# CONVERSATIONAL CO-SPEECH GESTURE GENERATION VIA MODELING DIALOG INTENTION, EMOTION, AND CONTEXT WITH DIFFUSION MODELS

*Haiwei Xue*[1], *Sicheng Yang*[1], *Zhensong Zhang*[2], *Zhiyong Wu*[1,3,*],
*Minglei Li*[4,*], *Zonghong Dai*[4], *Helen Meng*[3]

[1] Shenzhen International Graduate School, Tsinghua University, Shenzhen, China
[2] Huawei Noah's Ark Lab, Shenzhen, China
[3] The Chinese University of Hong Kong, Hong Kong SAR, China
[4] Huawei Cloud Computing Technologies Co., Ltd, Shenzhen, China

xhw22@mails.tsinghua.edu.cn, yangsc21@mails.tsinghua.edu.cn, zhangzhensong@huawei.com,
zywu@sz.tsinghua.edu.cn, liminglei29@huawei.com, daizonghong@huawei.com, hmmeng@se.cuhk.edu.hk

## ABSTRACT

Audio-driven co-speech human gesture generation has made remarkable advancements recently. However, most previous works only focus on single person audio-driven gesture generation. We aim at solving the problem of conversational co-speech gesture generation that considers multiple participants in a conversation, which is a novel and challenging task due to the difficulty of simultaneously incorporating semantic information and other relevant features from both the primary speaker and the interlocutor. To this end, we propose CoDiffuseGesture, a diffusion model-based approach for speech-driven interaction gesture generation via modeling bilateral conversational intention, emotion, and semantic context. Our method synthesizes appropriate interactive, speech-matched, high-quality gestures for conversational motions through the intention perception module and emotion reasoning module at the sentence level by a pretrained language model. Experimental results demonstrate the promising performance of the proposed method.

*Index Terms*— Co-speech gesture generation, interaction gesture, dialog intention and emotion, multi-agent conversational interaction

## 1. INTRODUCTION

Co-speech gesture is very important in daily communication [1], which complements the speech and makes the speech more engaging and vivid. For example, "you" word is often accompanied by an implicit gesture of pointing the hand at the listener. When saying the "cut" word, we may act the cutting gesture. What's more, when someone laughs, his body will have a rhythmic tremor in most cases, and some people may hold their stomach. In fact, there is an implicit connection between speech and gesture [2].

Recently, audio driven co-speech gesture generation has drawn much attention from the community [3, 4, 5, 6, 7], due to its wide applications in the industry, such as digital human, CyberVerse, game and movie etc. Besides audio, some methods also consider other modalities, such as text and speaker identity information [6].

