# Towards Inclusive HRI: Using Sim2Real to Address Underrepresentation in Emotion Expression Recognition

Saba Akhyani, Mehryar Abbasi, Mo Chen, and Angelica Lim

*Abstract*— Robots and artificial agents that interact with humans should be able to do so without bias and inequity, but facial perception systems have notoriously been found to work more poorly for certain groups of people than others. In our work, we aim to build a system that can perceive humans in a more transparent and inclusive manner. Specifically, we focus on dynamic expressions on the human face, which are difficult to collect for a broad set of people due to privacy concerns and the fact that faces are inherently identifiable. Furthermore, datasets collected from the Internet are not necessarily representative of the general population. We address this problem by offering a Sim2Real approach in which we use a suite of 3D simulated human models that enables us to create an auditable synthetic dataset covering 1) underrepresented facial expressions, outside of the six basic emotions, such as confusion; 2) ethnic or gender minority groups; and 3) a wide range of viewing angles that a robot may encounter a human in the real world. By augmenting a small dynamic emotional expression dataset containing 123 samples with a synthetic dataset containing 4536 samples, we achieved an improvement in accuracy of 15% on our own dataset and 11% on an external benchmark dataset, compared to the performance of the same model architecture without synthetic training data. We also show that this additional step improves accuracy specifically for racial minorities when the architecture's feature extraction weights are trained from scratch.

## I. INTRODUCTION

There has been an increasing interest in using robots in everyday social environments like hospitals, retail stores, and homes. As a result, it has become essential for robots to communicate and interact socially for various human-robot interaction (HRI) applications and understand people's nonverbal expressions. For instance, imagine a restaurant service robot seeing a look of confusion wash over your face as you look at the menu. It detects your hesitation and proactively offers assistance: "Do you have any questions I can answer?" This simple scenario plays out between humans each day all over the world, but robots are still far from capable of performing this kind of proactive assistance in a robust manner. Several major challenges prevent such systems from being deployed in the real world.

Firstly, a recent review has argued that state-of-the-art facial emotion classifiers cannot be applied effectively to human emotion analysis in the wild [3]. One underlying reason is that in HRI, as stated in a review by [22], "there are a wide range of possible affective levels expressed by

people in human-robot interaction that the robot needs to understand in order to participate in a natural bi-directional social interaction with humans." In other words, real-world interactions comprise a rich and subtle set of expressions, while most datasets focus on collecting the prototypical set of emotions [3] of happiness, sadness, anger, surprise, fear, disgust, and neutral [17], [34], [40], [19], Specifically, to the best of our knowledge, there is no public video dataset containing confusion [38] nor any comparison benchmark [15] for these dynamic facial expressions [16], [8]. Strategies are needed to effectively create video datasets for the many underrepresented emotional expression categories that *actually* occur in the wild with robots, such as the 28 social signals identified in Saheb Jam et al. [30], including confused, worried, skeptical, and so on. Indeed, prototypical sadness or fear were not seen in these real-world human-robot interactions, and we do not focus on these emotions in this paper.

Secondly, social robots should also have the ability to evaluate human affective expressions fairly, without discriminating against underrepresented groups. A recent survey on automatic multi-modal emotion recognition in the wild shows that inclusivity of all ethnicities remains a challenge in facial expression recognition systems and should be further investigated [32]. According to [27], [5], racial bias is apparent in current machine learning methods in general, especially those involving the face. One major cause of this bias is that major facial expression recognition (FER) datasets are underrepresentative of genders and non-Caucasian backgrounds [36]. As faces are inherently identifiable, ethical and privacy concerns arise when requiring real humans to provide their data [18], yet anonymizing the face can remove important facial features. Data collection to reduce bias in face-related algorithms is, therefore, a major ethical challenge.

Finally, mobile robots can potentially view humans from various viewing angles and in varied lighting conditions. Collecting and labeling large amounts of naturalistic videos for facial expression recognition is challenging due to several reasons. First, creating datasets for spontaneous (rather than posed) user affective states is very time-consuming [33]. Additionally, facial emotional expressions are difficult to label due to the subjective nature of annotation, compared to those in domains where deep learning has been most successful, such as in object recognition. Thus, data augmentation techniques are expected to be particularly useful to help facial expression recognition succeed in the wild.

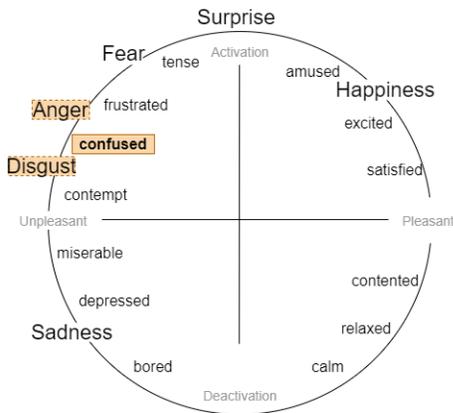In this paper, we employ a simulation to reality (Sim2Real)

Fig. 1: Circumplex Model of Affect adapted from [29].

approach to address the previously mentioned challenges. Sim2Real approaches have performed well in different domains such as hand tracking [23] and text detection [13]. For face recognition and facial feature detection, studies have attempted to address some of these problems by creating static, synthetic image datasets [11], [1], [14], [20]. We extend these approaches into the video emotion recognition domain by i) identifying facial movements and social signals of the desired dynamic emotional expressions from real data, ii) converting the identified social signals into renderable animations, iii) generating virtual human models of various and specific ethnicities, iv) rendering the animations into videos of all virtual human models from multiple viewing angles, v) pretraining a dynamic emotion classifier model on the synthetic videos, vi) retraining and testing the model on the real dataset.

