

ON THE ROLE OF LIP ARTICULATION IN VISUAL SPEECH PERCEPTION

Zakaria Aldeneh*, Masha Fedzechkina*, Skyler Seto, Katherine Metcalf,
Miguel Sarabia, Nicholas Apostoloff, Barry-John Theobald

Apple

ABSTRACT

Generating realistic lip motion from audio to simulate speech production is critical for driving natural character animation. Previous research has shown that traditional metrics used to optimize and assess models for generating lip motion from speech are not a good indicator of subjective opinion of animation quality. Devising metrics that align with subjective opinion first requires understanding what impacts human perception of quality. In this work, we focus on the degree of articulation and run a series of experiments to study how articulation strength impacts human perception of lip motion accompanying speech. Specifically, we study how increasing under-articulated (dampened) and over-articulated (exaggerated) lip motion affects human perception of quality. We examine the impact of articulation strength on human perception when considering only lip motion, where viewers are presented with talking faces represented by landmarks, and in the context of embodied characters, where viewers are presented with photo-realistic videos. Our results show that viewers prefer over-articulated lip motion consistently more than under-articulated lip motion and that this preference generalizes across different speakers and embodiments.

Index Terms— speech animation, audio-visual speech, lip sync, human-computer interaction

1. INTRODUCTION

Animating faces from speech has applications in interactive systems, entertainment, and accessibility. Most recent approaches for generating visual speech use neural networks, especially recurrent neural networks (RNNs) [1–6], by training the model to replicate ground-truth reference visual speech from acoustic speech features. Objective metrics typically used to optimize these models, such as the mean squared error (MSE), capture errors globally and can fail to capture *perceptually* significant errors [7–9]. Thus, assessing trained models usually involves some form of *subjective* assessment in the form of pairwise preference testing or mean opinion score (MOS) aggregation [4, 9–16], which are time-consuming and expensive.

Previous work has highlighted that commonly used objective measures for assessing generated visual speech quality are not indicative of subjective opinion of quality [7] as sequences with little difference in MSE can vary considerably in terms of subjective opinion of quality [8]. At issue is the *type* of error and *where* in the sequence errors occur are not taken into account. For example, missing a lip closure during a bilabial plosive (/b/) is more significant perceptually than parting the lips slightly more than is necessary for a velar plosive (/k/). Consequently, several approaches have been proposed to address the limitations of traditional objective metrics. For example, Zhou et al. [13] proposed a collection of metrics computed from various distances in lip/jaw position and velocity. Alternatively, Prajwal et al. used a pre-trained SyncNet model to quantify the quality of visual speech [17]. Chen et al. [18] used the distance between embeddings from a pre-trained lip reading model extracted from real sequences and generated equivalents as a measure of quality. Although previously proposed approaches can mitigate the limitations of traditional metrics, none of the approaches bridge the disconnect between human evaluations and objective measures.

Devising objective metrics that align with perceived quality of visual speech necessarily involves understanding how sources of error influence the perception of quality. In this work, we focus on the degree of articulation and ask whether the perceived quality of visual speech is influenced more by under-articulated (dampened) lip motion or by over-articulated (exaggerated) lip motion. We first study this effect from the perspective of pure motion by rendering the visual speech as facial landmarks, akin to point-light displays used to understand speech production. Then, we examine the impact of the degree of articulation in the context of a photo-realistic talking face created by a state-of-the-art approach [13]. We find that over-articulation impacts the perceived quality much less than under-articulation does. We conclude with a discussion about using human perception of visual speech quality to improve model development and evaluation.

2. DATASET AND DATA PREPARATION

We use the GRID corpus [19], which contains multimodal recordings of 34 native English speakers reciting fixed gram-

*Authors contributed equally.

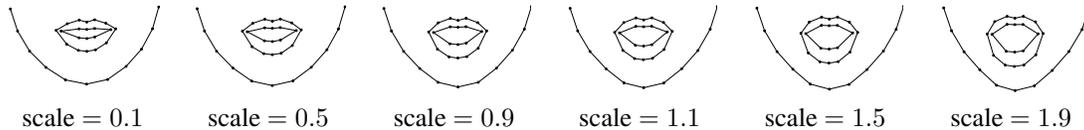


Fig. 1: Illustration of the effect of scaling the PCA values of the facial landmarks for the same video frame.

mar sentences of the form, *command-color-preposition-letter-number-adverb*, where each token in the grammar has a limited number of options (see [19] for details). An example utterance might be “place blue at A 3 now”. The dataset contains all combinations of the colors, letters, and numbers. Each speaker recited a total of 1000 sentences, and each video is three seconds long recorded at 25 frames per second.