While significant progress has been made, most of these works focus on single speaker speech-driven gesture generation. In real-life scenarios, conversations typically involve multiple persons, so speech-driven gesture generation in conversational interactions is
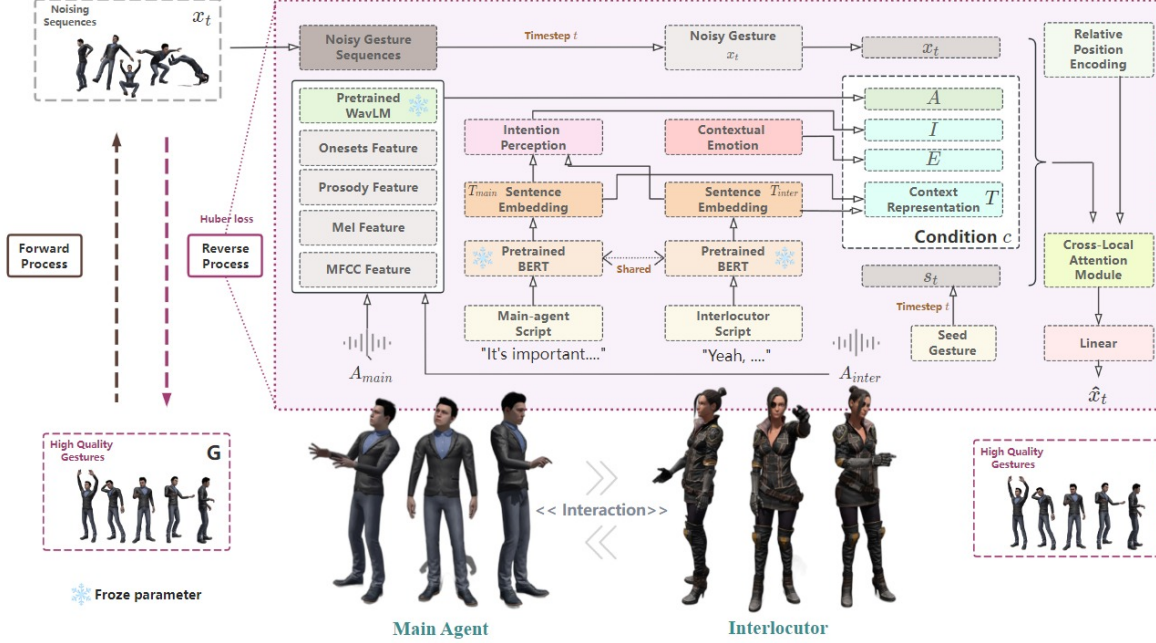
under-studied. Yonatan Shafir et al. [8] is an earlier work that implements text driven gesture generation for two-person interaction with diffusion model. Chopin et al. [9] introduces a bipartite graph approach to correlate multiple people's actions together to generate higher quality interaction actions. But these works do not take into account the semantic information of speech, while Evonne Ng et al. [10] considers two-person speech interaction, but they only model the face.

This paper aim at solving the novel task of multi-person conversational audio-driven gesture generation, which is challenging due to the following reasons: 1) It is difficult for the main speaker's generated gesture to simultaneously encompass both the main speaker's and the interlocutor's speech semantic information; 2) The gestures generated by the main speaker are struggling to appropriately accommodate the social characteristic of all participants in the conversation; 3) In multi-agent interactions, gesture generation lacks diversity and human-likeness.

To address the above challenges, we first employ a large pretrained language model to infer sentence-level contextual representations instead of the word-level static word embedding approach in previous works. Motivated by existing social psychology research [11, 12], we observe that each person enters an interactive environment with specific emotions, attitudes, intentions, and behavioral tendencies. Based on this observation, we then introduce conversational intention and contextual affective information as social characteristics and leverage the fine-tuned language models for the two downstream tasks of intention and emotion. The inferred information is treated as a generative gesture condition. Finally, we redesign the generation conditions and use a variant of the diffusion model to achieve more diverse and human-like gesture generation for conversational interaction.

The main contributions of our paper are: (1) We model contextual information and conversational intent, allowing the generated gestures to have more appropriate and human-like expressions across various conversational topics. (2) We propose to consider the emotional aspects of all participants, aiming to associate the speaker's emotions with their corresponding actions. (3) We are the first to apply the diffusion model to the multimodal conversational speech-driven interaction gesture generation task. We propose an interaction gesture generation framework via encoding speech semantics, dialog intents, and emotions of multiple participants to generate higher quality gesture motions through the diffusion model.

---

**Fig. 1**. (Left) Forward process of CoDiffuseGesture is to add noise to the motion sequence from $t = 0$ until $t = T$. (Right) Reverse process of CoDiffuseGesture is to learn the denoising ability. The script of the main speaker and the interlocutor is translated to the corresponding semantics embeddings by the pre-trained bert model. Then the semantic embedding is sent into the fine-tuned intent-aware module and emotion capture module respectively. A step $t$ and a noisy gesture sequence $X_t$ at the noising step conditioning on $c$ incorporating intent and emotion features $E, I$, audio $A$, seed gesture sequence $d$ are fed to the network. Cross-local attention and self-attention can better capture the correlation between speech and gesture of multi-agents for better quality of gesture generation.

