# Awakening Facial Emotional Expressions in Human-Robot

Yongtong Zhu[1], Lei Li[1], Iggy Qian[2], WenBin Zhou[3], Ye Yuan[1], Qingdu Li[1], Na Liu[1,*], Jianwei Zhang[4]

*Abstract*— The facial expression generation capability of humanoid social robots is critical for achieving natural and human-like interactions, playing a vital role in enhancing the fluidity of human-robot interactions and the accuracy of emotional expression. Currently, facial expression generation in humanoid social robots still relies on pre-programmed behavioral patterns, which are manually coded at high human and time costs. To enable humanoid robots to autonomously acquire generalized expressive capabilities, they need to develop the ability to learn human-like expressions through self-training. To address this challenge, we have designed a highly biomimetic robotic face with physical-electronic animated facial units and developed an end-to-end learning framework based on KAN (Kolmogorov-Arnold Network) and attention mechanisms. Unlike previous humanoid social robots, we have also meticulously designed an automated data collection system based on expert strategies of facial motion primitives to construct the dataset. Notably, to the best of our knowledge, this is the first open-source facial dataset for humanoid social robots. Comprehensive evaluations indicate that our approach achieves accurate and diverse facial mimicry across different test subjects.

## I. INTRODUCTION

The ability of humanoid social robots [1] to express emotions is crucial, as it not only enhances the naturalness of human-robot interactions but also improves emotional resonance and accuracy during these exchanges. Central to this capability is the skill of facial mimicry [2]–[4], which is essential for humanoid robots to effectively learn and replicate emotional responses. However, the current capabilities of these robots in mimicking expressions are primarily reliant on manually pre-set behavior patterns [5]–[8], a method that is both time-consuming and labor-intensive, thereby limiting the robots' flexibility and adaptability.

Current humanoid social robots face substantial challenges in efficiently and adaptively mimicking human facial expressions, primarily due to the absence of a universal learning framework for mastering these skills. Traditional approaches [9]–[11] often rely on predefined behavior patterns to search for the closest matching expressions; however, this method frequently fails to capture the rich diversity of human emotions. Although some frameworks [12], [13] based on self-supervised learning exist, they tend to exhibit low learning efficiency and struggle to mitigate the uncanny valley effect, which arises when a robot's appearance and

[1]Institute of Machine Intelligence, University of Shanghai for Science and Technology, Shanghai 200093, China
[2]College Of Information Engineering, Zhejiang University of Technology, Hangzhou 310023, China
[3]Shanghai Droid Robot Co., Ltd. Shanghai 200433, China
[4]Department of Informatics, University of Hamburg, 20146 Hamburg, Germany
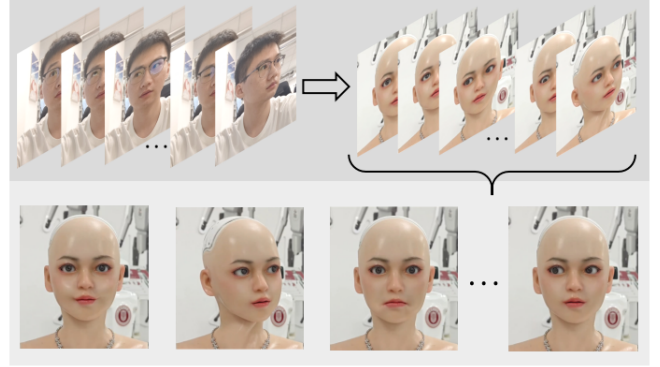*Corresponding author email: liuna@usst.edu.cn

Fig. 1: Rena is a general animatronic robotic face for emotional expressions. The robot achieves this by learning the correspondence between its facial feature representations and the control of servos. The entire learning process relies on the robot issuing motor commands based on specific expert strategies. The figure illustrates the robot's capability to replicate a variety of human expressions, with a particular emphasis on the realism of these expressions.

behavior closely resemble human characteristics but remain imperfectly aligned, leading to feelings of discomfort among humans. Consequently, these limitations restrict the practical application of such methods in human-robot interactions, particularly in contexts that demand a high degree of naturalness and accuracy in emotional expression.

In this paper, we introduce our robot named Rena, equipped with a 25 DoF control structure and highly biomimetic facial units. Furthermore, we present an innovative end-to-end learning framework based on KAN [14] and attention [15] mechanisms that focuses on key visual features to enhance facial imitation learning. Notably, due to the KAN's advantage of low parameter counts, our deployed robotic system achieves an inference speed of up to 50 frames per second. Moreover, we have carefully designed an automated data collection system based on facial motion primitives. Importantly, the constructed dataset aligns naturally with facial motion rules, enabling the model to learn expressions that effectively mitigate the uncanny valley effect.

Our experimental results demonstrate that our method surpasses existing self-supervised learning paradigms in the domain of facial mimicry, as illustrated in Fig. 1. We trained our approach on a dataset containing 9,000 samples and achieved a servo control error rate of only 4.4% on an independent test set. Additionally, we organized a substantial

number of participants to evaluate the effectiveness and generalizability of our method. Furthermore, our algorithm exhibits a real-time response time of just 0.02s on the robot, showcasing its exceptional real-time responsiveness, which is critical for interactive control.