We generate an inclusive synthetic facial expression dataset from virtual humans that specifically incorporates new expressions outside of the prototypical set of six basic emotions, underrepresented ethnicities, and varied camera angles. As a proof-of-concept, we detect the dynamic social signal of confusion. In order to evaluate the model, anger and disgust were chosen from the 6 basic emotions as challenging emotion expression confounders due to their similarity to confusion. Characteristics of these three emotions overlap, and expressions can be easily mistaken for one another, as illustrated in the Circumplex Model of Affect (CMA) shown in Fig. 1. The circumplex model is a graphical representation of affective states on the 2-D plane of arousal (vertical) and valence (horizontal).

Our contributions are as follows:

1) We propose the first video dataset of the understudied emotional expression of confusion (as well as nearby social signals "disgust" and "anger" as shown in the circumplex model in Fig. 1). We first gather a real dataset, then augment this limited data using Sim2Real to produce a much larger synthetic dataset. Our dataset also addresses the critical issue of racial bias, which is apparent in existing real-world data, by comprising faces of underrepresented ethnicities, including Black,

Asian and Hispanic individuals.

2) We explore the effect of adding synthetic data on improving fairness using a CNN+Time Series Classification (TSC) network architecture. Our experiments demonstrate that: 1) training on a combination of real-world data and a randomly selected portion of synthetic data (changing every epoch) achieves the highest performance and 2) fine-tuning on a pre-trained CNN with *unfrozen face feature extraction weights* decreases racial bias.

The Sim2Real approach would allow us to create even larger synthetic datasets in future, because of the flexibility to be applied to any desired emotion or ethnicity. This is feasible due to the fact that facial movements (action units) associated with any emotion can be extracted either automatically using OpenFace [2] or by manual observation. Additionally, the wide modification range of the MakeHuman toolkit [4] allows for the generation of several human models to incorporate other ethnicities.

## II. METHODOLOGY

The methodology behind our Sim2Real approach is explained in this section. An overview of our dataset generation and preparation, as well as an overview of our deep-learning-based dynamic facial expression recognition model, is provided. The synthetic data pretraining step of the Sim2real approach is explored in Section III.

### A. Dataset Generation

In this section, we describe the collection of in-the-wild emotionally expressive videos and the generation of synthetic videos using a suite of simulated humans. We focus this study on confusion, a dynamic social signal which is underrepresented in datasets and lacks examples on the web [15], yet is common in HRI [30].

*1) Collection of in-the-wild confusion, anger, and disgust videos:* Confusion is an affective state conveyed through varied and multiple expressions, as are both anger and disgust. Some of these expressions (such as frowning) are common between the three, resulting in some expressions being easily mistaken as another emotion, possibly due to these three emotions appearing very close together in the CMA [29]. We therefore focus on making a dataset for these three emotions.

We collected short video clips (1-3s) from YouTube.com and Giphy.com using search tags such as "angry", "confused", and "disgust" reactions to gather human facial expressions of our desired social signals, as shown in Fig. 2. This search resulted in 153 clips. Each video was then labeled for the conveying facial expression, by two annotators identifying with Canadian culture (inter-rater agreement kappa score=.88), and low confidence videos were discarded.

We created a multi-ethnicity dataset of real human videos expressing the three social signals of confusion (41 videos), anger (41 videos), and disgust (41 videos). The final dataset contains 123 videos, of which 26 are of non-Caucasian individuals.

Fig. 2: Examples of real videos: the first row is for angry, second for disgusted, and third row for confused expressions



Fig. 3: Sample of generated human-like models



(a) AU17    (b) AU23    (c) AU26    (d) AU61

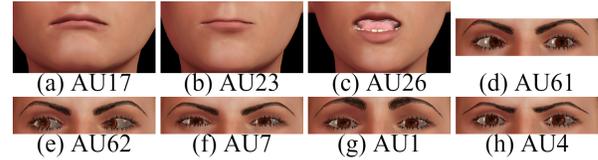(e) AU62    (f) AU7    (g) AU1    (h) AU4

Fig. 4: (a) Chin raiser, (b) Lip tightener , (c) Jaw drop, (d) Eyes left, (e) Eyes right, (f) Lid tightener, (g) Inner brow raiser, (h) Brow lowerer

*2) Generation of augmented dynamic social signal video dataset:* The task of creating desired social signals videos is made possible using the MakeHuman toolkit [4] and the FACSHuman plugin [12]. MakeHuman toolkit is an open-source and free 3D computer graphics toolset designed for prototyping human-like models. FACSHuman offers the possibility of manipulating the Action Units (AU) presented in the Facial Action Coding System (FACS) [9] on the 3D models created in the MakeHuman software. This manipulation of AUs is a key component of our Sim2Real process. FACSHuman enabled us to generate social signal animations that can be rendered into videos or frames on a selected human virtual model from any viewing angle. We created a script plugin that used this capability to render 4536 synthetic videos from the combination of 24 human virtual models, 21 social signals (facial movement animations), and 9 viewing angles.