## 2. METHODOLOGY

### 2.1. Problem Definition

We aim to synthesize the high-quality gestures of the main speaker with the multimodal information of the dialogue as input, including the main speaker's speech and script, the interlocutor's speech and script, the social feature of the dyadic non-verbal interaction. Given two speakers' (the main agent and the interlocutor) audio and their corresponding speech scripts as inputs, the goal of our model $F$ is to infer continuous diverse appropriate, human-like co-speech gestures as $G \in [G_1, ...., G_N]$, where $N$ denotes the length of speech audio $A$ and script $T$. Here, $G_i$ represents $J$ joints of the human body including two hands, which can be visualized in Blender.

Inspired by the Dyadic Nonverbal Interaction model [11], when two people have a conversation, the intention of the conversation and the emotion of the speech have a significant impact on the gestures of the speaker. Therefore, to obtain diverse co-speech gestures that are more human-like and appropriate for agent speech and interlocutor behavior, we introduce a conversational intention label $I$ and a multi-label emotion intensity value $E$ to provide supervision during the training phase. Formally, given a conversational sequence clip, there is a generation condition $c = [A, T, I, E]$ for synthesizing main agent gestures. Intent features are extracted from a fine-tuned pre-trained model of Distibert-Base-Uncased [13] trained on a massive dataset. Notice that the intent label actually is the one-hot vector in our approach. The emotion label $E$ contains positive, neutral, and negative emotion intensity values. In our approach, we employ a learning framework and procedure that resembles MDM [14]. The forward process adds Gaussian noise to the gesture sequence at each

diffusion step $t$. The reverse process is trained by optimizing the Huber loss in given condition $c$. The architecture of our proposed CoDiffuseGesture is shown in Figure 1.

### 2.2. Modeling Dialog Intention, Emotion, and Context

Inspired by Dyadic Nonverbal Interaction Systems (DNI) [11] and Basic Emotion Theory (BET) [12], we observe that each person enters an interactive environment with specific emotions, attitudes, intentions, and behavioral tendencies. Therefore, we model the intention, emotion, and context information of dyadic non-verbal interactions in generation condition $c$. The DNI model introduces a mental model in dyadic non-verbal interactions: 1) Perceptual processes ; 2) Cognitive resources + cognitive-affective; 3) Processes; 4) Goal; 5) Behavior; 6) Interaction. Through social psychological model [15], it is found that one of the important factors affecting the interaction and the resulting behavior is the goal, that is, the intention of the two parties during the conversation. This is also consistent with our intuition that when a conversation involves asking for directions, the person giving directions often combines gestures to indicate the correct path. And during casual conversations, there may be more body shakes caused by laughter.

Hence, modeling the intention during the conversation can make the generated gestures more appropriate to the interaction content of the current conversation. We use a fine-tuned version of Distilbert-Base-Uncased [13] on the massive dataset to reason the possible conversation intentions coarse-grainedly. The massive dataset includes nearly 60 different conversational intention labels. Such as transportation, music preferences, interest chat, daily socializing, takeout, etc. Based on the DNI architecture, the intention of the two parties

to chat is not invariable, and they may change during the interaction process. So the intention is continuously inferred from the short timestamp. It should be noted that the result of the final inference will be mapped to the one-hot vectors with the dimension of 60.

Another important influence on interactive gestures is emotion, which BET theory [12] states is fundamentally about instigating action and changing the probabilities of future actions. Some studies [16, 17] on emotion also point out that emotional expression conveys four kinds of information about interaction: (1) how the speaker is feeling at the moment; (2) what is happening in the current context; (3) the perceived action or behavior expected by the other person; and (4) the intention or plan that the two parties may make. The most common example is when a person is extremely angry and is likely to make fist gestures. Thus, the emotions of both parties in a conversation are one of the most important influence factors, which is confirmed in the cognitive-affective processes in the DNI architecture.

Based on the above research [12, 7], we decided to introduce the emotional information of both parties in the generation of interactive co-speech gestures. We employed the Roberta [18] model, which has been pre-trained by Twitter data, to extract the emotional information of both parties at the sentence level. Afterward, we deduce three kinds of emotional intensity of positive, neutral and negative, symboling emotional information $E$. Finally, to enhance the understanding of semantic content, we utilize a pre-trained language model [19] to represent contextual semantic information at the sentence level instead of FastText [20] in previous work [21]. Finally, the generative condition $c$ consists of speech information $A_{main}$ $A_{inter}$, dialogue intention features $I$, emotion features of both parties $E$, and contextual dialogue semantics $T_{main}$ $T_{inter}$.