Our contributions can be summarized as follows:

1) We present a facial emotion robot featuring highly biomimetic facial units and a high degree-of-freedom facial structure.
2) We propose an end-to-end learning framework based on KAN and feature attention mechanisms.
3) We have developed an automated data collection system grounded in expert strategies for facial action primitives.

## II. RELATED WORK

### A. Animatronic Robotics Face

Designing robots capable of mimicking human expressions serves as a foundational aspect of this work. Previous designs [12] typically featured a lower degree of freedom (DoF), which constrained their ability to present diverse facial expressions. Additionally, achieving a variety of expressions while maintaining a natural and lifelike appearance within limited structural space remains a significant challenge [13]. Recent research on speech-driven animatronic facial expression systems [16] employs a pneumatically driven system with 16 DoF; however, its limited range restricts expressive capability and lacks head movement flexibility. The Eva robot [12], with 22 DoF, also falls short due to its absence of highly biomimetic facial units, such as artificial eyeballs and skin. In contrast, our robot demonstrates superior expressiveness with 25 DoF combined with advanced biomimetic facial units.

### B. Synthetic Video Generation and Animation

Video synthesis animation is a crucial task in enabling emotional expression for 2D digital humans [17]–[19], involving the learning of mappings from the source image domain to the target image domain [20]–[22]. In contrast, our work focuses on learning the mapping from 2D facial image domains to the target robotic motion space. Recent research on speech-driven animatronic facial expression systems [16] has implemented a two-stage approach: first driving a 2D digital human face through voice inputs, followed by controlling a robot. However, its mapping, or redirection method, still relies on manually programmed predefined patterns, which limits the generalizability of its expression generation. Our end-to-end learning framework, by contrast, overcomes this limitation, offering a more flexible and generalizable approach to expression generation.

### C. Imitation Learning

Imitation learning, a method that involves learning new tasks by observing and mimicking others' behaviors, has been widely applied in areas such as robotic arm task execution [23]–[27]. In the field of facial emotion robots, imitation learning enables the robots to learn from human facial expressions, facilitating more natural and human-like emotional interactions. Recent studies utilizing imitation learning frameworks, including Eva 1.0, Eva 2.0, and XIN-REN [12], [13], [28], employ a key points displacement tracking approach. However, this method suffers from poor generalization capabilities and is difficult to fine-tune. A significant issue with these approaches is their reliance on random expressions as data sources, which inevitably leads to the uncanny valley effect. In contrast, our robot adopts a different strategy in dataset construction, inherently avoiding the uncanny valley effect. Using a carefully designed facial expression dataset, our method ensures that the learned expressions are natural and highly generalizable.

## III. DESIGN

Our robot design, as depicted in Fig. 2, incorporates highly realistic biomimetic facial modules and a high-degree-of-freedom structural control module. The structural control module consists of two main components: the head-eye motion unit and the mouth motion unit.
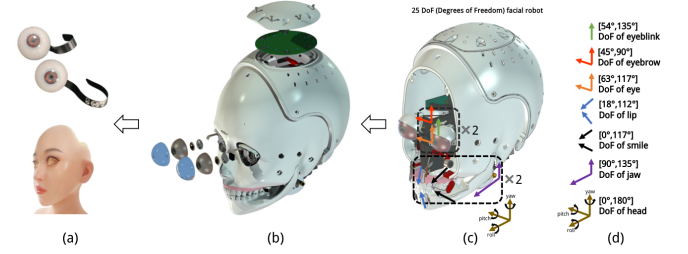


Fig. 2: The exploded view of the Rena robot is shown in (b), where (c) is a cross-sectional view of the robot structure, showing the 12 DoF motion trajectory of the mouth motion unit and the 8 DoF motion trajectory of the eyebrow area. (a) shows our bionic silicone skin and bionic eyeball.

### A. Structural Control Module

The head-eye motion unit is divided into the eyebrow-eye region and the head rotation region. Each eyeball and eyebrow are actuated by two servo motors connected via articulated linkages, which pull the external skin to achieve eyebrow raises, frowns, and horizontal and vertical eye movements.

The mouth motion unit comprises a total of 12 servo motors, which control the upper lip, lower lip, both corners of the mouth, and the jaw opening. Fig. 2(c) provides details on the control mechanisms of the mouth servos. In particular, the servo motors for the corners of the mouth and the jaw are coupled for synchronized movement, which poses significant challenges for traditional manual programming of facial expressions.

In the control module, the driving motors exclusively use the FEELECH-STS3032 servo model. This model is a closed-loop servo with a motion range of 180°. Considering the mechanical constraints across the overall structural design, we have additionally imposed specific motion limits on each servo within distinct control regions, as shown in

Fig. 2(d). This ensures that the actuating mechanisms driven by the various servos avoid interference while generating a wide range of expressions.

### B. Biomimetic Module

Facial emotion robots require not only precise motion structure design but also realistic bionic appearance to promote a more humanoid communication experience. Rena robot achieves this as shown in Fig. 2(a). The bionic eyeball we designed achieves unprecedented realism, even replicating the fine details of the blood vessels in the eyeball. This may make us forget that we are communicating with a robot during the interaction. Similarly, the skin is made of bionic silicone through a precise molding process, and its touch is almost the same as that of a real person. Therefore, Rena's appearance, composed of bionic skin and eyeballs, can provide users with an excellent experience during human-robot interaction.

