

Exploring Timeline Control for Facial Motion Generation

Yifeng Ma¹, Jinwei Qi², Chaonan Ji², Peng Zhang², Bang Zhang², Zhidong Deng¹, Liefeng Bo²

¹ Department of Computer Science and Technology, Tsinghua University

² Tongyi Lab, Alibaba Group

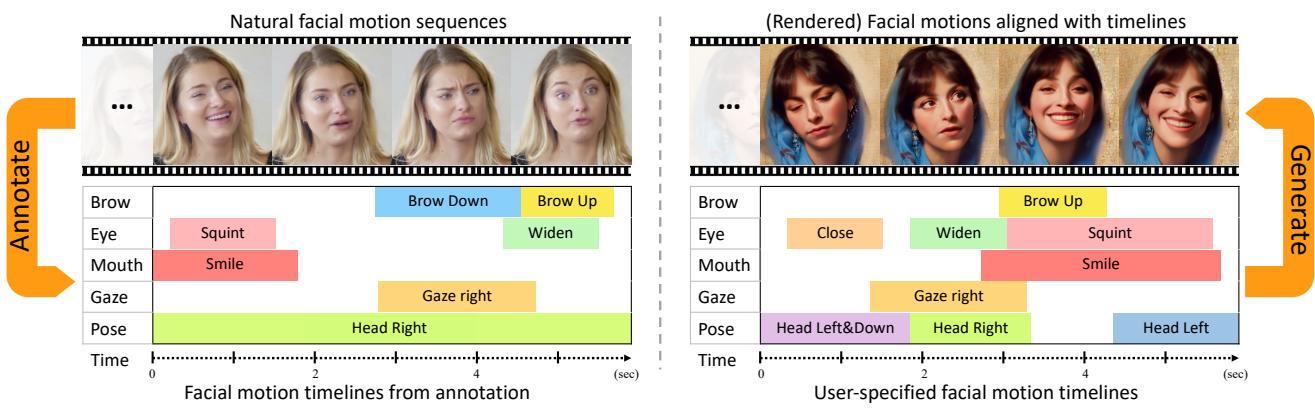


Figure 1. We introduce a new control signal for facial motion generation: **timeline control**. We first utilize a labor-efficient approach to annotate the time intervals of facial motion at a frame-level granularity. Using the annotations, we propose a model that can generate natural facial motions aligned with an input timeline. Compared to previous controls like audio and text, timeline control enables precise temporal control of facial motions. In this paper, facial motions are rendered into photorealistic videos for better visualization.

Abstract

This paper introduces a new control signal for facial motion generation: **timeline control**. Compared to audio and text signals, timelines provide more fine-grained control, such as generating specific facial motions with precise timing. Users can specify a multi-track timeline of facial actions arranged in temporal intervals, allowing precise control over the timing of each action. To model the timeline control capability, we first annotate the time intervals of facial actions in natural facial motion sequences at a frame-level granularity. This process is facilitated by Toeplitz Inverse Covariance-based Clustering to minimize human labor. Based on the annotations, we propose a diffusion-based generation model capable of generating facial motions that are natural and accurately aligned with input timelines. Our method supports text-guided motion generation by using ChatGPT to convert text into timelines. Experimental results show that our method can annotate facial action intervals with satisfactory accuracy, and produces natural facial motions accurately aligned with timelines.

1. Introduction

Generating vivid facial motions has drawn growing attention due to its broad applications, including digital human generation and filmmaking. To produce desired facial motions, current methods use audio or text to provide guidance. However, these methods lack a critical ability that users often require: generating specific facial motions with precise timing. For example, users may want to generate a brow raise between frames 10 and 30 while simultaneously generating a smile between frames 14 and 43. Audio-driven methods [43, 54] can only generate motions synchronized with the audio. Text signals offer only coarse-grained descriptions of facial motions and lack frame-level detail. Text-driven approaches [4, 44, 49, 51] rely on temporal adverbial cues (e.g., *then*) for rough timing guidance.

To achieve more fine-grained control, we introduce a novel control signal: timeline control for facial motion generation. In this setup, users can generate facial motions by inputting an intuitive timeline that contains several temporal intervals, each corresponding to a desired facial action. This setup enables users to manage the timing of each action.

Achieving timeline control is very challenging. Frame-

level control of facial motions requires the model to achieve exceptional precision in generating accurate actions. To model timeline control capability, it is crucial to annotate the precise start and end frames of facial actions, a challenge that existing methods have yet to overcome. Existing methods [44, 51] represent facial motion sequences as a time series of motion descriptors, such as Facial Action Units (AUs) [10] or blendshapes [46], and then leverage ChatGPT to summarize these time series for annotation. However, ChatGPT’s limited sensitivity to the temporal dynamics of facial motions prevents it from determining the exact start and end frames of facial actions. Another annotation approach is using thresholds, like labeling the eye as “closed” if the blendshape *eyeBlink* exceeds 0.4. However, setting thresholds for complex actions, such as brow motions, is challenging. Additionally, some actions are determined by the relative values of multiple motion descriptors rather than a single one.

To address this issue, we adopt Toeplitz Inverse Covariance-based Clustering (TICC) [18] to annotate the start and end frames of facial actions. TICC can segment the time series into a sequence of intervals, with each interval containing a single motion pattern. Each interval has clearly defined start and end frames. Once the start and end frames of the intervals are identified, the remaining task is to determine the facial action represented by each interval. This can be addressed by another feature of TICC: its ability to group the segmented intervals into several clusters based on the similarity of their motion patterns, with each cluster containing intervals that exhibit similar motion patterns. The automatic clustering process eliminates the need for manual threshold setting and considers the relationships between multiple motion descriptors. Once these clusters are obtained, we can determine the overall action (*e.g.* significantly raised brows) represented by each cluster by analyzing a few representative patterns within it. This process requires minimal labor. Using this procedure, we analyze motions in different facial regions individually and achieve frame-level annotation for the brows, eyes, mouth, gaze, and head motions in a labor-efficient manner.

Using frame-level facial motion annotations, we propose a diffusion-based generation model that produces natural facial motions accurately aligned with the input timeline. The model generates motions for each facial region (*e.g.* upper face, lower face) separately. This decoupling can help address the adverse impact of motion coupling in the data, thereby improving accuracy. However, some couplings are necessary for motion naturalness, such as a slight head lift during a brow raise. To properly manage coupling, the generation model is divided into a base network and multiple branch networks, with each branch network dedicated to a facial region. The base network encodes global facial motion couplings into base features. The branch network takes

base features and the timeline of the corresponding region as input to accurately generate facial motions for that region while maintaining natural couplings. We use FaceVerse[46] 3DMM coefficients to represent facial motions for generation. Since photorealistic facial motion enables more accurate assessments of motion realism than mesh-based motion [31], we render motions into photorealistic videos by a diffusion-based renderer.