**Creation of a suite of simulated humans**: The overarching vision of this work is to create a large, auditable suite of human models to represent people from many different backgrounds. As a first step, we create 24 simulated human adult models balanced on gender and four different ethnicities (Caucasian, Black, Asian, and Hispanic). The MassProduce plugin within the MakeHuman application was then used to create several randomly generated human models of multiple ethnicities and different ages and skin colors. Out of all those generated human-like models, we selected 24 models for our study based on the realism of action unit manipulation on the model (8 samples are shown in Fig. 3). We chose to use 3D models as we hope to eventually use them in HRI simulators, to create expressive virtual humans with facial expressions. These models are provided on Github[1] so that researchers can also import the 24 virtual humans with 21 dynamic expressions (7 social signals per emotion), to replay them in front of their virtual robot.

**Multiple social signals per emotion**: We used the FACSHuman software to create 21 different social signal animations, 7 for each emotional class. These animations convey multiple variations of social cues of confusion, anger, and disgust. These 21 social signals were manually animated over 25 frames and were created via inspection of the in-the-wild real human dataset. For example, Fig. 4 shows the AUs that were used to create dynamic confusion social signals. Varied AU combinations and sequences were used to animate the 21 social signals, validated by an annotator with Canadian culture. An example is a side-eye movement confusion state made by a timed sequence of the following AUs: AU61, AU62, AU61. While future work should automatically perform the animation creation process from video data, the manual animation creation step in this study allows us to validate the Sim2Real portion given human-level feature extraction, enabling us to identify specific underlying social signals for each emotion (7 for each emotion), which are now available for use in our 3D models. The dataset is available for download.[2]

**Multiple viewing angles**: As robots may view a human from varied angles, it is important that our generated dataset incorporate varied perspectives. Our dataset was therefore expanded by creating videos of the same facial gesture from 9 viewing angles, to make our network invariant to the face viewing angle, as shown in Fig. 5. The camera movement included horizontal rotations of $-40, -20, 0, 20, 40$ degrees and vertical rotations of $-30, 15, 0, 15, 30$ degrees. The

---

[1] https://github.com/sabaak95/confusionDetection

[2] https://www.rosielab.ca/datasets/confusion-in-the-wild

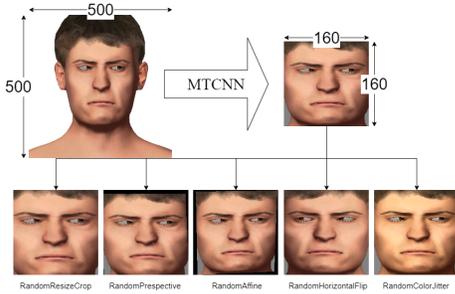Fig. 5: The 9 viewing angles used to generate our augmented dataset.



Fig. 6: Image preprocessing and augmentation.

nine combinations of $(H_{rotation}, V_{rotation})$={$(-40, -30)$, $(-20, -15)$, $(0, 0)$, $(20, 15)$, $(40, 30)$, $(40, -30)$, $(20, -15)$, $(20, -15)$, $(40, -30)$} were selected as our viewing angles in degrees, as shown in Fig. 5.

### B. Data preparation

In order to refine the data and remove any unimportant or unrelated information in the images, we used Multi-task CNN (MTCNN) to detect and crop the faces before feeding frames to our network [41] and resized images to $160 \times 160$. Additional transforms were also applied to the images randomly on each epoch, including cropping, perspective, affine transform, horizontal flip, and color transforms (shown in Fig. 6).We ensured that the same transformations were applied to all the frames from the same video.

### C. Model Architecture

We developed a basic framework for our video classification problem to test the Sim2Real strategy. Existing work can extract valuable frame-based facial features from a face image, such as FaceNet [31] and OpenFace [2]; one of these models can be used to first extract each frames' facial features. After frame-based feature extraction, we model the problem as Time Series Classification. Sections II-C.1, II-C.2 are dedicated for further exposition on our selections for this architecture.
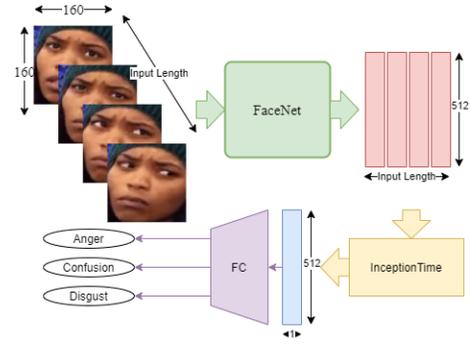


Fig. 7: Model architecture for FaceNet+InceptionTime

*1) Facial Feature Extraction Network:* We used the pre-trained FaceNet [31] architecture as our facial feature extractor. FaceNet uses an InceptionResnetV1 architecture trained on the VGGFace2 dataset. Each video is given to FaceNet frame by frame, and the output feature arrays are concatenated together across the time dimension to create a multivariate time series array.

