

RealSmileNet: A Deep End-To-End Network for Spontaneous and Posed Smile Recognition

Yan Yang¹[0000-0002-6246-1748], Md Zakir Hossain^{1,2}[0000-0003-1892-831X], Tom Gedeon¹[0000-0001-8356-4909], and Shafin Rahman^{1,3,4}[0000-0001-7169-0318]

¹ The Australian National University, Canberra ACT 0200, AU

² The University of Canberra, Bruce ACT 2617, AU

³ North South University, Dhaka, Bangladesh

⁴ Data61-CSIRO, Canberra ACT 0200, AU

{u6169130, zakir.hossain, tom.gedeon}@anu.edu.au,
shafin.rahman@northsouth.edu

Abstract. Smiles play a vital role in the understanding of social interactions within different communities, and reveal the physical state of mind of people in both real and deceptive ways. Several methods have been proposed to recognize spontaneous and posed smiles. All follow a feature-engineering based pipeline requiring costly pre-processing steps such as manual annotation of face landmarks, tracking, segmentation of smile phases, and hand-crafted features. The resulting computation is expensive, and strongly dependent on pre-processing steps. We investigate an end-to-end deep learning model to address these problems, the first end-to-end model for spontaneous and posed smile recognition. Our fully automated model is fast and learns the feature extraction processes by training a series of convolution and ConvLSTM layer from scratch. Our experiments on four datasets demonstrate the robustness and generalization of the proposed model by achieving state-of-the-art performances.

1 Introduction

Facial expression recognition is a process of identifying human emotion from videos, audios, and even the texts. Understanding facial expressions is essential for various forms of communication, such as the interaction between humans and machines. Also, the development of facial expression recognition contributes to the area of market research, health care, video game testings, and so on [1]. Meanwhile, people tend to hide their natural expression in different environments. Recognising spontaneous and posed facial expressions are necessary for social interaction analysis [2] because it can be deceptive and convey diverse meanings. The *smile* is the most common and easily expressible facial display, but still very hard to recognise. Because of the recurrence and cultural reasons, the study of cognitive and computer sciences broadly investigates the recognition of spontaneous (genuine/real/felt) and posed (fake/false/deliberate) smiles [2,3,4,5,6,7,8,9,10,11,12,13,14].

Previous efforts on recognizing spontaneous and posed smiles mostly follow a feature-based approach where machine learning (ML) models perform a

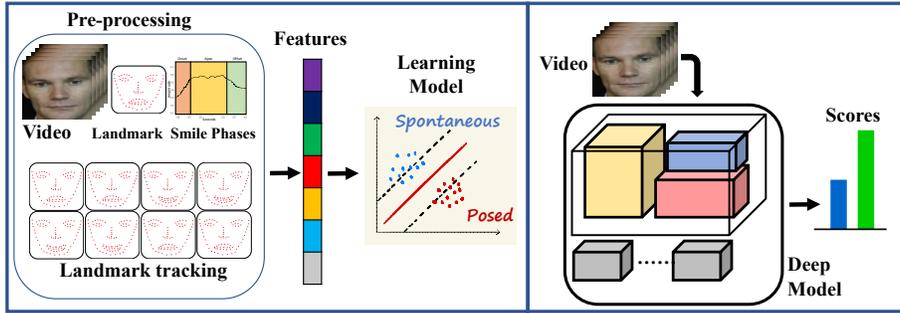


Fig. 1: Overview of different spontaneous smile recognition models. *(left)* Given a video as input, previous approaches [2,4,5,6,7,9] perform several manual or semi-automatic preprocessing steps like facial landmark detection [2,4,5,6,7], tracking [2,4,5,6,7], smile phases segmentation [2,5], and so on. across frames to calculate hand-engineered feature vectors (D-marker [2,4,15], HoG [5,9,16]), then feed the features to a learning model (SVM) for classification. The costly intermediate steps significantly increase the computation and limit the fully automatic process. *(right)* Our proposed end-to-end architecture takes video frames as input and recognizes spontaneous and posed smile by a simple forward pass.

binary classification based on the extracted visual features from a smile video [2,3,4,5,6,7,8,9,10,11,12,13,16,15]. We identify several limitations of such approaches. *(a) Manual annotation:* Many methods require manual annotation of facial landmarks for the first frame of a video [2,4,5,6,7,15]. It limits the automation of the recognition process. *(b) Landmark tracking:* Methods need to track face landmarks throughout the video [2,4,5,6,7,15]. It is a computationally expensive process, and the performance of the recognition broadly depends on it. *(c) Segmentation of temporal phases:* Some methods extract features from temporal stages of a smile (i.e., onset, apex, and offset) separately [2,5]. Automatic segmentation of a smile can be erroneous because, in many smile videos, these phases are not apparent, and methods need to assign zero values in the feature list to satisfy the constant length of the feature set. *(d) Limiting the maximum length of a smile:* Most traditional machine learning methods cannot handle the dynamic length of time series data. Traditional methods need to represent each smile by a fixed length. It decreases the robustness of the system because, in a real application, a smile video may come with variable length. *(e) Hand-engineered features:* Methods depend on hand-crafted features like D-marker [2,4,15], Histogram of Oriented Gradients (HoG) [5,9,16], Local Binary Pattern (LBP) like feature on region of interest [5,16]. The selection of such features sometimes requires extensive research and expert domain-specific knowledge [17]. Because of the issues mentioned above, traditional methods become slow, limits the automation process, and achieves poor generalization ability. Moreover, the overall performance of recognition broadly depends on the availability of many independent pre-processing steps.