Our method enables text-guided facial motion generation by using ChatGPT to convert text into timelines. ChatGPT effectively translates simple natural language descriptions into timelines, which our model transforms into facial motions. Users can modify these timelines, allowing precise timing control for text-driven applications.

In summary, our contributions are as follows:

- We are the first to develop a labor-efficient approach to annotate temporal intervals of facial actions. Such fine-grained annotations enable precise modeling of the temporal dynamics of facial motions.
- We are the first to achieve timeline control for facial motion generation. This enables users to generate specific facial motions with precise timing. Our method also supports generating actions from natural language text.
- Detailed evaluations show that our method generates accurate motion annotations and produces natural, precise facial motions from timeline inputs.

2. Related Work

Facial Motion Annotation. Compared to human motion annotation [15, 16, 29, 36], facial motion annotation is in its early stages. Existing methods [28, 44, 49, 51, 60] generate annotations that either overlook temporal changes or describe them only roughly with temporal adverbs. Some datasets [38, 56] label intervals of macro/micro expressions rather than facial actions used for facial motion generation. These datasets also require intensive manual labeling. Facial expression spotting methods [9, 17, 48, 57, 59] are also restricted to label macro/micro-expressions. In this paper, we aim to labor-efficiently generate frame-level facial action annotations.

Facial Motion Generation. Existing methods use audio or text as the control signal and cannot generate specific motions at precise timing. Audio-driven methods (2D [6, 14, 19, 23, 26, 32, 37, 42, 43, 47, 50, 53, 54, 58, 63, 64] and 3D methods [1, 7, 8, 11, 12, 24, 25, 34, 41, 52, 55, 65]) can only generate facial motions synchronized with the audio. Text-driven methods [4, 27, 28, 44, 49, 51] can only use temporal adverbial phrases to control timing roughly. Rule-based methods [5, 33] offer temporal control but often produce unnatural motions due to deviations from real movement distributions. This paper introduces timeline control to generate natural facial motions with precise timing.

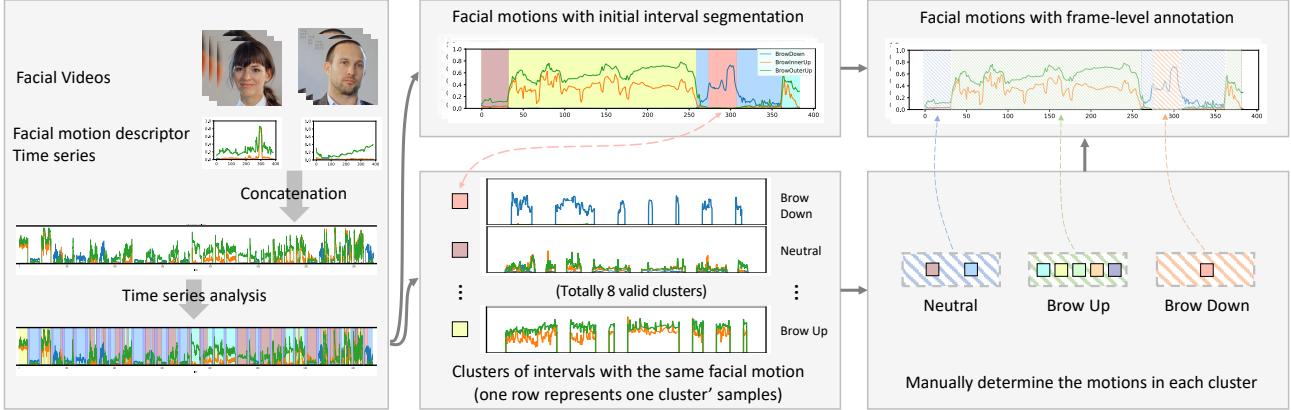


Figure 2. **The pipeline of frame-level facial motion annotation** (using brow motions as an example). We first extract facial motion descriptors (blendshapes) from natural facial motion videos and concatenate the results to create a facial motion time series for time series analysis. This analysis can simultaneously segment the sequence into a series of motion patterns and cluster similar patterns, resulting in multiple clusters, each containing consistent facial motion patterns. Then, by inspecting a few patterns, we identify the facial motions each cluster represents, thereby obtaining frame-level facial motion annotations for all videos.

Facial Region	Motion Categories
Brow	BrowUp, BrowDown, Neutral
Eye	EyeSquint, EyeWiden, EyeClose, Neutral
Mouth	SoftSmile, Smile, MouthFrown, Neutral
Gaze	Left, Right, Up, Down, Neutral
Head	Left, Right, Up, Down, Neutral

Table 1. Facial actions annotated by our method.

Timeline control has been utilized in human motion generation [2, 3, 35, 39] by first generating motion segments and then piecing them together. However, due to the rapid changes in facial movements, this approach is unsuitable, necessitating a new model structure for effective generation.

3. Method

In this work, we first annotate the temporal intervals of facial action for natural facial motion sequences. Based on these annotations, we develop a generation model that generates natural facial motions that are accurately aligned with an input timeline. Our method also supports text-guided motion generation by leveraging ChatGPT to convert text into timelines.

3.1. Annotating temporal intervals of facial actions

Problem Formulation. For each video V in dataset, our method aims to generate frame-level facial motion annotations $\mathbf{A} = [\mathbf{a}_i]_{i=1}^L$ for each video frame $[\mathbf{v}_i]_{i=1}^L$. Tab. 1 shows the annotated facial actions. The annotation \mathbf{a}_i for each frame is a vector of binary values, with each dimension representing a non-neutral facial action. A dimension is set to 1 if the action occurs, otherwise, it is 0.

Annotation Process Overview. We achieve the annotation through time series analysis TICC. Fig. 2 shows the pipeline: (1) We extract proper motion descriptors for each

facial region to obtain facial motion time series. (2) We apply TICC to simultaneously segment the time series into a sequence of motion patterns and cluster similar patterns. This results in several clusters of intervals, each containing similar facial motions. (3) Finally, for each cluster, we can infer the actions of all intervals by examining only a few intervals. Once the action for each cluster is identified, the actions for all intervals are determined, resulting in facial action annotations at a frame-level granularity.

Facial Motion Descriptor. We observe that previously used AUs lack sufficient precision, while ARKit blendshapes offer a more accurate representation. The publicly available MediaPipe blendshape detector struggles to detect certain blendshapes (*e.g.* eyeWide, mouthFrown), so we use an in-house blendshape detector. ARKit blendshapes are a set of coefficients that describe facial expressions, with each coefficient representing the movement of a specific facial region. Each coefficient ranges from 0 to 1, indicating the motion intensity.

We use the blendshapes corresponding to each facial region to construct motion time series for each region, resulting in time series data for the eye, brow, and mouth areas. For the eye region, we selected *eyeBlink*, *eyeSquint*, *eyeWide*, and for the brow region, we selected *browDown*, *browInnerUp*, *browOuterUp*. For the mouth region, we use only expression-related blendshapes (*mouthSmile*, *mouthStretch*, *mouthFrown*) to avoid interference from speech-related movements. ARKit blendshapes provide coefficients for the left and right sides. As facial motions in our dataset are largely symmetrical, we use only the left-side coefficients to simplify analysis.