*2) Time Series Classifier Network:* Further processing of this output requires a time series classification algorithm. K-Nearest Neighbor (KNN) algorithm with Dynamic Time Warping (DTW) [39] metric is one of the earliest techniques for this task that is still used, specially when working with relatively small datasets. Many machine learning algorithms have also been applied to this problem, such as ResNet and FCN [37]. However, we opted to use InceptionTime [10] as our classifier because it has proven to be a versatile and promising machine learning solution for many Time Series Classification tasks, based on its performance results on the UCR [7] benchmark collection datasets.

*3) Proposed and Baseline Architectures:* Our proposed DNN structure is the combination of FaceNet [31] and InceptionTime128 [10]. This structure is shown in Fig. 7. This model is referred to here as FN+INC25 or FN+INC64. The two numbers of 25 and 64 indicate the required number of frames for the video input. Videos shorter than the indicated number are padded to the required input length, and the longer videos are cropped. This procedure is explained in Section III. We included two instances of these models in our experiments. In one instance, the FaceNet weights are frozen. In the main instance, the FaceNet weights are not frozen, and are tuned alongside the InceptionTime weights. We hypothesized that the addition of synthetic data would allow us to tailor the FaceNet weights in favor of this dataset.

The second DNN architecture included in our experiments is I3D [6], an advanced video classification method that applies combined temporal and spatial processing using 3D convolutional layers. This addition enabled us to evaluate the problem using a model not already affected by previous facial information knowledge.

Finally, we included a KNN classifier architecture applicable to the small real dataset. This model uses FaceNet [31] for frame facial feature extraction and a KNN with DTW [39]

metric as the video classifier. We propose this baseline model, FN+KNN, for comparison.

## III. EXPERIMENTS AND RESULTS

In this section, we evaluate our network on the created real dataset. We performed several experiments varying the architecture, input length, and the use of only synthetic, synthetic plus real, or only real training data. We used 5-fold cross-validation on the *real dataset* to compare the performance of different approaches, with one fold consisting primarily of expressions by non-Caucasian individuals.

We explored three training strategies for our experiments. In the first strategy, the algorithms are only trained on the small real dataset. The baseline KNN model was tested under this strategy. In the second strategy, the networks are first trained on synthetic data, then fine-tuned on the real dataset. The second strategy was developed to add and assess the addition of synthetic data. Third, the strategy was to combine the real training data set with the synthetic dataset and pass them to the network alongside each other. This strategy was designed to explore if a higher performance could be achieved by creating ratioed synthetic and real data training. Instead of combining the whole synthetic dataset with the real data, we trained the network with one-fourth of the synthetic dataset. In the next test, the ratio of the synthetic dataset was set to half. Finally, in the last test, the ratio of the synthetic dataset was set to one, meaning the whole synthetic dataset was included.

One important factor in our training and testing processes is setting the input video length to a fixed number of frames $L$. Input videos shorter than the set length were looped until they reached length $L$. The way we dealt with longer videos differed depending on the training phase. In the training phase, we randomly selected $L$ consecutive frames from the lengthier videos. For a video with $N$ frames, frames $n$ to $n + L - 1$ are cropped. The $n$ is selected randomly on each epoch between 0 and $N - L$. However, in the testing phase, we only selected the middle $L$ frames as the representative sequence in each video. In our experiment, we set $L$ to two values: 25 and 64. We selected 25 because the number of frames in our synthetic videos was 25. The choice of 64 was reliant on two factors. First, 64 was long enough to include a majority of the input video while small enough to keep computational cost and time consumption adequate. Second, the I3D network used in some of our experiments was designed for 64 frame inputs. In the following subsections, we elaborate on the latter two training strategies: (i) fine-tuning the synthetic trained network on real data and combined synthetic. (ii) Combined synthetic and real data training.