Here, we propose an approach that elegantly solves the problems and encapsulates all the broken-up pieces of traditional techniques into a single, unified deep neural network called ‘*RealSmileNet*’. Our method is end-to-end trainable, fully automated, fast, and promotes real-time recognition. We employ shared convolution neural networks (CNN) layers to learn the feature extraction process automatically, Convolutional Long Short Term Memory network (ConvLSTM) [18] layers to track the discriminative features across frames and a classification sub-network to assign a prediction score. Our model adaptively searches, tracks, and learns both spatial and temporal facial feature representations across the frames in an end-to-end manner. In this way, the spontaneous smile’s recognition becomes as simple as a forward pass of the smile video through the network. In Fig. 1, we illustrate the difference between our method and the methods in the literature. Experimenting with four well-known smile datasets, we report state-of-the-art results without compromising the automation process. Our main contributions are summarized below:

- To the best of our knowledge, we propose the first end-to-end deep network for recognition of spontaneous and posed smiles.
- Our method is fully automated and requires no manual landmark annotation or feature engineering. Unlike traditional methods, the proposed network can handle variable length of smile videos leading to a robust solution.
- As a simple forward pass through the network can perform the recognition process, our approach is fast enough to promote a real-time solution.
- We present extensive experiments on four video smile datasets and achieve state-of-the-art performance.

2 Related Work

Dynamics of the spontaneous smile: The smile is the most common facial expression, and usually featured by Action Unit 6 (AU6) and Action Unit 12 (AU12) in the facial action coding system (FACS) [19]. The rise of cheek and pull of lip corners is commonly associated with a smile [19]. In terms of temporal dynamics, the smile can be segmented into the onset, apex, and offset phases. It corresponds to the facial expression variation from neutral to smile and then return to neutral. In physiological research on facial expressions, Duchenne defines the smile as the contraction of both the zygomatic major muscle and the orbicularis oculi muscle, which known as *D-Smile*. A *Non-D-smile* tends to be a polite smile where only the zygomatic muscle is contracted [17]. Recently, Schmidt et al. [15] proposed a quantitative metric called Duchenne Marker (D-Marker) to measure the enjoyment of smiles. Much research uses this (controversial) metric to recognise spontaneous and posed smiles [2,4,6,15]. Our end-to-end network for spontaneous smile recognition does not use the D-Maker feature.

Spontaneous smile engineering: The literature of spontaneous smile recognition usually follows a feature-based approach. Those methods extract features from each frame along the time dimension to construct a multi-dimensional signal. A statistical summary of the signal obtained from a smile video, such as

duration, amplitude, speed, accelerations, and symmetry, is considered in smile classification. The majority of competitive and notable research on Smile Classification relies on feature extraction by D-marker [2,4,15]. Dibeklioglu *et al.* [2] proposed a linear SVM classifier that uses the movement signal of eyelid, lip, and cheek. Mandal *et al.* [4] proposed a two-stream fusion method based on the movement of eyelid and lip and the dense optical flows with SVM. Pfister *et al.* [9] proposed feature extraction by using appearance-based local spatial-temporal descriptor (CLBP-TOP) for genuine and posed expression classifications. The CLBP-TOP is an extension of LBP, which able to extract the temporal information. Later, Wu *et al.* [5] used CLBP-TOP feature on the region of interests (eyes, lips and cheek) to train SVM for smile classifications. Valstar *et al.* [20] introduced the geometric feature-based smile classifications, using the movement of the shoulder and facial variation. We identify a few drawbacks of features based approaches. First, strongly dependence on accurate localization of the action units. Second, some approaches require manual labeling to track the changes in facial landmarks. Third, spontaneous smile recognition becomes a costly process - requiring laborious feature engineering and careful pre-processing.

End-to-end solution: In recent decades, the advancement of graphic processing units and deep learning methods allow end-to-end learning of deep neural networks that achieves unprecedented success in object re-identification [21], detection [22,23], segmentation [24] using the image, videos, and 3D point cloud data [25]. An end-to-end network takes the input (image/video/3D point cloud) and produces the output with a single forward pass. This network performs the feature engineering with the convolution layers and reduces the necessity of manual intervention and expert effort on the training process. In this vein, an end-to-end trainable deep learning model to automatically classify the genuine and posed smiles is the next step to solve the problems of feature-based solutions. With this motivation, Mandal *et al.* [16] extract features from pre-trained CNN networks (VGG Face Recognition model [26]) but eventually feed the features to a separate SVM. Instead, we propose the first fully end-to-end solution.

3 Our Approach

In contrast with available methods, an end-to-end deep learning model can provide a convenient solution by ensuring complete automation and saving computational cost after finishing the training. Because of the availability of enormous amounts of data in recent years, such end-to-end learning systems are gradually dominating research in AI. In this section, we describe an end-to-end solution for spontaneous and posed smile recognition.

3.1 Preliminaries

We consider a binary classification problem assigning labels to a sequence of images or video $\vec{X}_i = \langle \mathbf{x}_t | 1 \dots n_i \rangle$ by parameterizable models \mathcal{F}_θ where, n_i is number of frames associated with \vec{X}_i and $i \in \mathcal{T}$ and \mathcal{T} is total number of videos