For gaze and head pose motion, we use a 3DMM model FaceVerse [46] to detect coefficients. For gaze, we select eye coefficients, for head pose, we select angle coefficients.

Analyzing Facial Motion Time Series. We utilize

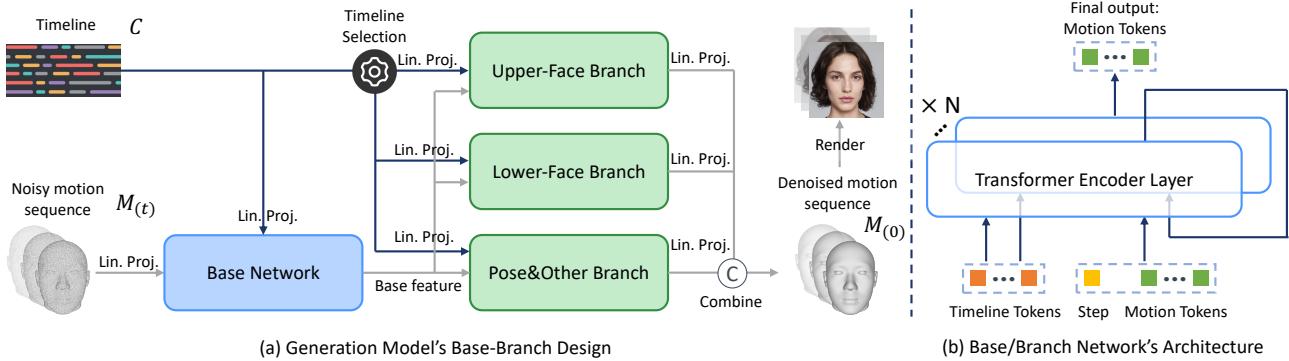


Figure 3. **Illustration of generation model.** (a) **Base-Branch Design.** The base network takes the timelines of all facial regions as input and outputs base features that model the global facial motion couplings. Through timeline selection, each region’s timeline is directed to its respective branch network. Since head pose is interconnected with all facial movements, the pose branch receives timelines of all regions. Each branch network takes the timeline of its corresponding region to generate the facial motions for that region. These motions are then combined to produce the overall motion of the entire face. *Lin. Proj.* denotes *Linear Projection*. (b) **Base/Branch Network’s Architecture.** Timeline control guides motion generation through cross-attention. The initial timeline tokens remain unchanged and are added at each layer. For clarity, the diffusion step (omitted in sub-figure (a) for clarity) is applied to each base and branch network.

TICC [18] for time series analysis. TICC is designed to handle a single long sequence (CubeMarker [22] can analyze multiple time series analysis but performs poorly.). However, our data consists of multiple short sequences. To address this issue, we concatenate all the short sequences into a single long sequence for analysis. To avoid the impact of abrupt transitions between different videos, we introduce a “null sequence” between videos to separate them. The null sequence is implemented as a length-100 sequence with values set to -1. These null sequences will eventually be grouped into a single cluster, which will not affect others.

After the time series analysis, the video is annotated with a sequence of motion pattern intervals and similar patterns are clustered into a few clusters. Therefore, for each cluster, we inspect a few intervals to determine the facial action they represent. This facial action is then considered to be the action for all intervals within that cluster. Different clusters may represent variations of the same facial action category, such as a high eyebrow raise versus a moderate one. We categorized the clusters into facial motion categories, such as eyebrow raise, frown, and neutral, to obtain frame-level annotations.

We observe that eye closure frequently produces sharp spikes in the time series, which disrupts the analysis of squinting and widening. Therefore, we analyze eye closure and squinting/widening separately. We also find that threshold-based annotation, which classifies actions using predefined thresholds (e.g., labeling the eye as “closed” if *eyeBlink* exceeds 0.4), achieves high accuracy for eye closure, gaze, and head pose. Consequently, we use this method for these regions, while applying time series analysis for others.

Detect Facial Motion for Unseen Videos. After learning various facial motion patterns from the data, TICC can also

be applied to previously unseen facial motion sequences and detect facial action intervals. We use this function to assess the accuracy of facial motion generation in our evaluation.

3.2. Facial Motion Generation from Timelines

Problem Formulation. Given a timeline control $C = [c_i]_{i=1}^L$ for each video frame, the generation model \mathcal{G} aims to generate natural facial motions $M = [m_l]_{l=1}^L$ that align with the timeline control. The control for each frame c_i follows the same format as the annotation. Generated facial motions are rendered by a diffusion-based renderer into photorealistic videos for better visualization.

Discussion: The coupling of facial motions and its impact on motion accuracy and naturalness. Facial motions across different regions are coupled. For example, when people smile, they might also squint and lower their eyebrows. These couplings can reduce the precision of generated motions. For instance, when generating a smile, the model may also learn to produce eyebrow-lowering actions, which could conflict with user-specified eyebrow motions (e.g. brow raise). To enhance accuracy, it is essential to decouple the generation of different facial regions. However, certain facial motion couplings are crucial for conveying naturalness—for instance, a subtle head movement that accompanies a shift in gaze. Therefore, the model must achieve a delicate balance, selectively decoupling motions to improve accuracy while preserving necessary couplings to enhance naturalness.

Generation Model \mathcal{G} . \mathcal{G} is a diffusion-based model. Diffusion models consist of two Markov chains [21, 40]: the forward chain incrementally injects Gaussian noise into the original signal, while the reverse chain sequentially reconstructs the original signal from the noise. \mathcal{G} predicts the original signal instead of noise, and the loss function is in-

troduced as:

$$\mathcal{L}_{\text{denoise}} = \mathbb{E}_{t \sim \mathcal{U}[1, T], M_{(0)}, C} (\|M_{(0)} - \mathcal{G}(M_{(t)}, t, C)\|^2), \quad (1)$$

where t denotes the diffusion step, $M_{(0)}$ is the original facial motion sequence, and $M_{(t)}$ is the noisy sequence produced by the diffusion forward function $q(M_{(t)} | M_{(t-1)}) = \mathcal{N}(\sqrt{\alpha_n} M_{(t-1)}, (1 - \alpha_n) I)$. C is the timeline condition.

To balance decoupling and coupling for both accuracy and naturalness, the generation model employs a base-branch design (Fig. 3). It consists of a base network and individual branch networks for each facial region (*e.g.* upper face, lower face). The base network takes timelines of all facial regions and noisy motions as input and encodes global motion couplings into base features. The branch network takes the base features the timelines relevant only to its designated facial region as inputs, and generates decoupled motions for each region.