### A. Fine-tuning the synthetic trained network on real data

In this experiment, the model was first trained on the synthetic dataset alone. The simulated human models were randomly divided into two sets of 19 and 5 models. All of the generated videos using simulated human models in the larger set were used for training, and those in the smaller set were used for validation. The respective numbers for the

TABLE I: Performance comparison of all models

| Network | Length | Syn[1] | %Prc[2] | %Rec[3] | %Fs[4] | %Acc[5] |
|---|---|---|---|---|---|---|
| FN(FZ*)+KNN | 25 | ✗ | 75 | 71 | 70 | 72 |
| FN(FZ*)+KNN | 64 | ✗ | 79 | 74 | 73 | 74 |
| FN(FZ*)+KNN | 283 | ✗ | 77 | 71 | 69 | 71 |
| **FN(FZ*)+IN25** | 25 | ✗ | 76 | 73 | 73 | 73 |
| **FN(FZ*)+IN25** | 25 | ✓ | 82 | 79 | 79 | 80 |
| **FN(FZ*)+IN64** | 64 | ✗ | 86 | 71 | 71 | 72 |
| **FN(FZ*)+IN64** | 64 | ✓ | 85 | 82 | 82 | 82 |
| **FN+IN25** | 25 | ✗ | 78 | 74 | 74 | 74 |
| **FN+IN25** | 25 | ✓ | **90** | **89** | **89** | **89** |
| **FN+IN64** | 64 | ✗ | 81 | 77 | 76 | 77 |
| **FN+IN64** | 64 | ✓ | 88 | 87 | 87 | 87 |
| I3D | 64 | ✗ | 66 | 63 | 60 | 66 |
| I3D | 64 | ✓ | 85 | 83 | 83 | 83 |

[*] Frozen FaceNet weights    [1] Synthetic-Data
[2] Precision    [3] Recall    [4] F-score    [5] Accuracy

TABLE II: Effect of the addition of synthetic data and weight unfreezing on the non-Caucasian Fold

| NET | Len | Synth | | %Acc | %Acc Increase |
|---|---|---|---|---|---|
| FN(FZ)+IN25 | 25 | ✗ | ‖ | 72 | Base |
| FN(FZ)+IN25 | 25 | ✓ | ‖ | 80 | +8 |
| FN+IN25 | 25 | ✗ | ‖ | 80 | +8 |
| FN+IN25 | 25 | ✓ | ‖ | 88 | +16 |
| FN(FZ)+IN64 | 64 | ✗ | ‖ | 64 | Base |
| FN(FZ)+IN64 | 64 | ✓ | ‖ | 72 | +8 |
| FN+IN64 | 64 | ✗ | ‖ | 80 | +16 |
| FN+IN64 | 64 | ✓ | ‖ | 88 | +24 |

videos in the training and validation data were 3591 and 945. The training was done over 20 epochs with the learning rate of $10^{-4}$ and the categorical cross-entropy loss function. The batch size was set to 8. After the training on synthetic data, we fine-tuned the model on the four selected training folds of the real dataset, over 50 epochs. The model is tested on the remaining single test fold. This operation is repeated 5 times each time a new fold is selected as the test fold. Learning rates and parameters were chosen empirically.

The results averaged over all 5 runs for these experiments are shown in Table I. FZ specifies the instances where FaceNet weights were frozen to treat FaceNet purely as a feature extraction network, with weights of the rest of the network updated during training. Table I also includes experiments in which the synthetic data training step was skipped to highlight its effect. Additionally, we compared our methods with the baseline FN+KNN classifier applied only to the real data. Our results show that the models trained on synthetic data outperformed their counterparts only trained on the real data. The unfrozen FN+IN25 model achieved an

89% accuracy on the real data when trained on synthetic data. In the frozen weights instances, the inception models perform similarly to the KNN models when trained only on the real data. However, the addition of synthetic data training improved the FZ model accuracy up to 83% in the case of FN(FZ)+INC64. This addition also significantly impacted the I3D model, and its accuracy of 83% outperforms all models not influenced by the synthetic data. Interestingly, this model even outperforms the FZ models trained on the synthetic data. This is quite impressive because I3D was designed for video action recognition tasks. Unlike the other models, the I3D had no prior information about the facial features.

In Table II we explored the effect of unfreezing the FaceNet weights and the addition of synthetic data on the non-Caucasian data fold. We created one test fold which included 25 videos of the underrepresented ethnicities. The FN(FZ)+IN models trained without synthetic data are highlighted as the base models in this table. This table shows that the addition of synthetic data combined with unfreezing of the pre-trained weights has the highest impact on the correct classification of the non-Caucasian data samples (24% increase). The addition of synthetic data alone has a limited beneficial impact; it can not alter the dataset bias effect of the original dataset on which the FaceNet was trained.

### B. Combining synthetic and real data for training

We designed another experiment to investigate how the accuracy changes with the addition of synthetic data. A portion of the synthetic data was randomly selected and combined with the real training data to create a new data set. The model was trained on this new training data and tested only on the real test data sample. We applied this training strategy to our most satisfactory model, input length 25 FaceNet + InceptionTime. In our first experiment, we set the ratio of selected synthetic data portion to 0.25. This ratio was doubled in the next experiment and doubled again in the last one. In each epoch $\lfloor ratio \times 24 \rfloor$ human-like models were randomly selected. For every selected human-like model, out of the nine videos of that human-like model expressing a specific expression from multiple angles, only one was chosen randomly. This selection method means that only $\lfloor ratio \times 24 \rfloor \times 21$ synthetic videos are used in the models training alongside the real data in that epoch. This number equals 126 for the ratio of 0.25, roughly equal to the number of real training videos. The selected human-like models and viewing angles were refreshed at the start of each epoch.