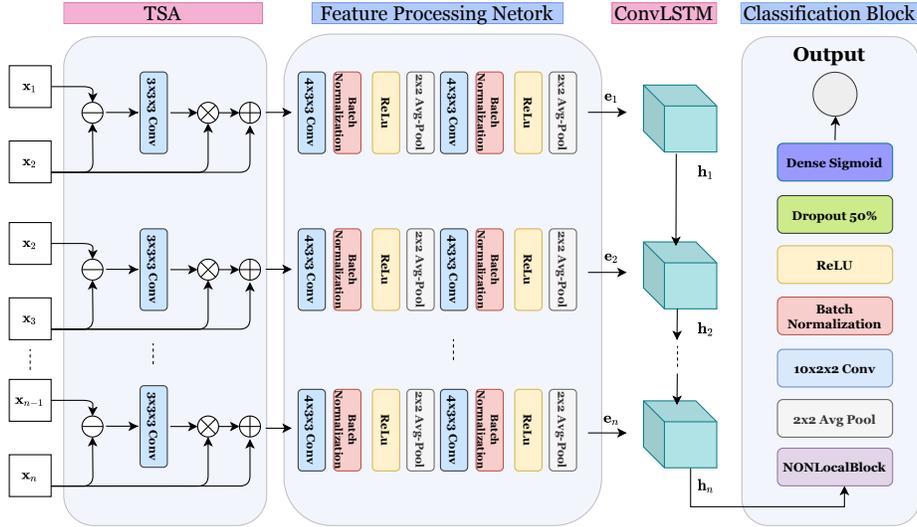


Fig. 2: Our proposed *RealSmileNet* architecture. The TSA layers guide the following feature processing network to extract discriminative facial representations. Then ConvLSTM tracks the face representation across the temporal direction to create a unified length video embedding. Finally, the classification block refines the video embedding and predicts a classification score. The number of kernels is denoted by the first number of Conv block, then the size of the kernel is followed.

in the dataset. The training dataset includes a set of tuples $\{(\vec{\mathbf{X}}_i, y_i) : i \in [0, T]\}$ where y_i represent the ground-truth label of the i th video. Here, $y_i = 1$ and $y_i = 0$ represent the class label spontaneous / posed smile respectively. Our goal is to train an end-to-end deep network model, \mathcal{F}_θ , that can assign a prediction label, \hat{y}_j , to all of \mathcal{K} testing videos, $\{\vec{\mathbf{V}}_j\}_{j=1}^{\mathcal{K}}$. We formulate \hat{y}_j as follows:

$$\hat{y}_j = \begin{cases} 1, & \text{if } \mathcal{F}_\theta(\vec{\mathbf{V}}_j) \geq 0.5 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

3.2 Architecture

We illustrate our proposed *RealSmileNet* architecture in Fig. 2. It has four components: Temporal based Spatial Attention (TSA), Feature Processing Network (FPN), ConvLSTM, and Classification block. TSA captures the motion of frames using two consecutive frames as input, FPN further processes the motion feature to generate a frame representation, ConvLSTM processes the temporal variation of frame features across different time frame to produce a video representation, and finally a classification block predicts a label for the input video. We train all components together from scratch as a single and unified deep neural network.

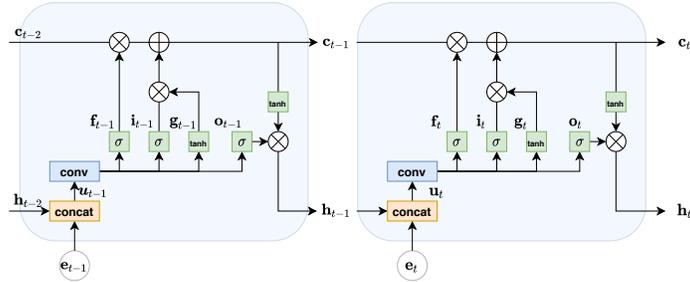


Fig. 3: State transitions of the ConvLSTM.

Temporal based Spatial Attention (TSA): We design the TSA network that learns the variation of pixels i.e. motion of a video by concentrating certain regions using residual attention. Previous research on video classification [27] showed that difference image of adjacent frames provides crude approximation of optical flow images that is helpful in action recognition. With this motivation, this network takes two consecutive frames of a video, applies a 2D convolution on difference map and performs some element-wise operations on the residual (skip) connections from \mathbf{x}_t . The overall TSA calculation is defined as:

$$TSA(\mathbf{x}_{t-1}, \mathbf{x}_t) = \left(\mathcal{C}(\mathbf{x}_t - \mathbf{x}_{t-1}) \otimes \mathbf{x}_t \right) \oplus \mathbf{x}_t \quad (2)$$

Where, \mathcal{C} represents a convolution layer that takes the difference between current frame \mathbf{x}_t and previous frame \mathbf{x}_{t-1} , \otimes and \oplus are the Hadamard product and element-wise addition respectively. The residual connections augment the output of the convolution and focus on certain area in the context of the current frame.

Feature Processing Network (FPN): We forward the output of the TSA network to the FPN layers to process the TSA features further. We design FPN with two sets of Conv, Batch Normalization, ReLU, and Avg-pooling layers. In FPN block, all the convolution layer and average pooling layer use 3x3 kernel size and 2x2 kernel size respectively. FPN learns a dense spatial feature representation of frames required to model the complex interplay of smile dynamics. During the experiment, we replace this FPN with popular ResNet18 and DenseNet like structure. However, we have achieved the best performance using our proposed implementation of an FPN. Besides, our FPN has less trainable parameters than its alternatives. In our model, TSA and FPN contribute together to get overall spatial information from frames. This representation plays the role of D-marker [2,4,15], HoG [5,9,16], LBP [5,16] of the traditional approach. The main difference is our model learns this representation, unlike conventional methods dependent on handcrafted and computationally intensive features.

ConvLSTM: We employ the ConvLSTM [18] to model the temporal dynamics of the video. We adaptively build up a saliency temporal representation of each video that contributes to the classification processes. Specifically, we concurrently learn the hidden states and input tensors by using convolution layers