We use the the expression, eye, and pose coefficients of FaceVerse [46] 3DMM model as target motion representation for generation. Its expression coefficients align with ARKit blendshapes, allowing separate representations of different facial regions. When dividing facial motions into different regions for separate generations, we observe that splitting them into the upper face, lower face, and pose & other regions strikes a better balance between decoupling and necessary coupling, rather than assigning a branch network for each individual region (*e.g.* further split eye and brow). Specifically, the upper face includes eye, brow, and gaze; the lower face includes mouth and jaw; and the pose & other regions cover the head pose and remaining regions like cheek and nose. It is important to note that, since pose is coupled with facial motions in all regions, the timeline control for the pose & other branch includes all facial regions (rather than just the control for pose).

As for the specific implementation, the timeline control for each frame c_i is implemented as a 16-dimensional vector composed of values 0 and 1, with each dimension representing an action. A value of 1 indicates that the action is performed in that frame, and a value of 0 indicates that the action is not present. The timeline control is transformed into timeline tokens through a linear projection. The base network and branch networks share the same structure. Their inputs consist of timeline tokens and motion tokens. The motion tokens are derived either from the noisy motion through linear projection or from base features. The motion tokens learn the temporal control from the timeline tokens through cross-attention in multiple transformer encoder layers. When incorporated into each encoder layer, the timeline tokens always use the initial timeline tokens rather than the output tokens from the previous layer. This prevents the temporal information in timeline tokens from being altered, thereby enhancing motion accuracy. The motion tokens for each layer come from the output of the motion tokens from

the previous layer. The network’s final output consists of the motion tokens from the last layer, which are then linearly projected to produce the facial motion. Each branch network generates the facial motion for its corresponding region. These motions are then combined to produce the motion for the entire face.

To improve flexibility and generalization, we use classifier-free guidance [20] to train our model. During training, we randomly drop the condition of each facial region timeline. we use a dropout probability of 0.5 for each condition independently, with a 0.1 probability of dropping all conditions and a 0.1 probability of maintaining all conditions. When the timeline condition for a certain region is dropped, its value is set to -1.

Rendering. We develop a diffusion-based portrait animation method to render motions as images. Since the renderer is not the main contribution of this paper, its details are reported in the *Supplementary Material*.

3.3. Translating Natural Language to Timelines

We observe that ChatGPT can generate a reasonable timeline from natural language descriptions, enabling our model to produce facial motions based on natural language input. The generated timeline can be further edited by users to fit their needs. To enhance the realism of the generated movements, we manually annotated several timeline descriptions as examples, allowing ChatGPT to perform few-shot learning. The details of the prompts and generated results are reported in *Supplementary Material*.

4. Experiments

Dataset. To effectively model natural facial motions, a dataset capturing authentic facial expressions is essential. For this purpose, we utilize the RealTalk dataset [13]. The RealTalk dataset contains 692 videos captured from genuine, unscripted conversations, resulting in highly natural facial movements with rich temporal dynamics. Previously used lab-recorded, scripted datasets, such as MEAD [45], lack these qualities. We use 1412 video clips extracted from Realtalk as the dataset. The size of the dataset is approximately 600,000 frames.

4.1. Facial Motion Annotation

Baselines. No previous methods have attempted annotating intervals of facial actions in a labor-efficient manner. We experimented with time series analysis methods such as AutoPlait [30], CubeMarker [22], and TICC. However, AutoPlait and CubeMarker failed to produce satisfactory results, as the motions of the segments within each cluster were inconsistent, making annotation impossible. Therefore, we primarily relied on TICC. Empirical results [18] indicate TICC’s annotation performance can be enhanced by tuning key hyperparameters, including the number of clusters and

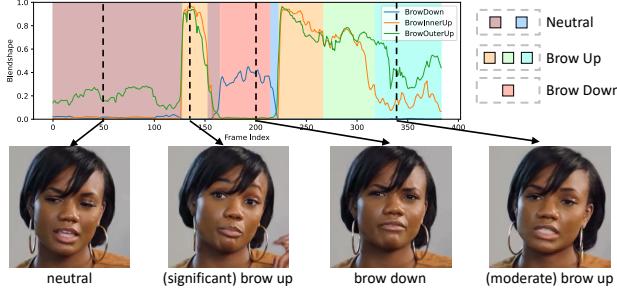


Figure 4. An example of brow motion annotation.

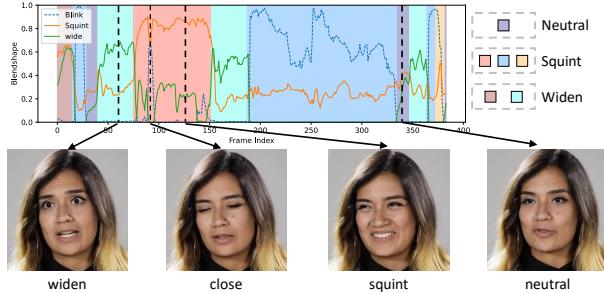


Figure 5. An example of eye motion annotation.

beta, a hyperparameter that facilitates temporal consistency. We examine the effects of these parameters.

Evaluation Metric. Annotating intervals of facial actions is essentially a multi-class classification problem, where each data point represents a frame and each class represents an action. To evaluate annotation accuracy, we manually annotate the facial actions on 50 videos to serve as ground truth, and use Macro-F1 as the metric for evaluation. For each class, the F1 score is the harmonic mean of the precision and recall of our estimate. Then, the macro-F1 score is the average of the F1 scores for all the classes. A higher Macro-F1 score indicates greater annotation accuracy. We calculated the macro-F1 score for the eye region only on segments when the eyes are not closed.

Qualitative results. Fig. 4 shows brow motion annotation. Our method can precisely segment different facial motion patterns and cluster similar ones. For a specific facial motion brow raise, our approach can distinguish multiple variants, such as a significant or moderate raise. In this work, these finer categories are grouped under one motion category for simplicity, but future methods could generate each subcategory separately. Annotation for the mouth is similar to that for the brow. Fig. 5 shows eye motion annotation. As stated in Sec. 3.1, only eye squint/widen are annotated using TICC. Eye closure is determined by a threshold on the blendshape *eyeBlink* and may occur within intervals of other eye motions. We also conduct a user study to evaluate annotation accuracy and the results are reported in Sec. 4.2.

Quantitative results. We select the optimal hyperparameters based on the highest macro-F1 score. Experiments

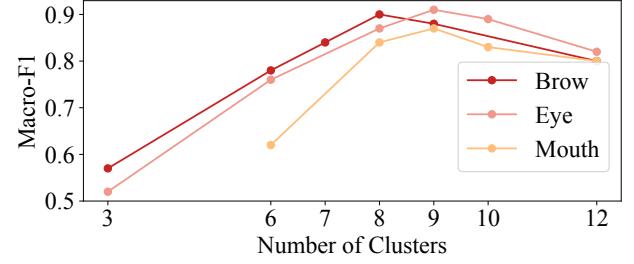


Figure 6. The annotation accuracy for different number of clusters.