We use Dlib [20] to extract 68 facial landmarks from each video frame and then reduce the effects of tracking noise by convolving the sequence across time with an averaging window of width three frames. We then align the facial landmarks to a common reference [1] to remove translation, rotation, and scaling variation.

To manipulate the degree of speech articulation, we project the landmarks onto a speaker-specific principal components analysis (PCA) model, and then for each utterance we randomly sample four scaling factors from $\{0.1, 0.3, 0.5, 0.7, 0.9, 1.0, 1.1, 1.3, 1.5, 1.7, 1.9\}$. These scaling factors were chosen to have paired values for under- and over-articulation. For example, the values of 0.9 and 1.1 are equivalently scaled either side of no change in articulation (1.0). The scaled PCA values were finally projected back to landmarks for use in generating the stimuli for the user-studies.

3. SUBJECTIVE ASSESSMENT OF VISUAL SPEECH QUALITY

User study setup. We use the perturbed landmarks from Section 2 to generate the stimuli for three experiments which assess the perceptual quality of visual speech. In these experiments, viewers were shown pairs of videos: a reference video (no scaling) and a (possibly) modified version of the same sequence containing under- or over-articulated lip motion. The order of the videos within a pair was always the same (reference video on the left and sample video on the right). Annotators were asked to indicate using a seven-point Likert scale if the sample animation looked: {Extremely Worse, Moderately Worse, Slightly Worse, The Same, Slightly Better, Moderately Better, Extremely Better} compared to the reference animation. All annotators spoke English, based on the crowdsourcing platform’s screening process.

Data analysis. We analyzed annotator perceptual judgments in all user studies using linear mixed-effects regres-

sion.¹ For the analysis, we decompose the scaling factor into two categorical variables: articulation type, representing the direction of scaling, and the scaling step. For example, a scaling step of 1 corresponds to scaling the PCA values by 1.1 for over-articulation and scaling the PCA values by 0.9 for under-articulation; a scaling step of 2 corresponds to scaling the PCA values by 1.3 for over-articulation and scaling the PCA values by 0.7 for under-articulation; and so on. Trials with a scaling factor of 1.0 (ground truth) were removed from this analysis since the scaling step did not apply. Scaling step was sliding difference coded.² Articulation type was sum-coded (under-articulation = 1). All models include the maximal converging random effects structure for the experimental design (a random intercept for annotator and a random intercept for item nested within a ground truth sequence, relating to the version of the ground truth sequence for Experiment 1; random intercepts for annotator, item, and speaker for Experiments 2 and 3).

3.1. Experiment 1: Rendering on Point-light Displays

Our first user study aims to measure the impact of under-articulated and over-articulated lip motion on the perceived quality of visual speech rendered as landmarks.

To remove inter-speaker variability effects from this initial analysis, we select 996 sequences for speaker $s-25$.³ This speaker was selected as they have the largest variance in lip-opening height, measured as the distance between the mid-points of the upper and lower lips. To ensure we focus exclusively on speech-related facial motion when assessing the quality of visual speech, we remove landmarks that correspond to the eyes, nose, and eyebrows. As a result, each frame is represented by 31 landmarks, where each landmark is a pair of x - y points, which were rendered to videos. See Figure 1 for example rendered video frames. The videos were presented to annotators as described in the setup.

¹We follow the recommendation by [21] on using linear regression for ordinal scale data.

²Sliding difference coding compares the mean of the dependent variable for one level of the categorical variable to the mean of the dependent variable for the preceding adjacent level (e.g., scaling step 2 vs. scaling step 1).

³The landmark extraction failed for videos: `prat5s`, `pwin2n`, `bbwq4n`, and `brwk8p`.

Table 1: Model summary for the analysis of perceptual judgments for landmarks-only representation. See Section 3.1 for model details. Step is scaling step (1-5, sliding difference coded); Art. is articulation type (sum-coded, under-articulation = 1).