## IV. PROPOSED APPROACH

We propose a learning-based framework for controlling a facial emotion robot to mimic various human facial expressions. First, an automated data collection system is employed to collect data in bulk form. Based on these data, the framework learns the mapping from robot facial expressions to servomotor control commands.

### A. Representation of Facial Expression

We represent facial expressions using the facial blendshape coefficients, a method that research has shown to be scientifically valid and effective [29]–[31]. This approach provides a standardized and normalized representation that can be used directly for data learning, as demonstrated in Equation 1. Specifically, we employ MediaPipe, which not only infers facial feature coefficients but also estimates head pose.

$$S = S_0 + \sum_{i=1}^{n} w_i S_i, \qquad (1)$$

Where,
- $S$: The final shape of the face.
- $S_0$: The base shape (neutral expression).
- $S_l$: The blendshape for the $l$-th expression.
- $w_l$: The blendshape coefficient for the $l$-th expression, representing the influence of that expression on the final shape.
- $n$: The total number of blendshapes, here $n = 52$.

### B. Dataset Construction

Self-supervised learning through random facial expressions is an effective approach, as demonstrated by the Eva robot [12]. However, in practice, many expressions generated during this process fall within the uncanny valley. To address this issue, we developed an expert strategy system based on our robot's control framework to generate expressions, with the guiding formulation provided in Algorithm 1. By following these manually designed rules, the system inherently avoids unnatural expressions at the data distribution

level, ensuring more natural and realistic training data. The detail is illustrated in the left part of Fig. 3, which provides a comprehensive overview of the data construction phase.

---

**Algorithm 1** Facial Expression Expert Strategy Module

---

1: **Define** Servo_Group = [25 DoF]
2: **Define** Constraints, $\Phi$ = {7 constraint groups}
3: **Define** Emotion_samples, $S$= {}
4: **For** iteration $n = 1, 2, ...$ **do**
5:   **Randomly select** constraint $\phi_i$ from the available
6:   constraint groups $\Phi$
7:   **Apply constraints:**
8:     1. Horizontal eye servos (left-right) synchronized
9:     2. Vertical eye servos (up-down) synchronized
10:     3. Blink servos synchronized
11:     4. Eye opening servos synchronized
12:     5. Only one eyebrow movement (frown or raised)
13:       allowed
14:     6. Head movement servos randomly moved
15:     7. Mouth servos:
16:       a. Only one action (smile, sadness, or mouth-
17:         corners up) allowed
18:       b. Smile and sadness servos synchronized
19:       c. Mouth corners up independent
20:   **Update** servo positions per selected constraints
21:   **Append** current facial image $s_i$ in $S$
22: **End for**

---

### C. Model Design

Previous studies (e.g., Eva [12], [13], XIN-REN [28]) have directly mapped facial features, such as key points, to servo commands. This approach, however, overlooks the fact that different facial regions contribute variably to the driving weights of distinct servos, as illustrated in Fig. 4.

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i}}\right), \qquad (2)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i}}\right), \qquad (3)$$

$$Self\_Atten = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V, \qquad (4)$$

Where $i$ and $pos$ represent the blendshape's shape, and $Q$, $K$, and $V$ refer to the input facial features in the self-attention mechanism.

In this study, our model regresses servomotor commands by utilizing blendshape coefficients. Inspired by the variation in driving weights across facial regions, we designed an attention module that enables the model to adaptively select relevant blendshape coefficients. Specifically, we first apply positional encoding to the 52 blendshape coefficients to enhance their distinctiveness. Subsequently, the encoded feature coefficients are processed through a non-linear mapping that incorporates an attention mechanism. The positional encoding and attention formulas are provided in Equations 2, 3, and 4. This allows the network to learn the relative
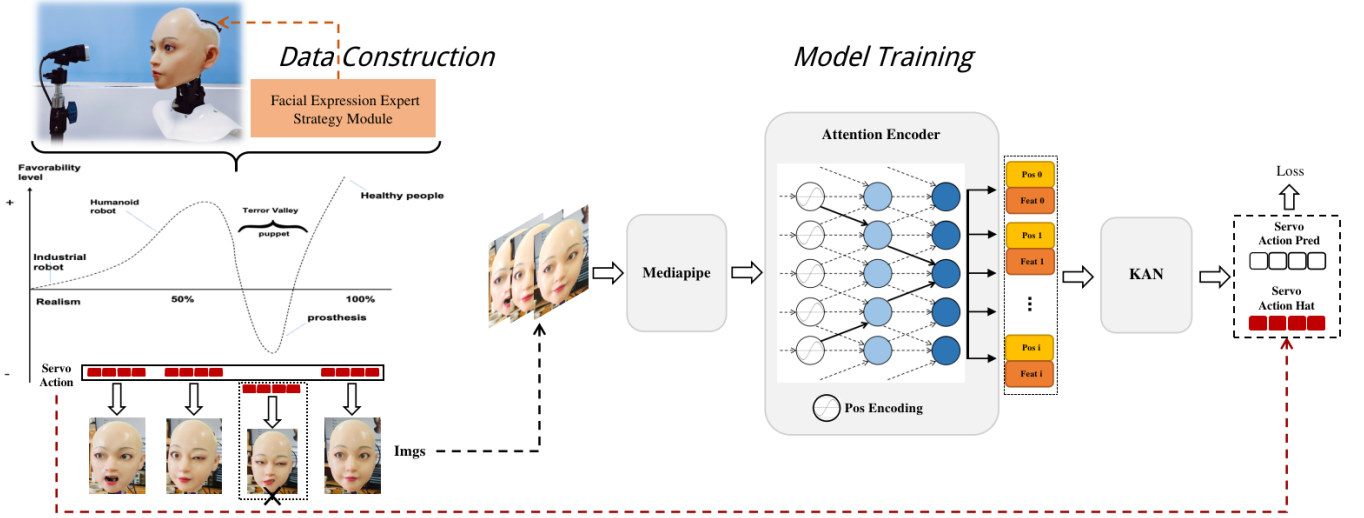
Fig. 3: Our method consists of two modules: dataset construction and network training. The dataset construction module automatically generates commands corresponding to random expressions using an expert strategy system, followed by capturing images with a conventional RGB camera. The network training module extracts facial representation blendshape coefficients using MediaPipe. Subsequently, our designed network performs regression to fit the servo commands.