The results for these experiments are shown in Table III. These results show that doubling the synthetic data ratio from 0.5 to 1 increases the model's performance. However, this does not apply to the change from 0.25 to 0.5. In the case of FaceNet+Inception64, the synthetic ratio of 0.25 results in the highest performing network. This model achieved a 94% accuracy, which shows an 18% increase over the performance of the same model trained without the synthetic data. The combined confusion matrix of all the folds for this highest performing network is shown in Fig. 8.

TABLE III: Performance comparison for the combined real and synthetic training method for the FN+INC models.

| Ratio[*] | Length‖ | %Precision | %Recall | %F-score | %Accuracy |
|---|---|---|---|---|---|
| 0.25 | 25 ‖ | 89 | 87 | 87 | 87 |
| 0.25 | 64 ‖ | **95** | **94** | **94** | **94** |
| 0.5 | 25 ‖ | 88 | 86 | 86 | 86 |
| 0.5 | 64 ‖ | 89 | 88 | 88 | 88 |
| 1 | 25 ‖ | 89 | 87 | 87 | 88 |
| 1 | 64 ‖ | 89 | 88 | 88 | 88 |

[*] The ratio for the selected proportion of the synthetic data. The 0.25 equals to 126 synthetic samples.
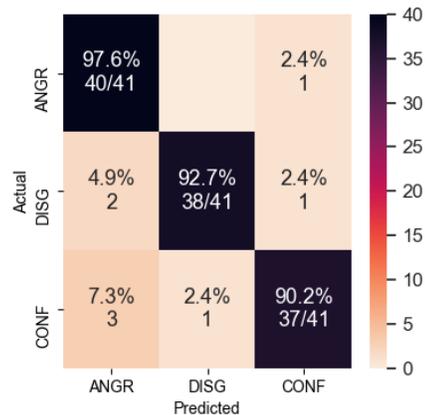


Fig. 8: Combined confusion matrix of all the folds for the FN+INC64

### C. Evaluating the Sim2Real approach on GIFGIF dataset

To evaluate the generalization of our Sim2Real approach and model, we selected an external dataset for validation, GIFGIF [28].[3] As previously noted, there are currently no video datasets with confusion samples [38]. GIFGIF [28] has video-level annotations and contains 2 of our emotions of interest ("anger" and "disgust").

This dataset is a collection of 3,858 cropped short videos with annotation scores for 17 emotions. We used GIFGIF API to get the first 400 highest-ranking videos for the "disgust" emotion. These videos were filtered down to 75 based on the following criteria: 1) contains a human face reaction video, 2) must not hold a higher score in other categories. Similarly, we chose the top 75 "anger" videos. The Arousal-Valence distribution of all these 150 samples is displayed in Fig. 9. The Arousal-Valence values are extracted using Emonet [35]. Fig. 9 suggests that this evaluation dataset is severely challenging.

For this evaluation, no additional training was performed.

[3] We also considered AffWild, EmoReact, ElderReact, but their data did not contain anger or confusion, or their annotation schemes were not directly comparable with our data (e.g., frame-based). CK+ [21], Oulu-Casia [42], and MMI [26] were also not selected since they are all acted/posed and we focus on in-the-wild interactions.

TABLE IV: Evaluation of Sim2Real effect on the model's performance on unseen data

| NET | Len | Synth[1] | Prc[2] | Rec[3] | Fs[4] | Acc[5] |
|---|---|---|---|---|---|---|
| FN+INC25 | 25 | ✗ | 74 | 64 | 67 | 64 |
| FN+INC25 | 25 | ✓ | 83 | 75 | 77 | 75 |

[1] Synthetic-Data  [2] Precision  [3] Recall  [4] F-score
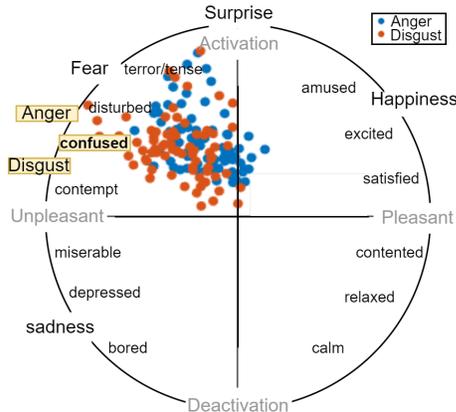[5] Accuracy



Fig. 9: Arousal-valence distribution of anger and disgust videos in GIFGIF [28] using EmoNet [35] on CMA.

We tested the two FN+IN25 models presented in Table I on this dataset. One model was trained on our real dataset, another model was pretrained on the synthetic dataset then trained on the real dataset. The GIFGIF dataset was used as a test dataset for these models. The results for this experiment are shown in Table IV. The FN+IN25 model pretrained on synthetic data achieved a 75% accuracy. Out of 150 videos, this model misclassified 7 disgust videos and 2 anger as confusion. The same model without synthetic pretraining achieved a 64% accuracy and misclassified 14 disgust and 4 anger videos as confusion. Therefore, without any additional transfer learning, we showed that our Sim2Real approach improved FN+INC25 performance on this dynamic FER task by 11%.