	Estimate	t	p
Art.	-0.472	-53.53	< 0.0001
Step 2 vs. 1	-0.09	-3.07	0.002
Step 3 vs. 2	-0.27	-9.75	< 0.0001
Step 4 vs. 3	-0.56	-20.19	< 0.0001
Step 5 vs. 4	-0.58	-21.21	< 0.0001
Step 2 vs. 1 * Art.	-0.16	-4.43	< 0.0001
Step 3 vs. 2 * Art.	-0.18	-6.39	< 0.0001
Step 4 vs. 3 * Art.	-0.32	-11.75	< 0.0001
Step 5 vs. 4 * Art.	-0.43	-15.76	< 0.0001

Results

Across all scaling steps, annotators preferred over-articulated lip motion (see Table 1). Every scaling step increase resulted in a lower perceptual score compared to the previous step, suggesting that the more the degree of speech articulation is modified, the more negative impact it has on the perception of its quality. However, articulation type interacted with scaling step—an increase in scaling had a more pronounced effect on under-articulation compared to over-articulation. Interestingly, over-articulation was rated as high as the ground truth, or even better, at all but the highest scaling steps (see Figure 2). These findings suggest that while viewers are sensitive to the degree of visual speech articulation, the sensitivity has less of an effect for over-articulation compared to under-articulation, and that viewers generally prefer over-articulation.

3.2. Rendering as Photo-realistic Sequences

We use the framework proposed by Zhou et al. [13] to generate photo-realistic sequences corresponding to dampened/exaggerated visual speech. This framework generates a video sequence corresponding to a talking face in two stages: first facial landmarks are predicted from speech, then the landmarks are used to animate a reference still-image using an image2image translation module similar to the model proposed by Zakharov et al. [4]. This two-stage pipeline allows first testing photo-realistic sequences generated from scaled ground-truth speech articulation and from scaled landmark sequences that are predicted from the entire talking face pipeline.

Table 2: Model summary for the analysis of perceptual judgments for photo-realistic talking face. See Section 3.2.2 for model details. Step is scaling step (sliding difference coded, 3 vs. 1; 5 vs. 3); Art. is articulation type (sum-coded, under-articulation = 1); Exp. is Experiment (sum-coded, Exp.2 = 1); speaker gender is sum-coded, male = 1.

	Estimate	t	p
Art.	-0.73	-16.41	< 0.0001
Step 3 vs. 1	-0.35	-15.69	< 0.0001
Step 5 vs. 3	-0.75	-32.94	< 0.0001
Exp.	0.01	0.44	0.66
Speaker gender	0.01	1.35	0.18
Step 3 vs. 1 * Art.	-0.59	-26.32	< 0.0001
Step 5 vs. 5 * Art.	-0.66	-28.83	< 0.0001
Exp. * Art.	0.05	6.67	< 0.0001
Step 3 vs. 1 * Exp.	0.02	1.33	0.18
Step 5 vs. 3 * Exp.	-0.07	-4.06	< 0.0001
Step 3 vs. 1 * Art. * Exp.	0.06	3.35	< 0.0001
Step 5 vs. 3 * Art. * Exp.	-0.01	-0.83	0.40

3.2.1. Experiment 2: Using ground-truth landmarks

We first randomly sample 30 utterances from each speaker in the GRID corpus to obtain 990 utterances.⁴ For the selected utterances, we scale the degree of articulation in the visual speech, as described in Section 2, but using only scaling factors: {0.1, 0.5, 0.9, 1.0, 1.1, 1.5, 1.9} since a power analysis of the results in Section 3.1 showed we would still maintain 100% power with this reduced set. The stimuli for the user-study are then created by reconstructing the landmarks from the scaled PCA values and inputting these landmarks with the first video frame from the original sequence to the image2image translation module.

3.2.2. Experiment 3: Using landmarks predicted from speech

For the same utterances in Section 3.2.1, we compute the predicted landmarks from acoustic speech using the pre-trained model provided by Zhou et al. [13]. These are then subject to the same scaling as in Section 2. The difference between this experiment and that in Section 3.2.1 is that the landmarks here are subject to all inaccuracies present in a state-of-the-art talking face system, and we can determine whether the under-vs. over-articulation findings still hold. In other words, does the conclusion that people prefer over-articulation still hold if the over-articulation emphasizes errors made by the model that predicts landmark sequences?