## 2.3. Diffusion Training Architecture

In our work, we follow the DiffuseStyleGesture [5] framework, but instead of restricting to a single agent only, we generalize to two-person dyadic interactions for conversational co-speech gesture generation. Similar to most diffusion model methods, both include the forward process $F_+$ and the reverse process $F_-$ for training.

**Forward Process** The forward process involves repeatedly adding a small amount of Gaussian noise to the real data until the data exhibits Gaussian noise. Formally, the forward process on a real sample from a real data distribution consists of a Markov chain with gradually increasing noise. The distribution will eventually resemble a standard Gaussian.

$$x_t, \ t = F_+(G, \mathcal{N}(0, I)) \tag{1}$$

**Reverse Process** The reverse process is also called the denoising module. This module needs to be trained together under $x_t$ pure noise, the noise steps $t$, the seed gesture $s_t$, and the conditions $c$. We predict the signal $\hat{x}$ itself rather than the noise in each noise step in denoising process.

$$\hat{x} = F_-(x_t, s_t, t, c) \tag{2}$$

where $c = [A, T, I, E]$. The audio features $A$, context $T$, and emotion intensity $E$ contain multimodal representations of the two individuals in a conversation, which are concatenated modal features of the main agent and the interlocutor, respectively. To train the reverse process for denoising, we optimize the Huber Loss between the generated and the real gestures. After the training is completed, in the inference procedure, it is only necessary to process the speech and text into corresponding features as conditions $c$ in the same way. And then randomly sample part of the gesture sequence from the

training data as the seed gestures $s$. Then pure noisy $x_t$, the seed gestures $s$, step $t$, and the condition $c$ can be feed into the trained denoising module to get the final predicted gesture result. Modeling the interaction information of two people can improve the quality of gesture generation, according to the experiment.

## 3. EXPERIMENTS

### 3.1. Dataset

Currently, there is still a shortage of high-quality datasets for conversational co-speech gesture generation. Fortunately, the GENEA2023 dataset [22] provides approximately 20 hours of excellent dyadic non-verbal interaction co-speech gesture data. This dataset is derived from the original Talking With Hands (TWD) [23] data, which is processed by the organizer to provide both sides' voice, speaker ID, and both sides' motion information saved through BVH motion capture files. Conversation topics are mostly related to our daily interactions. Each annotated clip at millisecond accuracy contains essentially only 1–2 words. Meanwhile, laughter is specially labeled, which is represented by the "#" symbol. The GENEA 2023 dataset contains a total of 483 dialogs including audio files, script and BVH motion capture files of both the main agent and the interlocutor. There are 372 dialogs in the training set, 41 in the validation data, and 70 in the test data. The total amount of the data is around 20 hours. The test data lacks the gesture file of the main agent speaker, which needs to be predicted by the model. In particular, 29 of the 70 test set of data do not match the main speaker and the interlocutor data, which can be used for the study of matching interlocutor consciousness.

### 3.2. Experiment Setup

We merge the training and validation data and split each data into 5-second clips of 30 FPS meta-training data for updating the model. The first second is treated as a seed gesture, and the subsequent duration is used to minimize the loss between the real and generated gestures to train the network. All feature is normalized before feeding into the model to reduce the risk of gradient explosion or dispersion.

We train our model on an NVIDIA A100 40G GPU with a batch size of 400 and AdamW optimizer with learning rate is $3 * 10^{-5}$. We train about a week on GENEA2023 for 500,000 epochs with $T = 1000$ noising steps and a cosine noise schedule. We evaluate the quality of dyadic nonverbal interaction gesture generation from three aspects: human-likeness, appropriateness for agent speech, appropriateness for interlocutor status.