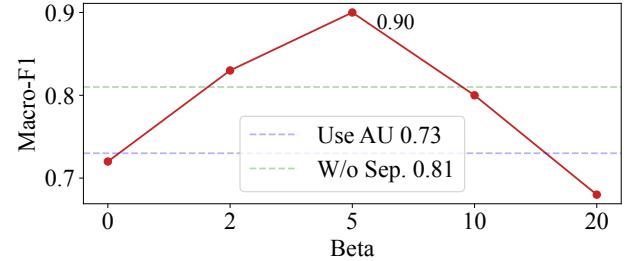


Figure 7. The annotation accuracy in brows for different β , concatenating videos without separating, and using AU as descriptors.

show that the optimal beta value for different regions is 5, with the best number of clusters being 8 for the eye region, 9 for the brow region, and 9 for the mouth region. (The number of clusters here refers to the count of valid clusters, excluding the one that contain only null sequences). The optimal scores for brow, eyes, and mouth are 0.90, 0.91, and 0.87, respectively. When using threshold-based annotation, the scores for eye closure, pose, and gaze are 0.95, 0.87, and 0.89, respectively.

The impact of different cluster numbers is shown in Fig. 6 (with beta set to 5). Too few clusters in TICC lead to under-segmentation, grouping distinct patterns into single clusters. Conversely, too many clusters decreases inter-cluster variability, making clusters less distinct and blurring boundaries.

Fig. 7 illustrates the impact of β , using brow as an example (with the number of clusters set to 8). β controls the smoothness of segment transitions. Too large β leads to over-smoothing. Conversely, a small β makes the model overly sensitive, leading to noisy segments that misrepresent the data's structure. Fig. 7 shows that when using AUs, the score decreases to 0.73. This highlights the importance of accurate facial motion descriptor quality for annotation. Fig. 7 also shows that separating different video clips when analyzing improves performance.

To conduct a preliminary test of TICC's effectiveness in detecting unseen facial motion sequences, we remove 50 manually annotated videos from the dataset, refit TICC, and then have it annotate these videos to calculate the macro-F1 score. The scores are similar to those obtained when the 50 videos are included (brow: 0.89; eyes: 0.88; mouth: 0.85).

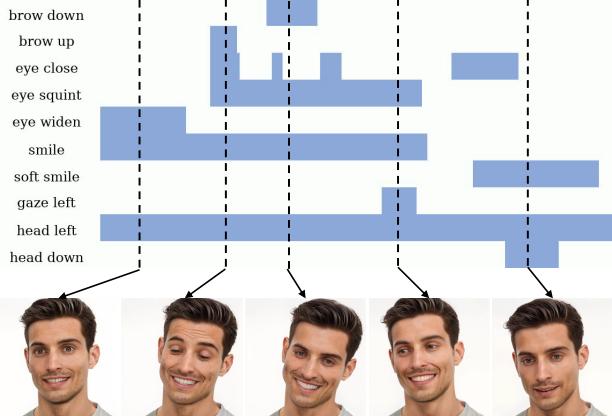


Figure 8. Qualitative results of facial motion generation from the timeline. Better viewed in *Supp. Video*. The timelines for the unploted regions are set to 0. The same applies to the subsequent figures.

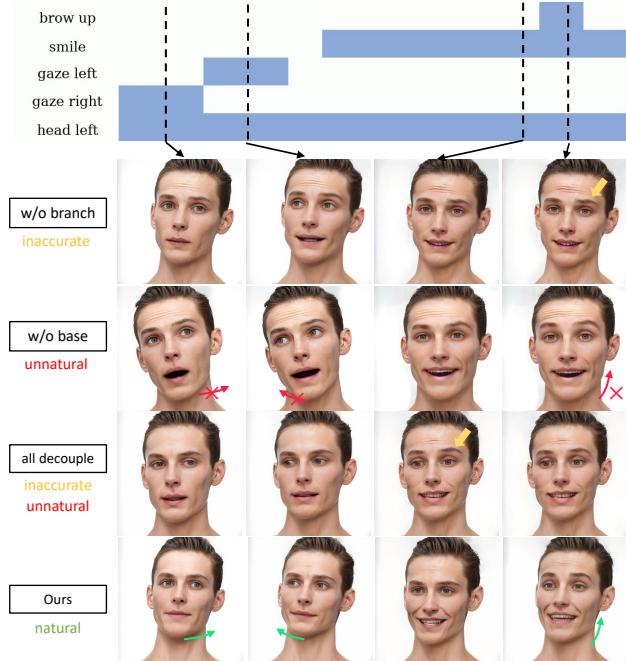


Figure 9. Qualitative results of ablation study. Green arrows indicate subtle natural head motions, red arrows indicate missing subtle natural head motions, and yellow arrows highlight inaccurate motions. Better viewed in *Supp. Video*.

4.2. Facial Motion Generation

Evaluation Metric. We utilize 100 timelines unseen during training to generate samples for evaluating the motion accuracy and naturalness. To evaluate accuracy, we use TICC or threshold-based approach to annotate facial motion intervals within the generated video. These intervals are then compared to the input timeline to calculate the macro-F1 score for each facial region. The average of these scores across all regions is used to assess the alignment between

Methods	Var	\rightarrow	$FID_{fm} \downarrow$	$FID_{\Delta fm} \downarrow$	$SND \downarrow$	$TAS \uparrow$
w/o branch	0.68	\rightarrow	7.39	0.14	7.53	0.66
w/o base	0.64	\rightarrow	12.4	0.18	12.58	0.81
all decoup.	0.41	\rightarrow	28.4	0.23	28.63	0.69
w/o time con.	0.71	\rightarrow	5.38	0.10	5.48	0.79
branchL1	<u>0.70</u>	\rightarrow	6.25	0.11	6.36	0.76
branchL3	0.68	\rightarrow	5.76	0.10	5.86	0.84
branchL4	0.69	\rightarrow	6.64	0.12	6.76	0.83
drop 0	0.62	\rightarrow	6.88	0.13	7.01	0.78
drop 0.3	0.67	\rightarrow	5.93	0.11	6.04	0.83
drop 0.7	0.78	\rightarrow	4.12	0.09	4.21	0.68
Ours	<u>0.70</u>	\rightarrow	4.54	0.09	4.63	0.84

Table 2. Results of ablation study. The unit of Var in the table is 10^{-2} , meaning that 0.70 in the table represents 0.70×10^{-2} . The unit of FID_{fm} and $FID_{\Delta fm}$ is 10^{-1} . Our results differ in magnitude from those of [62] because we use a different 3DMM (FaceVerse) for the calculations. ” \rightarrow ” means results are better if they are close to the variance of real data, which is 0.73.

the generated motion and the input timeline. We refer to this score as Timeline Alignment Score (TAS). To evaluate naturalness, we follow previous methods [44, 61] and use **Var** (variance of generated facial motions, with values closer to GT indicating better performance), **FID_{fm}** (FID score of 3DMM coefficients), **$FID_{\Delta fm}$** (FID score of 3DMM coefficient difference between consecutive frames), **SND** (sum of FID_{fm} and $FID_{\Delta fm}$).