importance of various facial features, which is illustrated in the overall workflow diagram shown in the right part of Fig. 3.
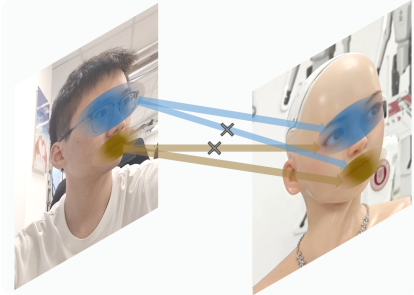


Fig. 4: A visual representation of how different facial regions contribute differently to the driving weights of different servos.

Considering the importance of response speed in robotic applications, this study employs a KAN instead of an MLP (Multilayer Perceptron) to learn the mapping regression from facial features to servomotor commands. KAN's learning parameters rely solely on the activation function, and as task complexity increases, the network can select higher-complexity spline activation functions to accommodate the task. Benefiting from KAN's ability to achieve similar fitting accuracy as MLP with fewer parameters, we integrated it into the overall framework to enhance inference speed. The final fitting formula of the KAN network is presented as $f(x)$, and define $f(x) = KAN(x)$:

$$f(\mathbf{x}) = \sum_{q=1}^{2n+1} \Phi_q \left( \sum_{p=1}^{n} \phi_{q,p}(x_p) \right), \qquad (5)$$

where $\phi_{q,p} : [0,1] \to \mathbb{R}$ and $\Phi_q : \mathbb{R} \to \mathbb{R}$. Specifically, we chose a basic cubic spline curve as the activation function and utilized a 3 layers' KAN network structure as follows:

$$KAN(x) = (\Phi_3 \circ \Phi_2 \circ \Phi_1)(x), \qquad (6)$$

$$\Phi_i(x) = w_b b(x) + w_s \text{spline}(x). \qquad (7)$$

We set $b(x) = \text{silu}(x) = \frac{x}{1+e^{-x}}$, based on the conclusion drawn from KAN, we parameterize the spline function $spline(x)$ as a linear combination of B-splines, such that $\text{spline}(x) = \sum_i c_i B_i(x)$. This representation allows for a flexible and accurate approximation of complex functions, leveraging the properties of B-splines to ensure smoothness and continuity across the domain of interest.

We also conducted a comparison with the method proposed by Eva [12]. Based on their approach, we designed an MLP model that utilizes key points regression to predict servomotor commands. However, directly extracting information from key points is not advisable due to the complexity and abstraction of the data [32]–[34]. The comparative experiments presented later further validate our analysis and demonstrate the effectiveness of our approach.

### D. Loss Design

Considering facial motion during expression generation, we propose an eye movement consistency loss based on our robot's structural control. This loss primarily ensures the consistency of the robot's eye motor movements.

$$L_{con} = |y_{\text{eyeleft}} - y_{\text{eyeright}}| + |y_{\text{browleft}} - y_{\text{browright}}|, \quad (8)$$

$$L_{MSE} = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2, \qquad (9)$$

Where $y_{\mathrm{eyeleft}}$ represents the command of the left eye, and the others have similar meanings, $\mathbf{y}$ and $\hat{\mathbf{y}}$ represent the servomotor commands predicted by the model and the actual servomotor commands.

The main objective function for this task is the MSE (Mean Squared Error) loss, which serves as a performance metric by evaluating the root mean square error between servomotor commands. The total loss function incorporates a hyperparameter $\lambda$, which, after extensive experimental validation, has been optimally set to 0.01. Further details on this experimental determination can be found in the experimental section.

$$L_{total} = L_{MSE} + \lambda L_{con}, \tag{10}$$

## V. EXPERIMENTS

### A. Dataset

Rena Facial Database: This dataset contains 9,000 sample images of the facial emotion robot, along with the corresponding servomotor commands that drive each expression. Of these, 8,000 images are used for training and 1,000 for testing.