## IV. DISCUSSION AND LIMITATIONS

In this section, we elaborate on insights that we found while doing experiments and after analyzing the results. Our experiments showed that additional synthetic data is similar to have an extensive dataset, and the generalization of the final model is increased.

An interesting finding in our experiments was that all the models with frozen FaceNet weights performed worse than their counterparts with unfrozen weights or even I3D. This was especially the case when considering non-Caucasian samples, which was an unexpected result because FaceNet was trained on a vast face recognition dataset. This shows



Fig. 10: Sample of wrongly classified videos from our dataset

that although the FaceNet feature embedding performs well on facial recognition tasks, it may not be entirely related to the facial changes of a specific emotional expression. However, more research is needed to investigate these hypotheses.

Another interesting point was the misclassification of specific samples that were revealed after we looked deeper into the results. These samples were classified wrongly even in our best model with an average accuracy of 94%, Fig. 10 show three examples of the eight wrongly classified videos from all folds. From left to right, each column corresponds to the first, middle, and last frame. The incorrect classification of the first video might be related to the minimal movement of the face. The generated synthetic dataset that we used lacks fully static samples. Additionally, the main concept behind the proposed model was the focus on dynamic movements. The incorrect prediction of the second video relates to its head movement. The dynamic movement of the expression is done over frames involved with head movement. This adds fluctuation to the inception model's multivariant time-series input that may not relate to the emotion. The OpenFace algorithm [2] uses perspective transform to make all of the input images have a frontal face view. The addition of this step may help deal with these types of videos. However, we believe the ultimate solution is in designing a model that can predict from shorter video snippets inputs. Poor prediction of the third video relates to the movement of hand midway through the video. Face occlusion is a challenge in FER, and even though new studies focus on reducing or removing the occlusions [25], [24], their advancement has been minimal.

An interesting notion observed in the annotation of the real dataset was that annotators had trouble distinguishing between disgust and confusion in some cases. However, when the audio was played alongside the video, this confusion was resolved. This could mean that the next step for more inclusive and accurate facial expression recognition systems could incorporate audiovisual data processing.

One of the main components of this work was the generation of synthetic data. The MakeHuman application limitation highly affects this component. More advanced applications can be used for this task to generate more

realistic synthetic datasets, and to explore other variations including age, non-binary gender, or conditions impacting facial development. Another point worth mentioning is that while we understand that the relatively small size of our dataset (synthetically generated dataset plus the real human dataset gathered from YouTube and Giphy) might be a limiting factor, this is sufficient to illustrate the proposed approach as a proof-of-concept.

## V. CONCLUSIONS AND FUTURE WORK

We showed that our Sim2Real approach improves FN+INC64 performance on our dynamic FER task by 11-18%, up to 94% on our internal test dataset, and up to 75% on a previously unseen dataset, compared to the performance of the same model architecture without synthetic training data. This performance was achieved by mixed synthetic and real data training. Additionally, it was shown that the proposed FN+INC model along with our Sim2Real approach is less sensitive to dataset ethnicity bias. This study was a first step towards emotion recognition in the wild, and future work can explore applying our approach and test the trained classifiers to data gathered from real-world HRI scenarios [30]. Another notion that can be explored in the future is to observe the effects of replacing the face feature extractor model with a static FER feature extractor. An increased number of simulated human models may improve the overall accuracy, especially if they can be made with photo-realistic 3D model generating engines such as MetaHuman creator from Unreal Engine[4]. Automated animation generation from real videos can also be the next step for this study.

## REFERENCES

[1] I. Abbasnejad *et al.*, "Using synthetic data to improve facial expression analysis with 3d convolutional networks," in *ICCVW*, 2017.

[2] T. Baltrušaitis *et al.*, "Openface: An open source facial behavior analysis toolkit," in *WACV*, 2016, pp. 1–10.

[3] L. F. Barrett *et al.*, "Emotional expressions reconsidered: Challenges to inferring emotion from human facial movements," *Psychological Science in the Public Interest*, vol. 20, no. 1, pp. 1–68, 2019.

[4] M. Bastioni *et al.*, "Ideas and methods for modeling 3d human figures: the principal algorithms used by makehuman and their implementation in a new approach to parametric modeling," in *Proceedings of the 1st Bangalore Annual Compute Conference*, 01 2008, p. 10.

[5] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in *FAT*, 2018.

[6] J. Carreira and A. Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *CVPR 2017*.

[7] H. A. Dau *et al.*, "The UCR time series archive," *CoRR*, vol. abs/1810.07758, 2018.

[8] F. Dornaika and B. Raducanu, "Efficient facial expression recognition for human robot interaction," in *IWANN*, 2007, pp. 700–708.

[9] P. Ekman and E. L. Rosenberg, *What the face reveals: Basic and applied studies of spontaneous expression using the Facial Action Coding System (FACS)*. Oxford University Press, USA, 1997.