During the experiment, participants score three aspects while watching gestures generated by blender renderings visualizations. The participants primarily consisted of graduate students and working professionsals, ranging in age from 20 to 40. Each evaluation aspect contains 1–5 points, which are excellent, good, fair, poor, and bad. We ensured that the participants remained unaware of the specific method variant associated with each video during the rating process. For the human-likeness assessment dimension, we asked the participant the question, "Does he resemble human movements? Do his actions resemble the kind of gestures a man might make?" We evaluate the appropriateness of main agent speech by asking the participants whether the main agent's gesture matches his speech and has a certain sense of rhythm. For the final evaluation indicator, we asked participants the question, "When the interlocutor is talking, what gesture might the main agent make? Is the listener's response human-like?" Based on the results report, it is found that
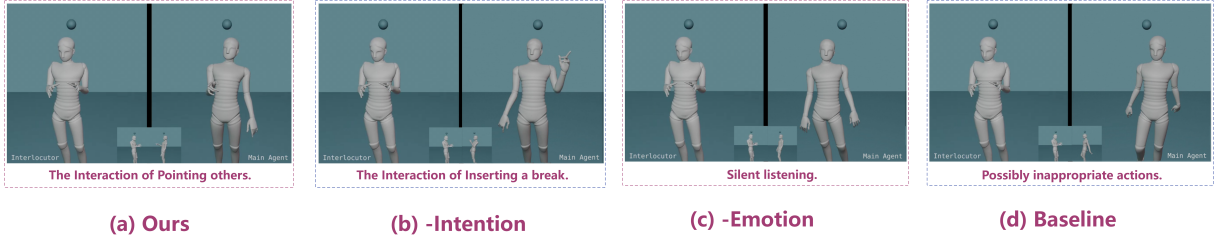
**(a) Ours**     **(b) -Intention**     **(c) -Emotion**     **(d) Baseline**

**Fig. 2**. Visualization of different gestures reflected by the four variational models.

**Table 1**. Evaluation results of different variant model. Human-likeness, Appropriateness for Agent Speech (AAS), and Appropriateness for Interlocutor status (AIS) are results of MOS with 95% confidence intervals. '-' is shorthand for 'without' in ablation studies. '*' denotes the proposed method.

|  | Human-likeness | AAS | AIS |
|---|---|---|---|
| Baseline | 3.625±0.1816 | 3.663±0.1573 | 3.494±0.1883 |
| Ours* | **3.875±0.1481** | 3.775±0.1298 | 3.606±0.1681 |
| - Intention | 3.763±0.1536 | **3.831±0.1510** | 3.550±.1517 |
| - Emotion | 3.781±0.1515 | 3.813±0.1366 | **3.650±0.1487** |

our proposed method has good quality for generating dyadic nonverbal interaction co-speech gestures in conversation. It can take into account both speaking and listening states to make appropriate and human-like gestures.

### 3.3. Experimental Results and Analysis

Table 1 shows the evaluation results of different variant models rated by invited participants. Each metric samples nearly 200 video ratings from nearly 30 participants, and the higher the score, the better the model performs.

**Quantitative Results** In the comparison of existing model experiments, the baseline was chosen to win second place [24] in the GENEA 2023 competition for comparison experiments and is easy for us to compare due to its open-source reproducibility. From the table above, our proposed method achieves a relatively good effect on human-likeness, with a score of $3.875 \pm 0.1481$ in the 95% confidence interval, which is significantly different from the baseline of $3.625 \pm 0.1816$. It is obvious that appropriateness for Agent Speech (AAS) and appropriateness for Interlocutor status (AIS) have improvement in comparison to the baseline.

**Ablation Study** We conduct ablation studies to demonstrate the effectiveness of each of the proposed components, as shown in the bottom of Table 1. The ablation studies utilize the symbol "-" as a shorthand for without. The experimental results shown in the table that the variation of without intention information has achieved the best effect among all experiments in the Appropriateness for Agent Speech (AAS), with a subjective score of $3.831 \pm 0.1510$, while only add the emotional information of the dialogue between the two parties to the conditions. This may indicate that emotional ups and downs have a certain implicit correlation with changes in gestures, for example, when happy, it may be accompanied by body and hand shaking, which vaguely reveals the happy mood of the speaker. On the other hand, under the variant model without emotional information, the matrix Appropriateness for Interlocutor status (AIS), which

involves the interlocutor, achieved the best performance with a score of $3.650 \pm 0.1487$ when the intention of dialogue between the two sides was taken into account. This is also very consistent with the characteristics of sociology, and the intention of dialogue between the two sides has a relatively important impact on the gestures of the two sides, so it may obtain relatively good results in the AIS matrix. Compared with the three indicators, it can be clearly seen that the AIS associated with the information of both sides has a low score, indicating that the current model needs to improve the modeling ability of the information of both sides.