MMI-Database: The MMI database consists of over 2,900 videos and high-resolution static images from 75 subjects, annotated with six basic emotions. We extracted 185 salient facial frames based on these six basic emotions for generalization testing.

Open Set Database: To further evaluate our model's generalization ability, we sampled four volunteers from the lab. Each volunteer participated in comparison experiments by mimicking facial expressions based on the six basic emotions.

The Rena Facial Database is a rigorously designed, comprehensive dataset that allows for detailed analysis of model training and qualitative error analysis, referred to as closed-set experiments. In contrast, the MMI Database and Open Set Database contain unlabeled facial expression samples used for qualitative testing of our model's performance. These samples are analyzed based on specific evaluation metrics and referred to as open-set experiments.

### B. Model Training

On the facial database, we compared three different models: an MLP framework based on facial landmarks, an MLP framework based on facial blendshape coefficients, and our proposed framework. The training loss and testing loss for these three models during the training process are shown in Fig. 6.
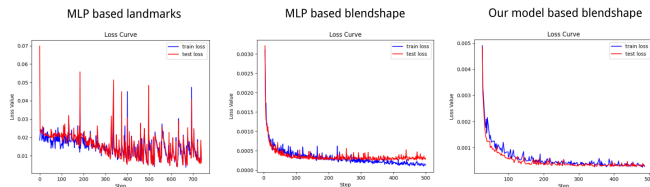


Fig. 6: Train loss and test loss during process for the three mentioned above.

The lowest test loss for the three models was 0.03, 0.0025, and 0.0020, respectively. Given that we have normalized the control range of each servo to [0,1], these results indicate average servo control error rates of 17%, 5%, and 4.4%. We also evaluated the overall error distribution of the models. To examine their performance across different error levels, we conducted full inference on the Facial Database using all three models. Subsequently, we assessed their performance using CED (Cumulative Error Distribution) curves, as shown in Fig. 7.
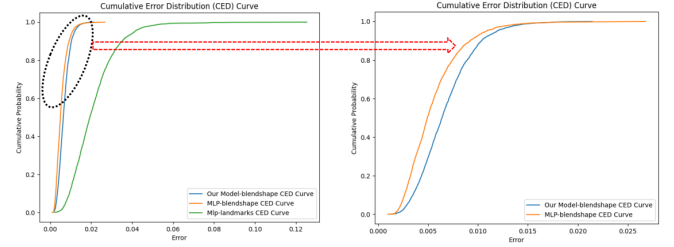


Fig. 7: CED curves for the three models mentioned above.

The test results of facial features based on blendshape representation are significantly better than those based on key points representation. Our model, which incorporates the designed KAN and attention mechanism, exhibits the best robustness in handling various facial expressions.

To evaluate the effectiveness of our designed model, we conducted ablation experiments to verify the utility of the attention mechanism. The evaluation metric for these experiments was the generalization error rate on the test set. As shown in Table I, the results confirm our hypothesis: incorporating the attention mechanism effectively enhances feature selection capability. Similarly, we designed a series of experiments to verify the efficacy of the proposed consistency loss. The results demonstrate that incorporating consistency loss optimizes model training. Additionally, the trend indicates that as the consistency loss increases from 0.001, the model's optimization performance begins to decrease.

TABLE I: Ablation Experiments on Hyperparameter $\lambda$ and Attention Mechanism.

| Backbone | Attention | $\lambda$ | Error(%) |
|---|---|---|---|
| MLP | ✓ | - | **0.0203** |
| MLP | ✗ | - | 0.0240 |
| Attention-KAN-bs | ✓ | 0.0 | 0.0275 |
| Attention-KAN-bs | ✓ | 0.1 | 0.0259 |
| Attention-KAN-bs | ✓ | 0.01 | **0.0234** |
| Attention-KAN-bs | ✓ | 0.001 | 0.0611 |