instead of maintaining different weight matrixes for the hidden state and input. We visualize the state transition in Fig. 3 that performs the following operations: *input vector*, $\mathbf{u}_t = \text{concatenate}(\mathbf{e}_t, \mathbf{h}_{t-1})$, *input gate*, $\mathbf{i}_t = \sigma(\mathbf{W}_i \otimes \mathbf{u}_t \oplus \mathbf{b}_i)$, *forget gate*, $\mathbf{f}_t = \sigma(\mathbf{W}_f \otimes \mathbf{u}_t \oplus \mathbf{b}_f)$, *output gate*, $\mathbf{o}_t = \sigma(\mathbf{W}_o \otimes \mathbf{u}_t \oplus \mathbf{b}_o)$, *cell gate*, $\mathbf{g}_t = \tanh(\mathbf{W}_g \otimes \mathbf{u}_t \oplus \mathbf{b}_g)$, *cell state*, $\mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} \oplus \mathbf{i}_t \otimes \mathbf{g}_t$, *hidden state*, $\mathbf{h}_t = \mathbf{o}_t \otimes \tanh(\mathbf{c}_t)$, where, σ and \tanh are the activation function of the sigmoid and hyperbolic tangent. \otimes , \otimes and \oplus represent convolution operator, Hadamard product, and element-wise addition. *concatenate* operator stands for concatenating the augments along the channel axis. $[\cdot]_t$ and $\mathbf{W}_{[\cdot]}$ denote the element at time slot t and the corresponding weight matrix respectively. Usually, ConvLSTM updates the input gate, forget gate, cell gate, and output gate by element-wise operations [18]. But, in our case, we learn more complex temporal characteristics by concatenating the input and hidden state as the input of the convolution layer. As nearby pixels of an image are both connected and correlated, using more complex flows within the ConvLSTM cell, the convolution layer can group the local features to generate robust temporal features while preserving more spatial information.

Classification Block: The hidden state of the last frame is passed to the classification block to assign a prediction label. This block is composed of dot product NonLocal block [28], average pooling (2×2 kernel size), convolution layer (with 2×2 kernel size), batch normalization, ReLU, Dropout (with 0.5 probability), and dense layers. The NonLocal block captures the dependency between any two positions [28]. Such reliance is critical because Ekman *et al.* [8] suggested the relative position of facial landmarks (such as symmetry) contributes to the smile classification. Then, we further trim the learned embedding of the video feature through the later layers of the classification block. In this way, the classification space is well-separated for binary classification (see Fig. 5).

Loss function: Given a video as an input, \vec{X}_i , our proposed network predicts a score, $\mathcal{F}_\theta(\vec{X}_i)$, which is compared with the ground-truth y_i to calculate the weighted binary cross-entropy loss:

$$\mathcal{L}_{CE} = -\frac{1}{\mathcal{T}} \sum_{i=1}^{\mathcal{T}} \left[\alpha y_i \log(\mathcal{F}_\theta(\vec{X}_i)) + \beta (1 - y_i) \log(1 - \mathcal{F}_\theta(\vec{X}_i)) \right], \quad (3)$$

where, α and β are the weights computed as the proportion of spontaneous and posed videos in the training dataset respectively.

Inference: For j th test video, \vec{V}_j , we perform a simple forward pass through the trained network and produce a prediction score, $\mathcal{F}_\theta(\vec{V}_j)$. Then, we apply Eq. 1 to assign the predicted label, \hat{y}_j for the input.

3.3 Analysis

We analyze and visualize different aspects of our model, which allows us to address many drawbacks of the traditional approaches.



Fig. 4: Visualization of FPN features for spontaneous (top two rows) and posed (bottom two rows) smiles from UVA-NEMO [2] by using score-cam [29,30]. Keeping equal temporal distance from each other, sample frames are selected for this visualization. The more 'warm' the color, the more important the area becomes during classification.

Our model automatically learns discriminative smile features that replace traditional handcrafted features. This learning process does not require manual landmark initialization and their tracking through the video. Our ConvLSTM block enables us to track the learned features automatically until the last frame. The iterative learning process of ConvLSTM does not have any restriction on the maximum length of smile videos, unlike traditional methods requiring maximum fixed video length. Our ConvLSTM block effectively manages the temporal aspect of features in the time dimension, which performs the role of face landmark tracking of other methods. Our classification block, coupled with the ConvLSTM learns to classify time series data. But, using the SVM like classifiers, which are commonly used in the area, are not an excellent fit to classify similar data. Therefore, traditional models perform hand-engineering to make the data fit for SVM. Instead, in our model, every component of traditional methods are embedded in the unified deep network. Thus, once learned, our model handles the intermediate process through a forward pass as a single unit. Such a strategy simplifies the process because that parallel implementation is easy for a deep learning model.

Visualization: In Fig. 4, we visualize the importance of different facial regions across various frames. We extract features (after FPN layers) and blend them on the input frame. This shows that our model extracts features where it finds discriminative information. One can notice that our model puts less emphasis on neutral faces (by assigning cooler color on the heatmap) because those frames have no role in the context of spontaneous or posed smile recognition. Moreover, the starting and ending frames (roughly, onset and offset regions) are promising

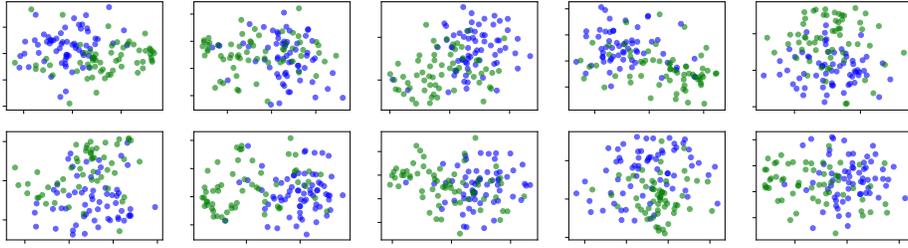


Fig. 5: 2D tSNE [31] visualization of smile video features extracted from the classification block (after ReLU layer) of our proposed model. **Blue** and **Green** represent posed and spontaneous smiles respectively. Each plot shows the test fold features from the 10-fold cross-validations data of UvA-NEMO database. Here, spontaneous and posed smiles are reasonably well-separated to recognize using the classification block.

to be the most discriminative for posed smiles, whereas middle frames (apex regions) are important for spontaneous smiles. To further visualize the feature embedding, we plot learned video features (from the output of ReLU layer of classification block (please see Fig. 2)) of the test fold belonging to the UvA-NEMO dataset in Fig. 5. We notice that spontaneous and posed smile features are properly separated for classification.