Qualitative Results. Fig. 8 shows the results generated by our model. Our model is capable of producing natural facial movements that are aligned with the input timeline. Note that because the training data includes videos of people speaking, our generated videos sometimes include speaking. Since FaceVerse 3DMM decouples the motion of the left and right face, we can generate motion separately for each side, allowing for generating asymmetric expressions. Our method can generate forehead wrinkles when the speaker raises brows. By introducing descriptors that capture more complex texture-related motions, we can extend our method to control them.

Quantitative Results & Ablation Study. To evaluate the impact of our design choices, we conduct an ablation study with several variants: (1) **w/o branch** remove all branches and only use the base network to predict the facial motions; (2) **w/o base** remove the base network and only use branch networks to generate each region independently; (3) **all decoupl.** decouple all facial and head motions into different branch networks rather than limiting the decoupling to the upper and lower halves of the face. (4) **w/o time con.** starting from the second layer of the base/branch network, the initial timeline token input is replaced with the timeline tokens from the previous layer. We investigate the optimal number of layers for the branch network. The entire net-

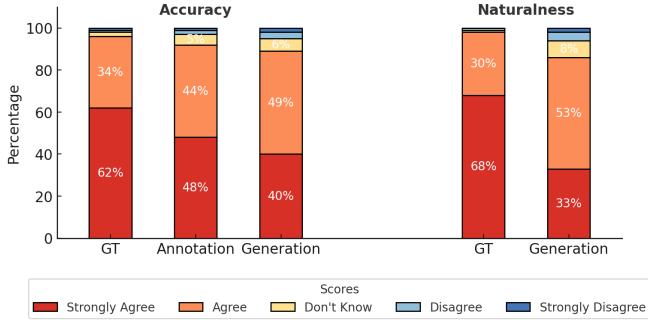


Figure 10. User study results.

work consists of 8 layers, with **branchL_x** indicating that the branch network uses x layers while the base network uses $8 - x$ layers. In the optimal version **Ours**, the branch network has 2 layers. We investigate the optimal drop probability. **drop x** indicates the drop probability is x , the optimal probability used in **Ours** is 0.5

Tab. 2 shows quantitative results and Fig. 9 shows qualitative results from a timeline. **w/o branch** fails to generate accurate motions, highlighting the necessity of using the branch network to enhance accuracy. When the head is oriented to the left, subtle movements may occur due to changes in gaze and brow. While these subtle motions do not affect the overall head direction, they are essential for motion naturalness. **w/o base** fails to generate these subtle motions (Fig. 9). Even though the pose branch receives the complete motion timeline, it still cannot produce such coupling. This highlights the necessity of using the base network to learn the coupling of facial movements. **all decoup.** produces stiff and inaccurate motion changes (brow raised too early in Fig. 9) due to overfitting to the movements of individual regions, resulting in weak generalization to timelines outside the training set. The performance of **w/o time con.** indicates that adding the initial timeline token at each layer can improve accuracy. **branchL_x**'s results show that increasing the number of branch layers excessively reduces naturalness. This is because having too few base layers is insufficient to model the natural coupling between motions. **drop x** 's results show that too low or too high drop probability both cannot generate good performance. Too low probability makes the model overly rely on the condition and hampers its generalization capability. Too high probability makes the model receive too few conditioning signals, reducing its accuracy. **drop 0.7** get better FID scores, but its accuracy is low.

Comparisons with previous methods. No prior methods have achieved fine-grained timeline control of facial motion generation. Similar approaches, such as AgentAvatar [44] and InstructAvatar [49], can only describe temporal changes coarsely using temporal adverbs. Neither of these methods provides open-source code. We can only conduct a qualitative comparison using their demos. Since the temporal

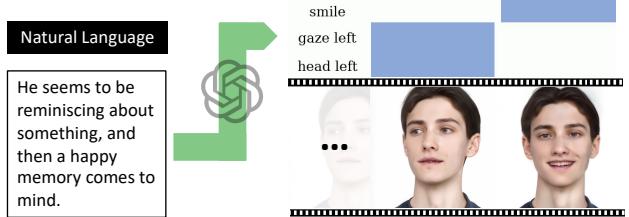


Figure 11. Our method supports text-guided facial motion generation by leveraging ChatGPT to convert text into timelines.

dynamics of facial motion are better demonstrated through video, these comparisons are included in *Supp. Video*.

User Study. We conduct a user study of 21 participants. The participants are required to rate three types of pairs: (1) **GT**: Real videos paired with manually annotated facial motion timelines. (2) **Annotation**: Real videos paired with timelines generated using our method. (3) **Generation**: Videos generated by our method paired with the input facial motion timelines. Participants need to evaluate their level of agreement (on a 5-point scale: strongly agree/disagree, agree/disagree, don't know) on two aspects: (1) Whether the motions in the video **accurately** match the corresponding timelines. (2) Whether the motions in the video appear **natural**. Each participant evaluates 15 pairs sampled from the test data for each type of pair, resulting in a total of 45 pairs per participant. Fig. 9 shows the results. 92% of the evaluations agree that our annotation is accurate. 89% and 86% of the evaluations consider the motions in our generated videos to be accurate and natural, respectively.

Generating Motions Using Natural Language. Our method supports text-guided facial motion generation by leveraging ChatGPT to convert text into timelines. As shown in Fig. 11, ChatGPT can infer that when a person is reminiscing, their head may turn to the side, and their gaze shifts accordingly. When a happy memory comes to mind, the person begins to smile.

5. Conclusion

In this paper, we explore a new control signal: timeline control for facial motion generation. Timeline allows for more fine-grained control than audio or text, enabling users to generate specific motions with precise timing. To model this capability, we first develop a labor-efficient approach to annotate the temporal intervals of facial actions. The annotation process relies on time series analysis of facial motion descriptors. Based on the annotations, we propose a generation model that can generate natural facial motions aligned with timelines. The generation model utilizes a base-branch design to effectively manage motion couplings across different facial regions. Our method supports text-guided generation by using ChatGPT to translate text into timelines. Experiments validate the effectiveness of our method.