### C. Evaluation Metrics

We deployed our model on the facial emotion robot and visualized the robot's facial expressions alongside the input human facial images. Our evaluation was conducted by comparing 185 salient frames. As shown in the Fig. 5, we present a visual comparison of selected frames between the robot's expressions and the baseline input facial images.
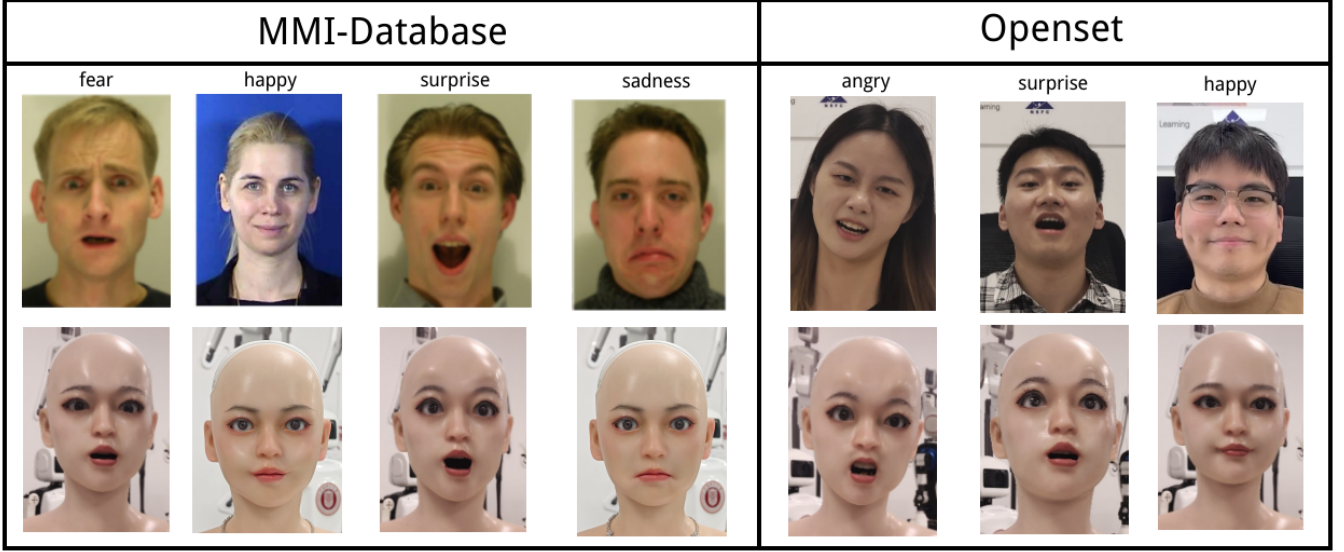
Fig. 5: We performed output servomotor commands on our Rena robot on both the MMI dataset and the open dataset to demonstrate that our method supports accurate simulation of various human expressions in multiple human subjects.

We utilized image distance(ID) and landmark distance(LD) as evaluation metrics, as shown in the following:

$$ID(x,y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)}, \quad (11)$$

$$LD = \frac{1}{n}\sum_{i=1}^{n}\sqrt{(\bar{x}_i - x_i)^2 + (\bar{y}_i - y_i)^2}, \quad (12)$$

In $ID$, $\mu_x$ is the mean value of image $x$. $\sigma_{xy}$ is the covariance of images $x$ and $y$. $C_1$ and $C_2$ are empirical constants, with values of 0.01 and 0.03 taken here. In $LD$, $x_i$ and $y_i$ are the coordinates of facial key points.

The former represents pixel-level accuracy, while the latter measures the similarity between facial key points. Our model outperforms the random baseline (RS) and Eva Method by a large margin, as shown in Table II. Note that the image distance is normalized by the total number of pixel values (480, 640, 3), which range from 0 to 1, where the landmark distance is normalized by the 63 landmarks.

TABLE II: Accuracy of Different Models

| Method | ID | LD |
|---|---|---|
| RS Method | 6.47 | 0.84 |
| Eva Method | 3.47 | 0.46 |
| Our Method | **2.96** | **0.074** |

Sequential indicators play a crucial role in the subjective experience of humans, as discussed in [35]. We further evaluate the performance of an online expression learned from a human performer. This performer executes a variety of facial expressions with neutral-peak-neutral variations. The expression shape vectors of the performer are captured at a rate of 30 frames per second and are used as the robot's target expression vectors.

The values of average space similarity and time similarity reflect the similarity of the imitation trajectory with the performer's facial actions, whereas the value of the average movement smoothness reflects the smoothness of continuous servo motions. The results in Table III show the progressiveness of our method compared with three baseline methods using the state-of-the-art humanoid expression generation systems (Jaekel method [36], Trovato method [37], and Habib method [38]).

TABLE III: Comparison of Sequential Indicators Versus Four Methods

| Sequential indicators | Space similarity | Time similarity | Movement smoothness |
|---|---|---|---|
| Jaekel Method | 85.4 | 85.2 | 83.3 |
| Trovato Method | 84.4 | 83.9 | 87.6 |
| Habib Method | 87.8 | 86.1 | 84.3 |
| Our Method | **91.8** | **88.1** | **92.7** |

The sequential indicators of space-similarity $G_s$, time-similarity $G_t$, and movement smoothness $G_d$, measured by servo hopping during $t_i$ to $t_L$, following the method proposed by Zhu et al. [39], are defined as follows:

$$G_s = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{1}{L}\sum_{k=1}^{L}F(d_i^H(t_k) - d_i^R(t_k), b_S)\right), \quad (13)$$

$$G_d = 1 - \frac{1}{L}\sum_{k=1}^{L}\frac{1}{m}\sum_{j=1}^{m}G(c_j(t_k)), \quad (14)$$

$$G_t = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{1}{L} \sum_{k=1}^{L} F \left( d_i^H(t_k) - d_i^H(t_{k-1}) \right. \right.$$
$$\left. \left. - (d_i^R(t_k) - d_i^R(t_{k-1})), b_T \right) \right), \qquad (15)$$