4 Experiment

4.1 Setup

Dataset: In this paper, we experiment on four popular smile datasets: Here, we briefly describe the data statistics. **(a) UVA-NEMO Smile Database** [2]: This dataset is recorded in 1920×1080 pixels at a rate of 50 frames per second. It composed of 597 spontaneous and 643 posed smile videos. The length of videos distributed from 1.8 seconds to 14.2 seconds. It contains over 400 participants (185 females and 215 males) with ages from 8 years to 76 years. There are 149 young people and 251 adults. **(b) BBC database**⁵ [2,5] This dataset contains 20 videos, recorded in 314×286 pixels with 25 frames per second. There are 10 spontaneous and 10 posed smiles. **(c) MMI Facial Expression Database** [32]: They provided spontaneous and posed facial expressions separately including 38 labeled posed smiles. Apart from these posed smiles, we identified 138 spontaneous and 49 posed smile videos from 9 and 25 participants, respectively. The age of participants ranges from 19 to 64. All of the videos contain frontal recordings. The part of the spontaneous smile is in 640×480 pixels at 29 frames per second, and the posed smile part is recorded in 720×576 pixels with 25 frames per second. **(d) SPOS database** [9]: It provides both gray and near-infrared sequences of images in 640×480 resolution with 25 frames per second.

⁵ <https://www.bbc.co.uk/science/humanbody/mind/surveys/smiles/>

Database	Video Spec.		Number of Videos		Number of Subjects	
	Resolution	FPS	Genuine	Posed	Genuine	Posed
UVA-NEMO	1920 x 1080	50	597	643	357	368
BBC	314 x 286	25	10	10	10	10
MMI	720 x 576	29	138	49	9	25
	640 x 480	25				
SPOS	640 x 480	25	66	14	7	7

Table 1: Summary of the smile datasets.

We use gray images in our experiments. The face region of each image has been cropped by the database owners. There are 66 spontaneous smiles, and 14 posed smiles from 7 participants. The age of participants distributed from 22 to 31, while 3 of them are male. Table 1 provides a summary of these datasets.

Train/test split: We use the standard train/test split protocol from [2] for UVA-NEMO database. Following the settings from [2,5], we perform 10-fold, 7-fold, and 9-fold cross-validation for BBC, SPOS, and MMI datasets, respectively, while maintains no subject overlap between training and testing folds.

Evaluation Processes: We have evaluated our model with prediction accuracy. The accuracy is the proportion of test data that is correctly predicted by our model. We report the average result of running ten trials.

Implementation details⁶: We train our model for 60 epoch with the mini-batch size 16. To optimize network parameters, we use Adam optimizer with a learning rate 10^{-3} and decay 0.005. We employ weighted binary cross-entropy loss for training where the weight is the ratio between spontaneous smiles and posed smiles in training data. To prepare the video to be manageable for the network, we sample 5 frames per second, crop the face using DLIB library [33] and resize each frame into the dimension 48×48 , which are purely automatic processes. We validate the sensitivity of these design choices in experiments. We implement our method using the *PyTorch* library [34].

4.2 Recognition Performance

In this subsection, we will compare our performance with other models, will show an ablation study, will design choice sensitivity, and will analyze the robustness of our approach.

Benchmark Comparisons: In Table 2, we compare our performance of spontaneous and posed smile classification with previously published approaches using four popular datasets. We divide all methods into two categories: semi-automatic and fully-automatic. Semi-automatic methods manually annotate facial landmark locations of the first frame of the video. In contrast, fully-automatic

⁶ Code and evaluation protocols available at: <https://github.com/Yan98/Deep-learning-for-genuine-and-posed-smile-classification>

Method	Process Type	UVA-NEMO	MMI	SPOS	BBC
Cohn'04 [7]	Semi-automatic	77.3	81.0	73.0	75.0
Dibeklioglu'10 [6]	Semi-automatic	71.1	74.0	68.0	85.0
Pfister'11[9]	Semi-automatic	73.1	81.0	67.5	70.0
Wu'14 [5]	Semi-automatic	91.4	-	79.5	90.0
Dibeklioglu'15 [2]	Semi-automatic	89.8	88.1	77.5	90.0
Mandal'17 [4]	Semi-automatic	80.4	-	-	-
Mandal'16 [16]	Fully-Automatic	78.1	-	-	-
Ours	Fully-Automatic	82.1	92.0	86.2	90.0

Table 2: Benchmark comparison of methods. ‘-’ means unavailable result.

methods require no manual intervention in the complete process. Our model successfully beats all methods in MMI, SPOS, and BBC datasets. For UVA-NEMO, we outperform the automatic method [16]. However, Wu *et al.* reported the best performance on the UVA-NEMO dataset.[5] For all these experiments, the same video and subjects are used during testing. But, being not end-to-end, previous methods apply many pre-processing steps (borrowed from different work) that are not consistent across methods. For example, the performance of [2,6,7,9] are adopted from the work [2] where same landmark initialization tracker [35], face normalization, etc. are used. However, the accuracy of [5] and [4] are reported from the original papers where they employed a different manual initialization and tracker [36,37]. Moreover, the number and position of landmarks used are also different across models. Because of these variations, performance of the semi-automatic methods are difficult to compare in a common framework. The automatic method [16] is our closest competitor because of the lack of requirement of landmark initialization or tracker and their best result can be gained in a fully automatic way. However, their feature extraction and learning model are still separated. Besides, to manage the variable length of video frames, they apply a frame normalization process (using fixed number of coefficients of Discrete Cosine Transform) to create fixed length videos. Our proposed model is fully-automated and end-to-end trainable as a single deep learning model. It requires no landmark detection, tracking, frame normalization to a fixed-length, etc.