**Visualization and Discussion** We visualized some gesture prediction results using the blender with python script from [25], as shown in Figure 2. We also strongly recommend the reader to visit our demo video[1] for better visualization. The first panel (a) of the figure shows the proposed approach, where an interaction is initiated from the main agent by pointing to the interlocutor for giving responses when the interlocutor is talking about himself, which is in line with our daily habits. In the second panel (b), the emotion-related variants also interact with each other over a wider range of movements while expressing content, as opposed to remaining quiet and even developing some joint problems in the two right panels (c-d). In the baseline, there is a more obvious hand back bend, which is not in line with our human habits.

## 4. CONCLUSION

Most previous works have focused on single-person gesture generation. In this paper, inspired by the work on social dyadic interactions, we propose a diffusion-based architecture to model the semantics, conversational intent, and emotions of both parties for high-quality, human-like, and appropriate conversational co-speech gesture generation. The results of the experiments demonstrate that taking into account the semantics of speech, intentions, and emotions of both participants in a conversational scene under the generation condition can enhance the system's performance, especially the method based on the diffusion model, capable of generating diverse and human-like gestures. While promising, there are still many questions and further research to be investigated. For instance, only the sociological characteristics of intention and emotion have been considered in our work. In fact, social knowledge, gender, and the surrounding environment have the potential to influence the gestures of both parties during their communication. Moreover, regarding design generation conditions, we are in the coarse-grained phase of global information extraction. It is possible to explore the application of named entity recognition to the semantics of a dialogue between two parties. These interesting problems are worth further research in the future.

---

[1]https://youtu.be/JHkyoI0qFNA

# 5. REFERENCES

[1] David Mcneill, "Hand and mind: What gestures reveal about thought," *Leonardo*, vol. 27, no. 4, 1992.

[2] Adam Kendon, "Movement coordination in social interaction: some examples described.," *Acta Psychol*, vol. 32, no. none, pp. 101–125, 1970.

[3] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, and Haolin Zhuang, "QPGesture: Quantization-based and phase-guided motion matching for natural speech-driven gesture generation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*. 2023, pp. 2321–2330, IEEE.

[4] Shenhan Qian, Zhi Tu, Yihao Zhi, Wen Liu, and Shenghua Gao, "Speech drives templates: Co-speech gesture synthesis with learned templates," in *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. 2021, pp. 11057–11066, IEEE.

[5] Sicheng Yang, Zhiyong Wu, Minglei Li, Zhensong Zhang, Lei Hao, Weihong Bao, Ming Cheng, and Long Xiao, "DiffuseStyleGesture: Stylized audio-driven co-speech gesture generation with diffusion models," in *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR, China*. 2023, pp. 5860–5868, ijcai.org.

[6] Youngwoo Yoon, Bok Cha, Joo-Haeng Lee, Minsu Jang, Jaeyeon Lee, Jaehong Kim, and Geehyuk Lee, "Speech gesture generation from the trimodal context of text, audio, and speaker identity," *ACM Trans. Graph.*, vol. 39, no. 6, pp. 222:1–222:16, 2020.

[7] Xingqun Qi, Chen Liu, Lincheng Li, Jie Hou, Haoran Xin, and Xin Yu, "Emotiongesture: Audio-driven diverse emotional co-speech 3d gesture generation," *CoRR*, vol. abs/2305.18891, 2023.

[8] Yonatan Shafir, Guy Tevet, Roy Kapon, and Amit H. Bermano, "Human motion diffusion as a generative prior," *CoRR*, vol. abs/2303.01418, 2023.