Where $L = T = 30$ is the frame rate of camera; $(x_{i0}^H, y_{i0}^H)$ and $(x_{i0}^R, y_{i0}^R)$ represent the $i$th feature point positions of the human and robot at the initial moment, respectively. $d_i^H(t) = \sqrt{(x_i^H(t) - x_{i0}^H)^2 + (y_i^H(t) - y_{i0}^H)^2}$ and $d_i^R(t) = \sqrt{(x_i^R(t) - x_{i0}^R)^2 + (y_i^R(t) - y_{i0}^R)^2}$ are the $i$th feature point displacements of the human and robot at moment $t$, respectively; $F(x, b) = e^{-x^2/b}$ is a fitting function that converts the deviation parameter $x$ to $(0, 1)$ similarity, and $b$ is the parameter used to control the mapping performance. $G(c_j(t_k))$ indicates whether the displacement of the $j$th servo exists unsmoothed hopping at moment $t_k$, and is measured as follows:

$$G(c_j(t_k)) =$$
$$\begin{cases} 1, & |c_j(t_k) - c_j(t_{k-1})| - |c_j(t_{k-1}) - c_j(t_{k-2})| > T_D, \\ 0, & \text{otherwise.} \end{cases}$$
$$(16)$$

### D. Implementation Details

The model was trained using a Titan RTX GPU, with an Intel-F200 camera for data capture. The batch size was set to 256, and the learning rate was 0.00001. The optimizer used for training was ADAM.

## VI. CONCLUSIONS

We propose a facial emotion robot featuring highly biomimetic facial units and a high-degree-of-freedom facial structure. Additionally, we introduce a lightweight end-to-end learning framework specifically designed for facial mimicry. Our experiments demonstrate that our approach can accurately and efficiently recognize facial features and drive the expression robot to replicate similar expressions. Imitation is a crucial step toward endowing robots with more complex skills and serves as a foundational meta-skill for interacting with the external world.

While we show that our robot can successfully imitate various human facial expressions through visual observation, it faces limitations in handling abrupt expression transitions during continuous facial motion. The robot focuses only on the current frame, making it less effective in responding to sudden changes in expressions. Thus, future work should focus on enhancing the model's ability to generate sequential expressions by refining the dataset and the model architecture.

## REFERENCES

[1] Robert Plutchik. Emotions: A general psychoevoiutionary theory. In *Approaches to emotion*, pages 197–219. Psychology Press, 2014.

[2] Nadja Reissland. Neonatal imitation in the first hour of life: Observations in rural nepal. *Developmental Psychology*, 24(4):464, 1988.

[3] Andrew N Meltzoff and M Keith Moore. Imitation in newborn infants: Exploring the range of gestures imitated and the underlying mechanisms. *Developmental psychology*, 25(6):954, 1989.

[4] Rick Van Baaren, Loes Janssen, Tanya L Chartrand, and Ap Dijksterhuis. Where is the love? the social aspects of mimicry. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 364(1528):2381–2389, 2009.

[5] Jun-Ho Oh, David Hanson, Won-Sup Kim, Young Han, Jung-Yup Kim, and Ill-Woo Park. Design of android type humanoid robot albert hubo. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1428–1433. IEEE, 2006.

[6] Takuya Hashimoto, Sachio Hitramatsu, Toshiaki Tsuji, and Hiroshi Kobayashi. Development of the face robot saya for rich facial expressions. In *2006 SICE-ICASE International Joint Conference*, pages 5423–5428. IEEE, 2006.

[7] Takuya Hashimoto, Sachio Hiramatsu, and Hiroshi Kobayashi. Dynamic display of facial expressions on the face robot made by using a life mask. In *Humanoids 2008-8th IEEE-RAS International Conference on Humanoid Robots*, pages 521–526. IEEE, 2008.

[8] Hisashi Ishihara, Rina Hayashi, Francois Lavieille, Kaito Okamoto, Takahiro Okuyama, and Koichi Osuka. Automatic generation of dynamic arousal expression based on decaying wave synthesis for robot faces. *Journal of Robotics and Mechatronics*, 36(6):1481–1494, 2024.

[9] David Loza, Samuel Marcos Pablos, Eduardo Zalama Casanova, Jaime Gómez García-Bermejo, and José Luis González. Application of the facs in the design and construction of a mechatronic head with realistic appearance. 2013.

[10] Chyi-Yeu Lin, Chun-Chia Huang, and Li-Chieh Cheng. An expressional simplified mechanism in anthropomorphic face robot design. *Robotica*, 34(3):652–670, 2016.

[11] Wagshum Techane Asheber, Chyi-Yeu Lin, and Shih Hsiang Yen. Humanoid head face mechanism with expandable facial expressions. *International Journal of Advanced Robotic Systems*, 13(1):29, 2016.

[12] Boyuan Chen, Yuhang Hu, Lianfeng Li, Sara Cummings, and Hod Lipson. Smile like you mean it: Driving animatronic robotic face with learned models. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2739–2746. IEEE, 2021.

[13] Yuhang Hu, Boyuan Chen, Jiong Lin, Yunzhe Wang, Yingke Wang, Cameron Mehlman, and Hod Lipson. Human-robot facial coexpression. *Science Robotics*, 9(88):eadi4724, 2024.

[14] Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljačić, Thomas Y Hou, and Max Tegmark. Kan: Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756*, 2024.