Design Choice Sensitivity: In Fig. 6, we report the sensitivity of the design choice of our method for different numbers of frames per second (FPS) and resolutions of the input frames. For all combination of FPS (1, 3, 5 and 7) and resolution (48×48 , 64×64 , 96×96 and 112×112) choices, we find FPS = 5 and resolution = 48×48 achieves the maximum performance. Note that image resolution is important because it decides the type of visual features extracted by the CNN layers. For example, a low resolution (48×48) lets the CNN kernels (of size 3×3) extracts coarse features that are the most discriminative for smile classification. Similarly, the choice of FPS also interacts with ConvLSTM layers to track the change of smile features across frames. The FPS = 5 and resolution = 48×48 is the trade-off to maximize the performance.

Ablation studies: In Tab. 3, we perform ablation studies by replacing part of our proposed network with a suitable alternative. Our observations from the ab-

Method	UVA-NEMO	MMI	SPOS	BBC
No TSA	78.5	92.0	81.5	80.0
miniResNet	73.8	84.9	80.5	90.0
miniDenseNet	77.0	71.2	82.2	90.0
No Weighted Loss	80.6	91.7	82.2	90.0
Softmax Function	79.2	71.6	82.2	70.0
Ours	82.1	92.0	86.2	90.0

Table 3: Ablation study. We experiment adding or removing parts of proposed method with reasonable alternatives.

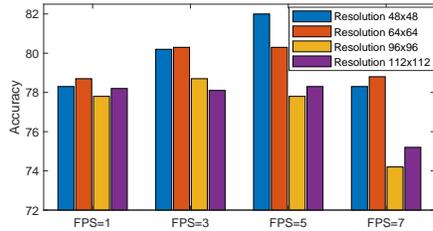


Fig. 6: Varying image size and FPS on UVA-NEMO dataset.

lation studies are the following. **(1)** We remove the TSA block from our model and directly forward the frame to FPN for feature extraction. In this situation, the network could not get optical flow information and motion from the spatio-temporal region. Smile features based on the difference of consecutive frames extract more discriminative features than a single frame. Thus, without using the TSA block, the performance degrades, especially on UVA-NEMO and SPOS datasets. **(2)** Now, we experiment on the alternative implementation of the FPN network, for example, ResNet12 [38] (composed by 3 layers) or DenseNet [39] (with growth rate 2 and 6,12, 24 and 16 blocks in each layer). We use a relatively small version of that well-known architecture because smile datasets do not contain enough instances (in comparison to large scale ImageNet dataset [40]) to train large networks. Alternative FPNs could not outperform our proposed FPN implementation. One reason could be that the smaller version of those popular architectures is still larger than our proposed FPN, and available smile data overfits the networks. Another reason is that, we could not use pre-trained weights for the alternatives, because of different input resolutions. **(3)** We try without the weighting version of the loss of Eq. 3, i.e. $\alpha = \beta = 1$. This impacts the performance of UVA-NEMO, MMI, and SPOS dataset because of the large imbalance in number of training samples of spontaneous and posed smiles. **(4)** We replace the dense sigmoid at the last layer of the classification block with softmax. In our sigmoid based implementation, we use one neuron at the last layer to predict a score within $[0, 1]$ and apply Eq. 1 for inference. For the softmax case, we add two neurons at the last layer, which increases the number of trainable parameters. We notice that the softmax based network does not perform better than our proposed sigmoid based case. This observation is in line with the recommendation of [41] that Softmax is better for the multi-class problem rather than a binary class case. **(5)** The performance of our final model outperforms all ablation alternatives consistently across datasets.

Effect of Age and Genders: To illustrate our approach’s robustness, we analyze the effect of model prediction for the different subject groups concerning age and gender. Firstly, we experiment on whether biologically similar appearance inserts any bias in the prediction. For this, we train and test our proposed model with only male/female/adult/young separately. In Fig. 7(*left*), we show

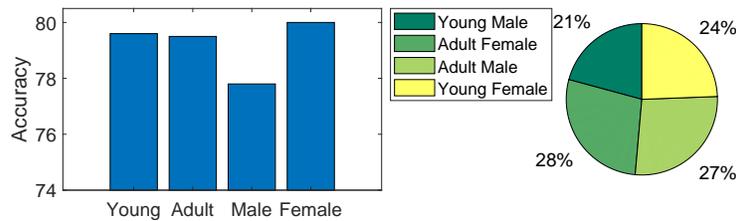


Fig. 7: (*left*) The accuracy of the model trained by individual group. (*right*) The normalized distribution of the total number of wrong predictions among different subject groups.

the results on each subgroup: male, female, adult, and young. We notice that the performance is similar to our overall performance reported in Table 2. This shows that our training has no bias on age- and gender-based subgroups. Secondly, in Fig. 7 (*right*), we visualize the normalized distribution of the total number of wrong predictions among adults/young and males/females. We find that there is no significant bias in misprediction distribution. In other words, as the mispredictions are similar across different groups, our model does not favor any particular group.

Cross-domain Experiments: We also perform experiments across datasets and subject groups. While training with UVA-NEMO and testing with BBC, MMI, and SPOS dataset, we get 80%, 92.5% & 82.5% accuracy, respectively. Moreover, training with adults and testing with young subjects got 75.9%, and conversely 74.8% accuracy. Again, training on female then testing on male subjects got 74.2% and conversely 75.1% accuracy. These experiments indicate the robustness of our method.

4.3 Discussion

Time Complexity: The processes of facial landmarks tracking/detection followed by handcrafted feature extraction are usually very computationally expensive. As evidence, when we re-implement D-marker feature-based approaches with DLIB library [33] to face normalization and facial landmark detection, it requires more than 28 hours for the processes using a single NVIDIA V100 GPU and one Intel Xeon Cascade Lake Platinum 8268 (2.90GHz) CPU. Although the training is efficient and effective, the pre-processing pipelines are costly. However, for our end-to-end learning models, the whole processing only spends up to eight hours using the same system configuration, which significantly saves time.