[9] Baptiste Chopin, Hao Tang, and Mohamed Daoudi, "Bipartite graph diffusion model for human interaction generation," *CoRR*, vol. abs/2301.10134, 2023.

[10] Evonne Ng, Hanbyul Joo, Liwen Hu, Hao Li, Trevor Darrell, Angjoo Kanazawa, and Shiry Ginosar, "Learning to listen: Modeling non-deterministic dyadic facial motion," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*. 2022, pp. 20363–20373, IEEE.

[11] Miles Patterson, "A systems model of dyadic nonverbal communication," *Journal of Nonverbal Behavior*, 2019.

[12] Alan Keltner, Jessica Cowen, "Emotional expression: Advances in basic emotion theory," *Journal of nonverbal behavior*, vol. 43, no. 2, 2019.

[13] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf, "Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter," *ArXiv*, vol. abs/1910.01108, 2019.

[14] Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H. Bermano, "Human motion diffusion model," *CoRR*, vol. abs/2209.14916, 2022.

[15] Lisa Feldman Barrett, Batja Mesquita, and Maria Gendron, "Context in emotion perception," *Current directions in psychological science*, vol. 20, no. 5, pp. 286–290, 2011.

[16] Andrea Scarantino, "How to do things with emotional expressions: The theory of affective pragmatics," *Psychological Inquiry*, vol. 28, no. 2-3, pp. 165–185, 2017.

[17] Agneta H Fischer and Disa A Sauter, "What the theory of affective pragmatics does and doesn't do," *Psychological Inquiry*, vol. 28, no. 2-3, pp. 190–193, 2017.

[18] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov, "Roberta: A robustly optimized bert pretraining approach," *arXiv preprint arXiv:1907.11692*, 2019.

[19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[20] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov, "Enriching word vectors with subword information," *Transactions of the association for computational linguistics*, vol. 5, pp. 135–146, 2017.

[21] Che-Jui Chang, Sen Zhang, and Mubbasir Kapadia, "The IVI lab entry to the GENEA challenge 2022 - A tacotron2 based method for co-speech gesture generation with locality-constraint attention mechanism," in *International Conference on Multimodal Interaction, ICMI 2022, Bengaluru, India, November 7-11, 2022*, Raj Tumuluri, Nicu Sebe, Gopal Pingali, Dinesh Babu Jayagopi, Abhinav Dhall, Richa Singh, Lisa Anthony, and Albert Ali Salah, Eds. 2022, pp. 784–789, ACM.

[22] Taras Kucherenko, Rajmund Nagy, Youngwoo Yoon, Jieyeon Woo, Teodor Nikolov, Mihail Tsakov, and Gustav Eje Henter, "The GENEA challenge 2023: A large scale evaluation of gesture generation models in monadic and dyadic settings," *CoRR*, vol. abs/2308.12646, 2023.

[23] Gilwoo Lee, Zhiwei Deng, Shugao Ma, Takaaki Shiratori, Siddhartha S. Srinivasa, and Yaser Sheikh, "Talking with hands 16.2m: A large-scale dataset of synchronized body-finger motion and audio for conversational motion analysis and synthesis," in *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019*. 2019, pp. 763–772, IEEE.

[24] Sicheng Yang, Haiwei Xue, Zhensong Zhang, Minglei Li, Zhiyong Wu, Xiaofei Wu, Songcen Xu, and Zonghong Dai, "The diffusestylegesture+ entry to the GENEA challenge 2023," *CoRR*, vol. abs/2308.13879, 2023.

[25] Youngwoo Yoon, Pieter Wolfert, Taras Kucherenko, Carla Viegas, Teodor Nikolov, Mihail Tsakov, and Gustav Eje Henter, "The GENEA challenge 2022: A large evaluation of data-driven co-speech gesture generation," in *International Conference on Multimodal Interaction, ICMI 2022, Bengaluru, India, November 7-11, 2022*, Raj Tumuluri, Nicu Sebe, Gopal Pingali, Dinesh Babu Jayagopi, Abhinav Dhall, Richa Singh, Lisa Anthony, and Albert Ali Salah, Eds. 2022, pp. 736–747, ACM.