

# Audio2Gestures: Generating Diverse Gestures from Speech Audio with Conditional Variational Autoencoders

Jing Li<sup>1</sup>, Di Kang<sup>2</sup>, Wenjie Pei<sup>1</sup>, Xuefei Zhe<sup>2</sup>, Ying Zhang<sup>2</sup>, Zhenyu He<sup>\*1</sup>, and Linchao Bao<sup>2</sup>

<sup>1</sup>Harbin Institute of Technology, Shenzhen

<sup>2</sup>Tencent AI Lab

## Abstract

Generating conversational gestures from speech audio is challenging due to the inherent one-to-many mapping between audio and body motions. Conventional CNNs/RNNs assume one-to-one mapping, and thus tend to predict the average of all possible target motions, resulting in plain/boring motions during inference. In order to overcome this problem, we propose a novel conditional variational autoencoder (VAE) that explicitly models one-to-many audio-to-motion mapping by splitting the cross-modal latent code into shared code and motion-specific code. The shared code mainly models the strong correlation between audio and motion (such as the synchronized audio and motion beats), while the motion-specific code captures diverse motion information independent of the audio. However, splitting the latent code into two parts poses training difficulties for the VAE model. A mapping network facilitating random sampling along with other techniques including relaxed motion loss, bicycle constraint, and diversity loss are designed to better train the VAE. Experiments on both 3D and 2D motion datasets verify that our method generates more realistic and diverse motions than state-of-the-art methods, quantitatively and qualitatively. Finally, we demonstrate that our method can be readily used to generate motion sequences with user-specified motion clips on the timeline. Code and more results are at <https://jingli513.github.io/audio2gestures>.

## 1. Introduction

In the real world, co-speech gestures help express oneself better, and in the virtual world, it makes a talking avatar act more vividly. Attracted by these merits, there has been a growing demand for generating realistic human motions for given audio clips recently. This problem is very challenging because of the complicated one-to-many relationship between audio and motion. A speaker may act different gestures when speaking the same words due to different mental and physical states.

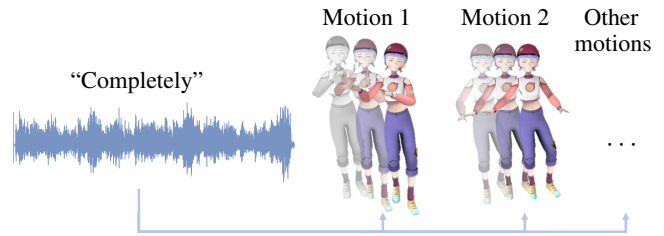


Figure 1. Illustration of the existence of one-to-many mapping between audio and motion in Trinity dataset [10]. Different gestures are performed when the subject says “completely”. Similar phenomena broadly exist in co-speech gestures. The character used for demonstration is from Mixamo [1].

Existing algorithms developed for audio to body dynamics have some obvious limitations. For example, [12] adapts a fully convolutional neural network to co-speech gesture synthesis tasks. Nevertheless, their model tends to predict averaged motion and thus generates motions lacking diversity. This is due to the underlying one-to-one mapping assumption of their model, which ignores that the relationship between speech and co-speech gesture is one-to-many in nature. Under such an overly simplified assumption, the model has no choice but to learn the averaged motion when several motions match almost the same audio clips in order to minimize the error. The above evidence inspires us to study whether or not explicitly modeling this multimodality improves the overall motion quality. To enhance the regression capability, we introduce an extra motion-specific latent code. With this varying *full* latent code, which contains the same shared code and varying motion-specific code, the decoder can regress different motion targets well for the same audio, achieving one-to-many mapping results. Under this formulation, the shared code extracted from audio input serves as part of the control signal. The motion-specific code further modulates the audio-controlled motion, enabling multimodal motion generation.

Although this formulation is straightforward, it is not trivial to make it work as expected. Firstly, there exists an easy degenerated solution since the motion decoder could utilize only the motion-specific code to reconstruct the motion. Secondly, we need to generate the motion-specific

\*Corresponding author: zhenyuhe@hit.edu.cn

code since we do not have access to the target motion during inference. Our solution to the aforementioned problems is providing *random noise* to the motion-specific code so that the decoder has to utilize the deterministic information contained in the shared code to reconstruct the target.

But under this circumstance, it is unsuitable for forcing the motion decoder to reconstruct the exact original target motion anymore. So a *relaxed motion loss* is proposed to apply to the motions generated with random motion-specific code. Specifically, it only penalizes the joints deviating from their targets larger than a threshold. This loss encourages the motion-specific code to tune the final motion while respecting the shared code’s control.

Our contributions can be summarized as:

- We present a co-speech gesture generation model whose latent space is split into shared code and motion-specific code to better regress the training data and generate diverse motions.
- We utilize random sampling and a relaxed motion loss to avoid degeneration of the proposed network and enable the model to generate multimodal motions.
- The effectiveness of the proposed method has been verified on 3D and 2D gesture generation tasks by comparing it with several state-of-the-art methods.
- The proposed method is suitable for motion synthesis from annotations since it can well respect the pre-defined actions in the timeline by simply using their corresponding motion-specific code.

## 2. Related Work

**Audio to body dynamics.** Early methods generate human motion for specified audio input by blending motion clips chosen from a motion database according to hidden Markov model [28] or conditional random fields [27]. Algorithms focusing on selecting motion candidates from a pre-processed database usually cannot generate motions out of the database and does not scale to large databases.