Human Response vs. Our Model: Several recent works estimate the ability of the average human to recognize spontaneous smiles [11,12]. In Fig. 8, we show the comparison of our work with human responses using the experiment set-up mentioned in [11]. In an experiment with 26 human and 20 videos selected from the UVA-NEMO dataset, Hossain *et al.* [11] reported a 59% average accuracy of the human. In this set-up, our trained *RealSmileNet* (without using any of

those 20 videos and their subjects during training) successfully achieves 95% of accuracy on the same test set. In another experiment with 36 humans and 30 videos from the same UVA-NEMO dataset, Hossain *et al.* [11] reported 70% for humans, whereas our proposed model achieves 90% accuracy. These comparisons show that our end-to-end model is already capable of beating human-level performance.

Limitations: One notable drawback of deep learning-based solutions is the dependence on large-scale balanced data. Being a deep model, our proposed model has also experienced this issue during training with UVA-NEMO dataset, which includes smiles of subjects with a wide range of ages, e.g., child (0–10 years), young (11–18 years), adult (19–69 years), aged people (≥ 70 years). However, among the 1,240 videos, the distributions are 20.24%, 18.47%, 60.16%, and 1.45% respectively. The imbalanced/skewed distribution usually cannot be well modeled in the deep models [42], and can lead to unexpected bias in the kernel of the convolution layer. Here, our model performs less well than the semi-automatic method of Wu *et al.* [5] (See Table 2). In future, one can collect more data to handle such shortcomings.

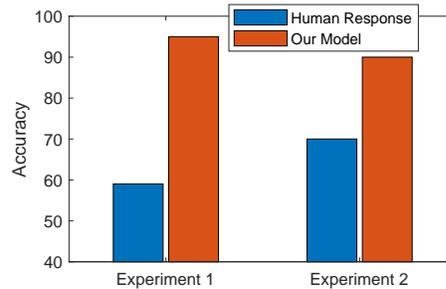


Fig. 8: Human response vs. our model

5 Conclusion

Traditional approaches for recognizing spontaneous and posed smiles depend on expert feature engineering, manual labeling, and numerous costly pre-processing steps. In this paper, we introduce a deep learning model, *RealSmileNet*, to unify the broken-up pieces of intermediate steps into a single, end-to-end model. Given a smile video as input, our model can generate an output prediction by a simple forward pass through the network. Our proposed model is not only fast but also removes the expert intervention (hand engineering) in the learning process. Experimenting with four large scale smile datasets, we establish state-of-the-art performances on three datasets. Our experiment previews the applicability of *RealSmileNet* to many real-world applications like polygraphy, human-robot interactions, investigation assistance, and so on.

Acknowledgment: This research was supported by the National Computational Infrastructure (NCI), Canberra, Australia.

References

1. Li, S., Deng, W.: Deep facial expression recognition: A survey. *IEEE Transactions on Affective Computing* (2020)
2. Dibeklioglu, H., Salah, A., Gevers, T.: Recognition of genuine smiles. *IEEE Transactions on Multimedia* **17** (2015) 279–294
3. Frank, M., EKMAN, P.: Not all smiles are created equal: The differences between enjoyment and nonenjoyment smiles. *Humor-international Journal of Humor Research - HUMOR* **6** (1993) 9–26
4. Mandal, B., Ouarti, N.: Spontaneous vs. posed smiles - can we tell the difference? *International Conference on Computer Vision and Image Processing* **460** (2017)
5. Wu, P., Liu, H., Zhang, X.: Spontaneous versus posed smile recognition using discriminative local spatial-temporal descriptors. In: *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2014, Florence, Italy, May 4-9, 2014, IEEE* (2014) 1240–1244
6. Dibeklioglu, H., Valenti, R., Salah, A.A., Gevers, T.: Eyes do not lie: spontaneous versus posed smiles. In Bimbo, A.D., Chang, S., Smeulders, A.W.M., eds.: *Proceedings of the 18th International Conference on Multimedia 2010, Firenze, Italy, October 25-29, 2010, ACM* (2010) 703–706
7. Cohn, J.F., Schmidt, K.L.: The timing of facial motion in posed and spontaneous smiles. *International Journal of Wavelets, Multiresolution and Information Processing* **02** (2004)
8. Ekman, P., Hager, J., Friesen, W.: The symmetry of emotional and deliberate facial actions. *Psychophysiology* **18** (1981) 101 – 106
9. Pfister, T., Li, X., Zhao, G., Pietikäinen, M.: Differentiating spontaneous from posed facial expressions within a generic facial expression recognition framework. In: *IEEE International Conference on Computer Vision Workshops, ICCV 2011 Workshops, Barcelona, Spain, November 6-13, 2011, IEEE Computer Society* (2011) 868–875
10. Hossain, M.Z., Gedeon, T.D.: An independent approach to training classifiers on physiological data: An example using smiles. In Cheng, L., Leung, A.C.S., Ozawa, S., eds.: *Neural Information Processing, Cham, Springer International Publishing* (2018) 603–613
11. Hossain, M.Z., Gedeon, T., Sankaranarayana, R.: Using temporal features of observers’ physiological measures to distinguish between genuine and fake smiles. *IEEE Trans. Affect. Comput.* **11** (2020) 163–173
12. Hossain, M.Z., Gedeon, T.: Discriminating real and posed smiles: human and avatar smiles. In Brereton, M., Soro, A., Vyas, D., Ploderer, B., Morrison, A., Waycott, J., eds.: *Proceedings of the 29th Australian Conference on Computer-Human Interaction, OZCHI 2017, Brisbane, QLD, Australia, November 28 - December 01, 2017, ACM* (2017) 581–586
13. Hossain, M.Z., Gedeon, T.: Observers’ physiological measures in response to videos can be used to detect genuine smiles. *Int. J. Hum. Comput. Stud.* **122** (2019) 232–241
14. Gao, R., Islam, A., Gedeon, T., Hossain, M.Z.: Identifying real and posed smiles from observers’ galvanic skin response and blood volume pulse. *The 27th International Conference on Neural Information Processing, LNCS*, (2020)
15. Schmidt, K., Bhattacharya, S., Denlinger, R.: Comparison of deliberate and spontaneous facial movement in smiles and eyebrow raises. *Journal of nonverbal behavior* **33** (2009) 35–45