References

- [1] Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. Facetalk: Audio-driven motion diffusion for neural parametric head models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21263–21273, 2024. [2](#)
- [2] Nikos Athanasiou, Mathis Petrovich, Michael J Black, and GÜl Varol. Teach: Temporal action composition for 3d humans. In *2022 International Conference on 3D Vision (3DV)*, pages 414–423. IEEE, 2022. [3](#)
- [3] German Barquero, Sergio Escalera, and Cristina Palmero. Seamless human motion composition with blended positional encodings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 457–469, 2024. [3](#)
- [4] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024. URL <https://openai.com/research/video-generation-models-as-world-simulators>, 3, 2024. [1, 2](#)
- [5] Justine Cassell, Catherine Pelachaud, Norman Badler, Mark Steedman, Brett Achorn, Tripp Becket, Brett Douville, Scott Prevost, and Matthew Stone. Animated conversation: rule-based generation of facial expression, gesture & spoken intonation for multiple conversational agents. In *Proceedings of the 21st annual conference on Computer graphics and interactive techniques*, pages 413–420, 1994. [2](#)
- [6] Zhiyuan Chen, Jiajiong Cao, Zhiqian Chen, Yuming Li, and Chenguang Ma. Echomimic: Lifelike audio-driven portrait animations through editable landmark conditions. *arXiv preprint arXiv:2407.08136*, 2024. [2](#)
- [7] Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. Capture, learning, and synthesis of 3d speaking styles. In *CVPR*, pages 10101–10111, 2019. [2](#)
- [8] Radek Danček, Kiran Chhatre, Shashank Tripathi, Yandong Wen, Michael J Black, and Timo Bolkart. Emotional speech-driven animation with content-emotion disentanglement. *arXiv preprint arXiv:2306.08990*, 2023. [2](#)
- [9] Yicheng Deng, Hideaki Hayashi, and Hajime Nagahara. Spotformer: Multi-scale spatio-temporal transformer for facial expression spotting. *arXiv preprint arXiv:2407.20799*, 2024. [2](#)
- [10] Paul Ekman. Facial action coding system (facs). *A Human Face, Salt Lake City*, 2002. [2](#)
- [11] Xiangyu Fan, Jiaqi Li, Zhiqian Lin, Weiyue Xiao, and Lei Yang. Unitalker: Scaling up audio-driven 3d facial animation through a unified model. *arXiv preprint arXiv:2408.00762*, 2024. [2](#)
- [12] Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. Faceformer: Speech-driven 3d facial animation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18770–18780, 2022. [2](#)
- [13] Scott Geng, Revant Teotia, Purva Tendulkar, Sachit Menon, and Carl Vondrick. Affective faces for goal-driven dyadic communication. *arXiv preprint arXiv:2301.10939*, 2023. [5](#)
- [14] Jiazhi Guan, Zhiliang Xu, Hang Zhou, Kaisiyuan Wang, Shengyi He, Zhanwang Zhang, Borong Liang, Haocheng Feng, Errui Ding, Jingtuo Liu, et al. Resyncer: Rewiring style-based generator for unified audio-visually synced facial performer. *arXiv preprint arXiv:2408.03284*, 2024. [2](#)
- [15] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao Sun, Annan Deng, Minglun Gong, and Li Cheng. Action2motion: Conditioned generation of 3d human motions. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2021–2029, 2020. [2](#)
- [16] Chuan Guo, Shihao Zou, Xinxin Zuo, Sen Wang, Wei Ji, Xingyu Li, and Li Cheng. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5152–5161, 2022. [2](#)
- [17] Xupeng Guo, Xiaobiao Zhang, Lei Li, and Zhaoqiang Xia. Micro-expression spotting with multi-scale local transformer in long videos. *Pattern Recognition Letters*, 168:146–152, 2023. [2](#)
- [18] David Hallac, Sagar Vare, Stephen Boyd, and Jure Leskovec. Toeplitz inverse covariance-based clustering of multivariate time series data. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 215–223, 2017. [2, 4, 5](#)
- [19] Tianyu He, Junliang Guo, Runyi Yu, Yuchi Wang, Jialiang Zhu, Kaikai An, Leyi Li, Xu Tan, Chunyu Wang, Han Hu, et al. Gaia: Zero-shot talking avatar generation. *arXiv preprint arXiv:2311.15230*, 2023. [2](#)
- [20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. [5](#)
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851. Curran Associates, Inc., 2020. [4](#)
- [22] Takato Honda, Yasuko Matsubara, Ryo Neyama, Mutsumi Abe, and Yasushi Sakurai. Multi-aspect mining of complex sensor sequences. In *2019 IEEE International Conference on Data Mining (ICDM)*, pages 299–308. IEEE, 2019. [4, 5](#)
- [23] Jianwen Jiang, Chao Liang, Jiaqi Yang, Gaojie Lin, Tianyun Zhong, and Yanbo Zheng. Loopy: Taming audio-driven portrait avatar with long-term motion dependency. *arXiv preprint arXiv:2409.02634*, 2024. [2](#)
- [24] Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. *ACM Transactions on Graphics (TOG)*, 36(4):1–12, 2017. [2](#)
- [25] Haojie Li and Shuangping Huang. Kmtalk: Speech-driven 3d facial animation with key motion embedding. 2024. [2](#)
- [26] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Talkinggaussian: Structure-persistent 3d talking head synthesis via gaussian splatting. In *European Conference on Computer Vision*, pages 127–145. Springer, 2025. [2](#)
- [27] Jun Ling, Yiwen Wang, Han Xue, Rong Xie, and Li Song. Posetalk: Text-and-audio-based pose control and motion refinement for one-shot talking head generation. *arXiv preprint arXiv:2409.02657*, 2024. [2](#)