[15] A Waswani, N Shazeer, N Parmar, J Uszkoreit, L Jones, A Gomez, L Kaiser, and I Polosukhin. Attention is all you need. In *NIPS*, 2017.

[16] Boren Li, Hang Li, and Hangxin Liu. Driving animatronic robot facial expression from speech. *arXiv preprint arXiv:2403.12670*, 2024.

[17] Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu Zhu, Jun He, Hongyan Liu, and Zhaoxin Fan. Emotalk: Speech-driven emotional disentanglement for 3d face animation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20687–20697, 2023.

[18] Wei Zhao, Yijun Wang, Tianyu He, Lianying Yin, Jianxin Lin, and Xin Jin. Breathing life into faces: Speech-driven 3d facial animation with natural head pose and detailed shape. *arXiv preprint arXiv:2310.20240*, 2023.

[19] Tao Liu, Chenpeng Du, Shuai Fan, Feilong Chen, and Kai Yu. Diffdub: Person-generic visual dubbing using inpainting renderer with diffusion auto-encoder. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3630–3634. IEEE, 2024.

[20] Yifeng Ma, Shiwei Zhang, Jiayu Wang, Xiang Wang, Yingya Zhang, and Zhidong Deng. Dreamtalk: When expressive talking head generation meets diffusion probabilistic models. *arXiv preprint arXiv:2312.09767*, 2023.

[21] Chenxu Zhang, Chao Wang, Jianfeng Zhang, Hongyi Xu, Guoxian Song, You Xie, Linjie Luo, Yapeng Tian, Xiaohu Guo, and Jiashi Feng. Dream-talk: diffusion-based realistic emotional audio-driven method for single image talking face generation. *arXiv preprint arXiv:2312.13578*, 2023.

[22] Chao Liang, Qinghua Wang, Yunlin Chen, and Minjie Tang. Wav2liphr: Synthesising clear high-resolution talking head in the wild. *Computer Animation and Virtual Worlds*, 35(1):e2226, 2024.

[23] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024.

[24] Li-Heng Lin, Yuchen Cui, Amber Xie, Tianyu Hua, and Dorsa Sadigh. Flowretrieval: Flow-guided data retrieval for few-shot imitation learning. *arXiv preprint arXiv:2408.16944*, 2024.

[25] Joey Hejna, Chethan Bhateja, Yichen Jian, Karl Pertsch, and Dorsa Sadigh. Re-mix: Optimizing data mixtures for large scale imitation learning. *arXiv preprint arXiv:2408.14037*, 2024.

[26] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.

[27] Boyuan Zheng, Sunny Verma, Jianlong Zhou, Ivor W Tsang, and Fang Chen. Imitation learning: Progress, taxonomies and challenges. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[28] Fuji Ren and Zhong Huang. Automatic facial expression learning method based on humanoid robot xin-ren. *IEEE Transactions on Human-Machine Systems*, 46(6):810–821, 2016.

[29] Erika Chuang and Chris Bregler. Performance driven facial animation using blendshape interpolation. *Computer Science Technical Report, Stanford University*, 2(2):3, 2002.

[30] Rachel McDonnell, Katja Zibrek, Emma Carrigan, and Rozenn Dahyot. Model for predicting perception of facial action unit activation using virtual humans. *Computers & Graphics*, 100:81–92, 2021.

[31] John P Lewis, Ken Anjyo, Taehyun Rhee, Mengjie Zhang, Frederic H Pighin, and Zhigang Deng. Practice and theory of blendshape facial models. *Eurographics (State of the Art Reports)*, 1(8):2, 2014.

[32] Rahib H Abiyev. Facial feature extraction techniques for face recognition. *Journal of Computer Science*, 10(12):2360, 2014.

[33] Lance Williams. Performance-driven facial animation. In *Acm SIGGRAPH 2006 Courses*, pages 16–es. 2006.

[34] Inkyu Park and Jaewoong Cho. Said: Speech-driven blendshape facial animation with diffusion. *arXiv preprint arXiv:2401.08655*, 2023.

[35] Hui Yu and Honghai Liu. Regression-based facial expression optimization. *IEEE transactions on human-machine systems*, 44(3):386–394, 2014.

[36] Peter Jaeckel, Neill Campbell, and Chris Melhuish. Facial behaviour mapping—from video footage to a robot head. *Robotics and Autonomous Systems*, 56(12):1042–1049, 2008.

[37] Gabriele Trovato, Massimiliano Zecca, Tatsuhiro Kishi, Nobutsuna Endo, Kenji Hashimoto, and Atsuo Takanishi. Generation of humanoid robot's facial expressions for context-aware communication. *International Journal of Humanoid Robotics*, 10(01):1350013, 2013.

[38] Christian Becker-Asano and Hiroshi Ishiguro. Evaluating facial displays of emotion for the android robot geminoid f. In *2011 IEEE workshop on affective computational intelligence (WACI)*, pages 1–8. IEEE, 2011.

[39] Niyati Rawal, Dorothea Koert, Cigdem Turan, Kristian Kersting, Jan Peters, and Ruth Stock-Homburg. Exgennet: Learning to generate robotic facial expression using facial expression recognition. *Frontiers in Robotics and AI*, 8:730317, 2022.