16. Mandal, B., Lee, D., Ouarti, N.: Distinguishing posed and spontaneous smiles by facial dynamics. In Chen, C., Lu, J., Ma, K., eds.: *Computer Vision - ACCV 2016 Workshops - ACCV 2016 International Workshops*, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part I. Volume 10116 of *Lecture Notes in Computer Science.*, Springer (2016) 552–566
17. Duchenne, B.: *The Mechanism of Human Facial Expression*. Cambridge: Cambridge University Press (1990)
18. Shi, X., Chen, Z., Wang, H., Yeung, D., Wong, W., Woo, W.: Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R., eds.: *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*, December 7-12, 2015, Montreal, Quebec, Canada. (2015) 802–810
19. Littlewort, G., Frank, M., Lainscsek, C., Fasel, I., Movellan, J.: Automatic recognition of facial actions in spontaneous expressions. *Journal of Multimedia* **1** (2006)
20. Valstar, M.F., Pantic, M., Ambadar, Z., Cohn, J.F.: Spontaneous vs. posed facial behavior: automatic analysis of brow actions. In Quek, F.K.H., Yang, J., Massaro, D.W., Alwan, A.A., Hazen, T.J., eds.: *Proceedings of the 8th International Conference on Multimodal Interfaces, ICMI 2006*, Banff, Alberta, Canada, November 2-4, 2006, ACM (2006) 162–170
21. Zheng, L., Yang, Y., Hauptmann, A.G.: Person re-identification: Past, present and future. *CoRR* **abs/1610.02984** (2016)
22. Rahman, S., Khan, S.H., Porikli, F.: Zero-shot object detection: Joint recognition and localization of novel concepts. *International Journal of Computer Vision* **128** (2020) 2979–2999
23. Rahman, S., Khan, S., Barnes, N., Khan, F.S.: Any-shot object detection. *arXiv preprint arXiv:2003.07003* (2020)
24. Rochan, M., Rahman, S., Bruce, N.D.B., Wang, Y.: Weakly supervised object localization and segmentation in videos. *Image Vis. Comput.* **56** (2016) 1–12
25. Yang, H., Shi, J., Carlone, L.: TEASER: fast and certifiable point cloud registration. *CoRR* **abs/2001.07715** (2020)
26. Parkhi, O.M., Vedaldi, A., Zisserman, A.: Deep face recognition. In Xie, X., Jones, M.W., Tam, G.K.L., eds.: *Proceedings of the British Machine Vision Conference 2015*, BMVC 2015, Swansea, UK, September 7-10, 2015, BMVA Press (2015) 41.1–41.12
27. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1*. NIPS’14, Cambridge, MA, USA, MIT Press (2014) 568–576
28. Wang, X., Girshick, R.B., Gupta, A., He, K.: Non-local neural networks. In: *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018*, Salt Lake City, UT, USA, June 18-22, 2018, IEEE Computer Society (2018) 7794–7803
29. Wang, H., Wang, Z., Du, M., Yang, F., Zhang, Z., Ding, S., Mardziel, P., Hu, X.: Score-cam: Score-weighted visual explanations for convolutional neural networks. *CoRR* (2019)
30. Ozublak, U.: Pytorch cnn visualizations. <https://github.com/utkuozbulak/pytorch-cnn-visualizations> (2019)
31. Rauber, P.E., Falcão, A.X., Telea, A.C.: Visualizing time-dependent data using dynamic t-sne. In Bertini, E., Elmqvist, N., Wischgoll, T., eds.: *Eurographics Conference on Visualization, EuroVis 2016, Short Papers*, Groningen, The Netherlands, 6-10 June 2016, Eurographics Association (2016) 73–77

32. Valstar, M., Pantic, M.: Induced disgust, happiness and surprise: An addition to the mmi facial expression database. *Proc. Int'l Conf. Language Resources and Evaluation, Workshop EMOTION (2010)* 65–70
33. King, D.E.: Dlib-ml: A machine learning toolkit. *J. Mach. Learn. Res.* **10** (2009) 1755–1758
34. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In Wallach, H.M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E.B., Garnett, R., eds.: *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada.* (2019) 8024–8035
35. Tao, H., Huang, T.S.: Explanation-based facial motion tracking using a piecewise bézier volume deformation model. In: *1999 Conference on Computer Vision and Pattern Recognition (CVPR '99)*, 23-25 June 1999, Ft. Collins, CO, USA, IEEE Computer Society (1999) 1611–1617
36. Nguyen, T.D., Ranganath, S.: Tracking facial features under occlusions and recognizing facial expressions in sign language. In: *8th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2008)*, Amsterdam, The Netherlands, 17-19 September 2008, IEEE Computer Society (2008) 1–7
37. Tomasi, C., Kanade, T.: Detection and tracking of point features. Technical report, Carnegie Mellon University, Technical Report CMU-CS-91-132 (1991)
38. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, IEEE Computer Society (2016) 770–778
39. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, IEEE Computer Society (2017) 2261–2269
40. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)* **115** (2015) 211–252
41. Nwankpa, C., Ijomah, W., Gachagan, A., Marshall, S.: Activation functions: Comparison of trends in practice and research for deep learning. *CoRR* **abs/1811.03378** (2018)
42. Vandal, T., Kodra, E., Dy, J.G., Ganguly, S., Nemani, R.R., Ganguly, A.R.: Quantifying uncertainty in discrete-continuous and skewed data with bayesian deep learning. In Guo, Y., Farooq, F., eds.: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, ACM (2018) 2377–2386