- [28] Yifeng Ma, Suzhen Wang, Yu Ding, Bowen Ma, Tangjie Lv, Changjie Fan, Zhipeng Hu, Zhidong Deng, and Xin Yu. Talkclip: Talking head generation with text-guided expressive speaking styles. *arXiv preprint arXiv:2304.00334*, 2023. 2
- [29] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5442–5451, 2019. 2
- [30] Yasuko Matsubara, Yasushi Sakurai, and Christos Faloutsos. Autoplait: Automatic mining of co-evolving time sequences. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 193–204, 2014. 5
- [31] Evonne Ng, Javier Romero, Timur Bagautdinov, Shaojie Bai, Trevor Darrell, Angjoo Kanazawa, and Alexander Richard. From audio to photoreal embodiment: Synthesizing humans in conversations. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2024. 2
- [32] Federico Nocentini, Thomas Besnier, Claudio Ferrari, Sylvain Arguillere, Stefano Berretti, and Mohamed Daoudi. Scantalk: 3d talking heads from unregistered scans. *arXiv preprint arXiv:2403.10942*, 2024. 2
- [33] Catherine Pelachaud, Norman I Badler, and Mark Steedman. Generating facial expressions for speech. *Cognitive science*, 20(1):1–46, 1996. 2
- [34] 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. 2
- [35] Mathis Petrovich, Or Litany, Umar Iqbal, Michael J Black, Gul Varol, Xue Bin Peng, and Davis Rempe. Multi-track timeline control for text-driven 3d human motion generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1911–1921, 2024. 3
- [36] Matthias Plappert, Christian Mandery, and Tamim Asfour. The kit motion-language dataset. *Big data*, 4(4):236–252, 2016. 2
- [37] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. A lip sync expert is all you need for speech to lip generation in the wild. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 484–492, 2020. 2
- [38] Fangbing Qu, Su-Jing Wang, Wen-Jing Yan, He Li, Shuhang Wu, and Xiaolan Fu. Cas(me)² : A database for spontaneous macro-expression and micro-expression spotting and recognition. *IEEE Transactions on Affective Computing*, 9(4):424–436, 2018. 2
- [39] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H Bermano. Human motion diffusion as a generative prior. *arXiv preprint arXiv:2303.01418*, 2023. 3
- [40] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 4
- [41] Zhiyao Sun, Tian Lv, Sheng Ye, Matthieu Lin, Jenny Sheng, Yu-Hui Wen, Minjing Yu, and Yong-jin Liu. Diffposetalk: Speech-driven stylistic 3d facial animation and head pose generation via diffusion models. *ACM Transactions on Graphics (TOG)*, 43(4):1–9, 2024. 2
- [42] Shuai Tan, Bin Ji, Mengxiao Bi, and Ye Pan. Edtalk: Efficient disentanglement for emotional talking head synthesis. In *European Conference on Computer Vision*, pages 398–416. Springer, 2025. 2
- [43] Linrui Tian, Qi Wang, Bang Zhang, and Liefeng Bo. Emo: Emote portrait alive-generating expressive portrait videos with audio2video diffusion model under weak conditions. *arXiv preprint arXiv:2402.17485*, 2024. 1, 2
- [44] Duomin Wang, Bin Dai, Yu Deng, and Baoyuan Wang. Agentavatar: Disentangling planning, driving and rendering for photorealistic avatar agents. *arXiv preprint arXiv:2311.17465*, 2023. 1, 2, 7, 8
- [45] Kaisiyuan Wang, Qianyi Wu, Linsen Song, Zhuoqian Yang, Wayne Wu, Chen Qian, Ran He, Yu Qiao, and Chen Change Loy. Mead: A large-scale audio-visual dataset for emotional talking-face generation. In *ECCV*, pages 700–717. Springer, 2020. 5
- [46] Lizhen Wang, Zhiyuan Chen, Tao Yu, Chenguang Ma, Liang Li, and Yebin Liu. Faceverse: a fine-grained and detail-controllable 3d face morphable model from a hybrid dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20333–20342, 2022. 2, 3, 5
- [47] Suzhen Wang, Yifeng Ma, Yu Ding, Zhipeng Hu, Changjie Fan, Tangjie Lv, Zhidong Deng, and Xin Yu. Styletalk++: A unified framework for controlling the speaking styles of talking heads. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 2
- [48] Su-Jing Wang, Ying He, Jingting Li, and Xiaolan Fu. Mesnet: A convolutional neural network for spotting multi-scale micro-expression intervals in long videos. *IEEE Transactions on Image Processing*, 30:3956–3969, 2021. 2
- [49] Yuchi Wang, Junliang Guo, Jianhong Bai, Runyi Yu, Tianyu He, Xu Tan, Xu Sun, and Jiang Bian. Instructavatar: Text-guided emotion and motion control for avatar generation. *arXiv preprint arXiv:2405.15758*, 2024. 1, 2, 8
- [50] Huawei Wei, Zejun Yang, and Zhisheng Wang. Aniportrait: Audio-driven synthesis of photorealistic portrait animation. *arXiv preprint arXiv:2403.17694*, 2024. 2
- [51] Sijing Wu, Yunhao Li, Yichao Yan, Huiyu Duan, Ziwei Liu, and Guangtao Zhai. Mmhead: Towards fine-grained multi-modal 3d facial animation. *arXiv preprint arXiv:2410.07757*, 2024. 1, 2
- [52] Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. Codetalker: Speech-driven 3d facial animation with discrete motion prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12780–12790, 2023. 2
- [53] Mingwang Xu, Hui Li, Qingkun Su, Hanlin Shang, Liwei Zhang, Ce Liu, Jingdong Wang, Luc Van Gool, Yao Yao, and Siyu Zhu. Hallo: Hierarchical audio-driven visual synthesis for portrait image animation. *arXiv preprint arXiv:2406.08801*, 2024. 2

- [54] Sicheng Xu, Guojun Chen, Yu-Xiao Guo, Jiaolong Yang, Chong Li, Zhenyu Zang, Yizhong Zhang, Xin Tong, and Baining Guo. Vasa-1: Lifelike audio-driven talking faces generated in real time. *arXiv preprint arXiv:2404.10667*, 2024. [1](#) [2](#)
- [55] Karren D Yang, Anurag Ranjan, Jen-Hao Rick Chang, Raviteja Vemulapalli, and Oncel Tuzel. Probabilistic speech-driven 3d facial motion synthesis: New benchmarks methods and applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27294–27303, 2024. [2](#)
- [56] Chuin Hong Yap, Connah Kendrick, and Moi Hoon Yap. Samm long videos: A spontaneous facial micro-and macro-expressions dataset. In *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 771–776. IEEE, 2020. [2](#)
- [57] Chuin Hong Yap, Moi Hoon Yap, Adrian Davison, Connah Kendrick, Jingting Li, Su-Jing Wang, and Ryan Cunningham. 3d-cnn for facial micro-and macro-expression spotting on long video sequences using temporal oriented reference frame. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 7016–7020, 2022. [2](#)
- [58] Zhenhui Ye, Ziyue Jiang, Yi Ren, Jinglin Liu, JinZheng He, and Zhou Zhao. Geneface: Generalized and high-fidelity audio-driven 3d talking face synthesis. *arXiv preprint arXiv:2301.13430*, 2023. [2](#)
- [59] Jun Yu, Zhongpeng Cai, Zepeng Liu, Guochen Xie, and Peng He. Facial expression spotting based on optical flow features. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 7205–7209, 2022. [2](#)
- [60] Jianhui Yu, Hao Zhu, Liming Jiang, Chen Change Loy, Weidong Cai, and Wayne Wu. CelebV-Text: A large-scale facial text-video dataset. In *CVPR*, 2023. [2](#)
- [61] Zhentao Yu, Zixin Yin, Deyu Zhou, Duomin Wang, Finn Wong, and Baoyuan Wang. Talking head generation with probabilistic audio-to-visual diffusion priors. *arXiv preprint arXiv:2212.04248*, 2022. [7](#)
- [62] Zhentao Yu, Zixin Yin, Deyu Zhou, Duomin Wang, Finn Wong, and Baoyuan Wang. Talking head generation with probabilistic audio-to-visual diffusion priors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7645–7655, 2023. [7](#)
- [63] Longhao Zhang, Shuang Liang, Zhipeng Ge, and Tianshu Hu. Personatalk: Bring attention to your persona in visual dubbing. *arXiv preprint arXiv:2409.05379*, 2024. [2](#)
- [64] Wenxuan Zhang, Xiaodong Cun, Xuan Wang, Yong Zhang, Xi Shen, Yu Guo, Ying Shan, and Fei Wang. Sadtalker: Learning realistic 3d motion coefficients for stylized audio-driven single image talking face animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8652–8661, 2023. [2](#)
- [65] Qingcheng Zhao, Pengyu Long, Qixuan Zhang, Dafei Qin, Han Liang, Longwen Zhang, Yingliang Zhang, Jingyi Yu, and Lan Xu. Media2face: Co-speech facial animation generation with multi-modality guidance. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–13, 2024. [2](#)