.

3.2 Training the neural network

We used the following [Colab notebook](#) to perform transfer learning/train the neural network on the sorted data folder we obtained above. This method required the use of the *fastai* library and ResNet-34, a pretrained Convolutional Neural Network (CNN) model trained on the ImageNet dataset used for image recognition. Our goal was to achieve training a model that moderately detects the correct emotion in any voice recording regardless of the content of speech.

We used the *fit_one_cycle()* function to train our model. We used five epochs every cycle and modified the learning rate through reading and interpreting the *lr_find()* graph (example depicted in Figure 2). We selected the sections with the steepest decline in accuracy every cycle and observed how that affected overall accuracy of the model.

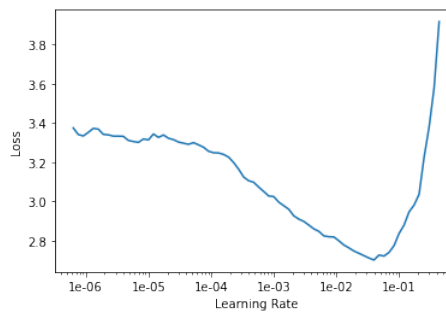


Fig. 2. An example of learning rate accuracy graph.

3.3 Creating the API

... ..

3.4 Developing the game and connecting to API

The largest hurdle in our project was to figure out the more technical aspects in regards to machine learning and the backend flow of our model to a game. Therefore, we decided that our game should be simple but succinct in demonstrating the capability of emotion detection in voice. We ended up using Unity to create a simple 2D side-scrolling game where the player controls a character who goes on an emotional and difficult journey towards realizing they need therapy (a screenshot displayed in Figure 3).

We were inspired to pursue this theme because of the various stages involved in the journey to acceptance, such as denial, anger, depression, bargaining, and acceptance. Each one of those stages could be narratively related to one of the emotions present in our emotion detection model (neutral, calm, happy, sad, angry, fearful, disgust, and surprised). To invoke the player to emotionally react in a more natural way, we embedded emotion detection into the game's mechanics. To "win" the game, the player must get past their mental blocks (represented by actual physical blocks in the game) by emotionally speaking in a way that's appropriate for the current situation presented in the game. For example, if the character just got unexpectedly fired, the player would have to attempt to sound surprised ("How could you fire me?!"). If the player was able to vocally emote in the way the game wants, the player progresses.

Regarding the backend of this emotional voice game controller, the player has to tap "Record", speak into their device, then tap "Send". A .wav file is then created which is sent to our Emotion Prediction API. Our model then sends the detected emotion in that .wav file back to the game in Unity in text form, which is displayed to the player and unlocks new areas of the game (if correct). The player can attempt this as many times as needed until they reach the end of the game.

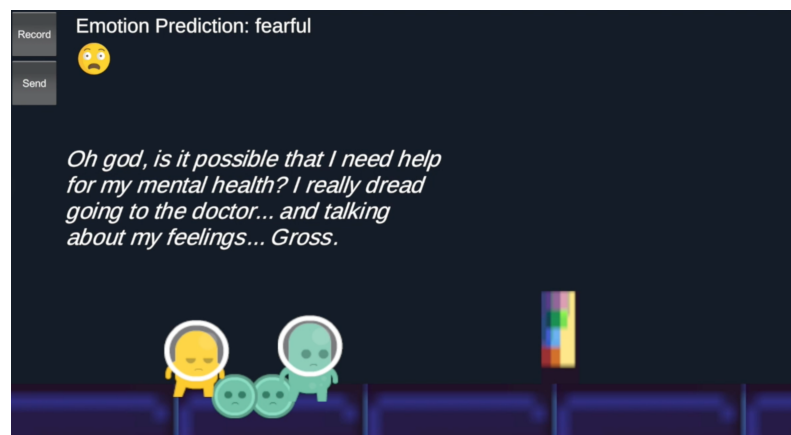


Fig. 3. A screenshot of the journey to therapy game we created.

3.5 Resources

All of our code (Colab notebooks and Unity game) can be found in this [Github repository](#).

4 RESULTS

4.1 Model accuracy

After training our model, we were able to achieve about 80% accuracy for emotion detection in voice. Though this is far from perfect, this was as best as we could get it improve with our current knowledge. Hopefully, we can find methods to improve this in the future.

4.2 Game and emotion prediction performance

We were pleasantly surprised at how seamless recording the player's voice and obtaining the emotion prediction from our model was after it was properly set up. The prediction wasn't instantaneous but only took 1-2 seconds each time, despite this process involving creating a .wav file, converting that into a spectrogram, using the model to predict emotion from that image, and then sending that to Unity to affect the gameplay. Despite this, we are still hoping to improve the backend connection process.

5 CONCLUSION

REFERENCES

- [1] Robin Hunicke, Marc LeBlanc, and Robert Zubek. "MDA: A formal approach to game design and game research". In: *Proceedings of the AAAI Workshop on Challenges in Game AI*. Vol. 4. 1. San Jose, CA. 2004, p. 1722.
- [2] Petra Sundström. "Exploring the affective loop". PhD thesis. 2005.
- [3] Jenova Chen. "Flow in games (and everything else)". In: *Communications of the ACM* 50.4 (2007), pp. 31–34.
- [4] Nicole Lazzaro. "Why we play: affect and the fun of games". In: *Human-computer interaction: Designing for diverse users and domains* 155 (2009), pp. 679–700.
- [5] Julia Ayumi Bopp, Elisa D Mekler, and Klaus Opwis. "Negative emotion, positive experience? Emotionally moving moments in digital games". In: *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. 2016, pp. 2996–3006.