Results

Overall, the results from the experiments with photo-realistic sequences have similar trends to the results from the exper-

⁴Video data for s-21 is missing from the corpus.

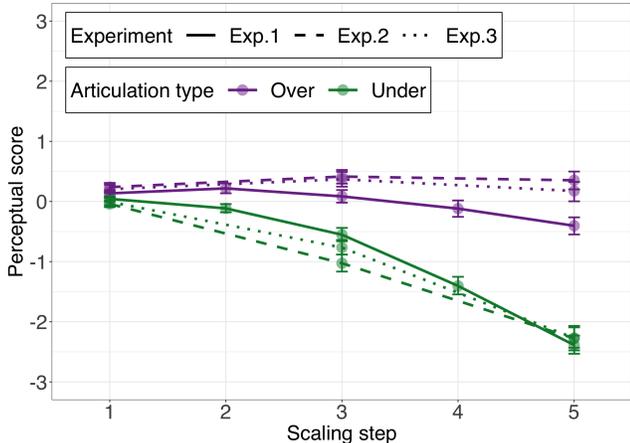


Fig. 2: Mean perceptual scores for over- (purple) and under-articulated (green) visual speech under varying scaling steps in the three experiments. A perceptual score of -3 , -2 , -1 , 0 , 1 , 2 , and 3 on the y -axis corresponds to “Extremely Worse”, “Moderately Worse”, “Slightly Worse”, “The Same”, “Slightly Better”, “Moderately Better”, and “Extremely Better”, respectively. The solid lines represent scores for point-light displays (Exp.1) (Section 3.1), the dashed lines represent scores for the photo-realistic model driven by ground-truth landmarks (Exp.2) (Section 3.2.1), and the dotted lines represent the photo-realistic model driven by landmarks predicted from speech (Exp.3) (Section 3.2.2). The error bars represent bootstrapped 95% confidence intervals.

iment that used point-light displays (see Table 2 and Figure 2). Applying an increasingly larger magnitude scaling to increase the degree of under- or over-articulation resulted in lower perceived quality. Over-articulation was preferred overall, and the preference becomes more pronounced at each increased scaling step. Overall the perceptual scores for the photo-realistic sequences did not differ when either ground-truth landmarks or predicted landmarks were used. However, the difference between under- and over-articulation was more pronounced for the ground-truth landmarks, but this difference was only present at higher scaling steps, suggesting that this difference is somewhat small. Speaker gender did not influence perceptual scores.

We run an additional analysis where we compare the perceptual scores for the point-light display videos (using scaling steps 1, 2, and 3) and the photo-realistic sequences (collapsing over the predicted and ground-truth landmarks). Photo-realistic sequences were rated higher overall. Additionally, the difference between under- and over-articulation was larger for photo-realistic sequences than for videos using point-light displays ($p < 0.001$ for all).

The findings from these experiments suggest that the preference for over-articulation generalizes from point-light displays to photo-realistic stimuli of different speakers, and that

this preference is stronger for photo-realistic stimuli. Additionally, we find that the perception of photo-realistic sequences does not differ based on whether image2image model is conditioned on ground-truth or predicted landmarks.

4. CONCLUSIONS

In this work, we studied the relationship between the degree of articulation and the perceived quality of visual speech. We showed that varying the strength of over- and under-articulation affects the perceived quality. Specifically, we found that individuals consistently preferred visually over-articulated speech to visually under-articulated speech for all scaling steps. We also found that while the perceptual score goes down for high scaling steps, the reduction in the score is more pronounced for visually under-articulated speech. The observed preference for over-articulation in point-light displays is even more pronounced for more photo-realistic sequences across different speakers. Finally, the perception of photo-realistic sequences does not differ based on whether the image2image network is driven by ground-truth or predicted landmarks.

5. FUTURE WORK

The findings from this work impact lip animation models in two aspects. First, the findings from our perceptual studies can be incorporated into training algorithms to yield more natural animations. Specifically, loss functions for neural-based lip-sync systems can be re-weighted to account for the perceived effects of articulation errors (i.e., err on the side of over-articulation for the best perceived outcome). Second, our findings shed light on how lip motion can influence subjective evaluation during model development and benchmarking. Future work will investigate the effect of other common augmentations (e.g., jitter, synchronization) on the perceived quality of lip motion and will also investigate the automatic prediction of the perceptual scores.