[10] H. I. Fawaz *et al.*, "Inceptiontime: Finding alexnet for time series classification," *Data Mining and Knowledge Discovery*, vol. 34, no. 6, pp. 1936–1962, 2020.

[11] T. Gerig *et al.*, "Morphable face models-an open framework," in *FG 2018*.

[12] M. Gilbert *et al.*, "Facshuman a software to create experimental material by modeling 3d facial expression," in *IVA 2018*.

[13] A. Gupta *et al.*, "Synthetic data for text localisation in natural images," in *CVPR*, 2016, pp. 2315–2324.

[14] J. Han *et al.*, "Improving face detection performance with 3d-rendered synthetic data," *arXiv:1812.07363*, 12 2018.

[15] M. Hucko *et al.*, "Confusion detection dataset of mouse and eye movements," in *UMAP*, July 2020.

[16] R. E. Jack and P. G. Schyns, "Toward a social psychophysics of face communication," *Annual Review of Psychology*, vol. 68, no. 1, pp. 269–297, Jan. 2017.

[17] M. Jangid *et al.*, "Video-based facial expression recognition using a deep learning approach," in *Advances in Computer Communication and Computational Sciences*, 2019, pp. 653–660.

[18] M. I. Jordan and T. M. Mitchell, "Machine learning: Trends, perspectives, and prospects," *Science*, no. 6245, pp. 255–260, July 2015.

[19] B.-K. Kim *et al.*, "Fusing aligned and non-aligned face information for automatic affect recognition in the wild: a deep learning approach," in *CVPR Workshop 2016*.

[20] A. Kortylewski *et al.*, "Empirically analyzing the effect of dataset biases on deep face recognition systems," in *CVPR Workshops*, 2018, pp. 2093–2102.

[21] P. Lucey *et al.*, "The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression," in *CVPR workshops*, 2010, pp. 94–101.

[22] D. McColl *et al.*, "A survey of autonomous human affect detection methods for social robots engaged in natural HRI," *JINT*, vol. 82, no. 1, pp. 101–133, Aug. 2015.

[23] F. Mueller *et al.*, "Ganerated hands for real-time 3d hand tracking from monocular rgb," in *CVPR 2018*.

[24] N. Naik and M. A. Mehta, "An improved method to recognize hand-over-face gesture based facial emotion using convolutional neural network," in *CONECCT*, July 2020.

[25] B. Nojavanasghari *et al.*, "Hand2face: Automatic synthesis and recognition of hand over face occlusions," in *ACII*, oct 2017, pp. 209–215.

[26] M. Pantic *et al.*, "Web-based database for facial expression analysis," in *ICME*, 2005, pp. 5–pp.

[27] L. Rhue, "Racial influence on automated perceptions of emotions," *SSRN Electronic Journal*, 2018.

[28] T. Rich, K. Hu, and B. Tome, Available at http://gifgif.media.mit.edu/results/.

[29] J. A. Russell, "A circumplex model of affect." *Journal of Personality and Social Psychology*, vol. 39, no. 6, pp. 1161–1178, 1980.

[30] G. Saheb Jam *et al.*, "Developing a data-driven categorical taxonomy of emotional expressions in real world human robot interactions," in *HRI*, 2021, p. 479–483.

[31] F. Schroff *et al.*, "FaceNet: A Unified Embedding for Face Recognition and Clustering," *arXiv:1503.03832 [cs]*, June 2015.

[32] G. Sharma and A. Dhall, "A survey on automatic multimodal emotion recognition in the wild," in *Advances in Data Science: Methodologies and Applications*, Aug. 2020, pp. 35–64.

[33] S. D. Sims and C. Conati, "A neural architecture for detecting user confusion in eye-tracking data," in *ICMI*, Oct. 2020.

[34] Y. Tang, "Deep learning using support vector machines," *CoRR, abs/1306.0239*, vol. 2, 2013.

[35] A. Toisoul *et al.*, "Estimation of continuous valence and arousal levels from faces in naturalistic conditions," *Nature Machine Intelligence*, vol. 3, no. 1, pp. 42–50, 2021.

[36] T. Tommasi *et al.*, "A deeper look at dataset bias," in *Domain adaptation in computer vision applications*, 2017, pp. 37–55.

[37] Z. Wang *et al.*, "Time series classification from scratch with deep neural networks: A strong baseline," in *IJCNN*, May 2017.

[38] F. I. Yasser *et al.*, "Detection of confusion behavior using a facial expression based on different classification algorithms," *ETJ*, vol. 39, no. 2A, pp. 316–325, Feb. 2021.

[39] B.-K. Yi *et al.*, "Efficient retrieval of similar time sequences under time warping," in *ICDE 1998*.

[40] Z. Yu and C. Zhang, "Image based static facial expression recognition with multiple deep network learning," in *ICMI*, 2015, p. 435–442.

[41] K. Zhang *et al.*, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, 2016.

[42] G. Zhao *et al.*, "Facial expression recognition from near-infrared videos," *Image and Vision Computing*, vol. 29, no. 9, pp. 607–619, 2011.

[4] https://www.unrealengine.com