6. ACKNOWLEDGEMENTS

Thanks to our colleagues, Vikram Mitra, Russ Webb, and Ahmed Hussen Abdelaziz, for their insightful feedback on this work. Also, thanks to Lukas Michelsen, the annotation team, and the annotators for their help in conducting the perceptual studies.

7. REFERENCES

- [1] Sefik Emre Eskimez, Ross K Maddox, Chenliang Xu, and Zhiyao Duan, “Generating talking face landmarks from speech,” in *International Conference on Latent Variable Analysis and Signal Separation*, 2018.

- [2] David Greenwood, Iain Matthews, and Stephen Laycock, "Joint learning of facial expression and head pose from speech," *Proc. Interspeech*, 2018.
- [3] Olivia Wiles, A Koepke, and Andrew Zisserman, "X2face: A network for controlling face generation using images, audio, and pose codes," in *Proceedings of the European conference on computer vision (ECCV)*, 2018.
- [4] Egor Zakharov, Aliaksandra Shysheya, Egor Burkov, and Victor Lempitsky, "Few-shot adversarial learning of realistic neural talking head models," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019.
- [5] Amir Jamaludin, Joon Son Chung, and Andrew Zisserman, "You said that?: Synthesising talking faces from audio," *International Journal of Computer Vision*, 2019.
- [6] Ting-Chun Wang, Arun Mallya, and Ming-Yu Liu, "One-shot free-view neural talking-head synthesis for video conferencing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [7] Barry-John Theobald and Iain Matthews, "Relating objective and subjective performance measures for aam-based visual speech synthesis," *IEEE transactions on audio, speech, and language processing*, 2012.
- [8] Ahmed Hussen Abdelaziz, Barry-John Theobald, Paul Dixon, Reinhard Knothe, Nicholas Apostoloff, and Sachin Kajareker, "Modality dropout for improved performance-driven talking faces," in *Proceedings of the 2020 International Conference on Multimodal Interaction*, 2020.
- [9] Danny Websdale, Sarah Taylor, and Ben Milner, "Speaker-independent speech animation using perceptual loss functions and synthetic data," *IEEE Transactions on Multimedia*, 2021.
- [10] Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia Rodriguez, Jessica Hodgins, and Iain Matthews, "A deep learning approach for generalized speech animation," *ACM Transactions on Graphics (TOG)*, 2017.
- [11] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman, "Synthesizing obama: learning lip sync from audio," *ACM Transactions on Graphics (ToG)*, 2017.
- [12] Ahmed Hussen Abdelaziz, Barry-John Theobald, Justin Binder, Gabriele Fanelli, Paul Dixon, Nick Apostoloff, Thibaut Weise, and Sachin Kajareker, "Speaker-independent speech-driven visual speech synthesis using domain-adapted acoustic models," in *2019 International Conference on Multimodal Interaction*, 2019.
- [13] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li, "Makelttalk: speaker-aware talking-head animation," *ACM Transactions on Graphics (TOG)*, 2020.
- [14] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic, "Realistic speech-driven facial animation with gans," *International Journal of Computer Vision*, 2020.
- [15] Xinya Ji, Hang Zhou, Kaisiyuan Wang, Wayne Wu, Chen Change Loy, Xun Cao, and Feng Xu, "Audio-driven emotional video portraits," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
- [16] Ahmed Hussen Abdelaziz, Anushree Prasanna Kumar, Chloe Seivwright, Gabriele Fanelli, Justin Binder, Yannis Stylianou, and Sachin Kajareker, "Audiovisual speech synthesis using tacotron2," in *Proceedings of the 2021 International Conference on Multimodal Interaction*, 2021.
- [17] 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*, 2020.
- [18] Lele Chen, Guofeng Cui, Ziyi Kou, Haitian Zheng, and Chenliang Xu, "What comprises a good talking-head video generation?," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020.
- [19] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu Shao, "An audio-visual corpus for speech perception and automatic speech recognition," *The Journal of the Acoustical Society of America*, 2006.
- [20] Davis E King, "Dlib-ml: A machine learning toolkit," *The Journal of Machine Learning Research*, 2009.
- [21] Jeff Sauro and James Lewis, *Quantifying the user experience: Practical statistics for user research*, Elsevier, 2012.