Recently, deep generative models, such as VAEs [23] and GANs [14], have achieved great success in generating realistic images, as well as human motions [39, 17, 29]. For example, [36] utilizes a classic LSTM to predict the body movements of a person playing the piano or violin given the sound of the instruments. However, the body movements of a person playing the piano or violin show regular cyclic pattern and are usually constrained within a small pose space.

In contrast, generating co-speech gestures is more challenging in the following two aspects – the motion to generate is more complicated and the relationship between the speech and motion is more complicated. As a result, Speech2Gesture [12] proposes a more powerful fully convolutional network, consisting of a 8-layer CNN audio encoder and a 16-layer 1D U-Net decoder, to translate log-mel audio feature to gestures. And this network is trained with

14.4 hours of data per individual on average in comparison to 3 hours data in [36]. Other than greatly enlarged network capacity, this fully convolutional network better avoids the error accumulation problem often faced by RNN-based methods. However, it still suffers from predicting the averaged motion due to the existence of one-to-many mapping in the training data. The authors further introduce adversarial loss and notice that the loss helps to improve diversity but degenerates the realism of the outputs. In contrast, our method avoids learning the averaged motion by explicitly modeling the one-to-many mapping between audio and motion with the help of the extra motion-specific code.

Due to lack of 3D human pose data, the above deep learning based methods [36, 12] have only tested 2D human pose data, which are 2D key point locations estimated from videos. Recently, [10] collects a 3D co-speech gesture dataset named Trinity Speech-Gesture Dataset, containing 244 minutes motion capture (MoCap) data with paired audio, and thus enables deep network-based study on modeling the correlation between audio and 3D motion. This dataset has been tested by StyleGestures [16], which is a flow-based algorithm [22, 16]. StyleGestures generates 3D gestures by sampling poses from a pose distribution predicted from previous motions and control signals. However, samples generated by flow-based methods [22, 16] are often not as good as VAEs and GANs. In contrast, our method learns the mapping between audio and motion with a customized VAE. Diverse results can be sampled since VAE is a probabilistic generation model.

**Human motion prediction.** There exist many works focus on predicting future motion given previous motion [17, 35, 39]. It is natural to model sequence data with RNNs [11, 19, 31, 39]. But [17] has pointed out the RNN-based methods often suffer from error accumulation and thus are not good at predicting long-term human motion. So they proposes to use a fully convolutional generative adversarial network and achieves better performance at long-term human motion prediction. Similarly, we also adopt a fully convolutional neural network since we need to generate long-term human motion. Specific to 3D human motion prediction, another type of error accumulation happens along the kinematic chain [35] because any *small* joint rotation error propagates to all its descendant joints, e.g. hands and fingers, resulting in *considerable* position error especially for the end-effectors (wrists, fingers). So QuaterNet [35] optimizes the joint position which is calculated from forward kinematics when predicting long-term motion. Differently, we optimize the joint rotation and position losses at the same time to help the model learn the joint limitation at the same time.

**Multimodal generation tasks.** Generating data with multimodality has received increasing interests in various tasks, such as image generation [18, 42], motion generation [38,

43]. For image generation, MUNIT [18] disentangles the embedding of images into content feature and style feature. BicycleGAN [42] combined cVAE-GAN [26] and cLR-GAN [6, 9] to encourage the bijective consistency between the latent code and the output so that the model could generate different output by sampling different codes. For video generation, MoCoGAN [38] and S3VAE [43] disentangle the motion from the object to generate videos in which different objects perform similar motions. Different from [38, 43], our method disentangle the motion representation into the audio-motion shared information and motion-specific information to model the one-to-many mapping between audio and motion.

### 3. Preliminaries

In this section, we first briefly introduce the variational autoencoder (VAE) [23], which is a widely used generative model. Then we describe 3D motion data and the most commonly used motion losses.

#### 3.1. Variational autoencoder

Compared to autoencoder, VAE additionally imposes constraints on the latent code to enable sampling outputs from the latent space. Specifically, during training, the distribution  $P$  of the latent code is constrained to match a target distribution  $Q$  with KL divergence as follows:

$$\mathcal{D}(Q(z) \| P(z|X)) = E_{z \sim Q} [\log Q(z) - \log P(z|X)], \quad (1)$$

where the  $X$  represents the input of the corresponding encoder (audio or motion in our case), and  $z$  represents its corresponding latent code. The above goal can be achieved by minimizing the Evidence Lower Bound (ELBO) [8]:

$$\log P(X|z) - \mathcal{D}[Q(z|X) \| P(z)]. \quad (2)$$

The second term of Eq. 2 is a KL-divergence between two Gaussian distributions (with a diagonal covariance matrix). The prior distribution  $P$  is set to Gaussian distribution (with a diagonal covariance matrix in our model, thus, the KL-divergence can be computed as:

$$\mathcal{D} = \frac{1}{2} \left( \text{tr}(\Sigma(X)) + \mu(X)^T \mu(X) - k - \log \det(\Sigma(X)) \right), \quad (3)$$

where  $k$  is the dimension of the distribution [8].

#### 3.2. Motion reconstruction loss

In our method, the generated motion is supervised with *motion reconstruction loss*, consisting of rotation loss, position loss, and speed loss. Formally, it is defined as follows:

$$L_{\text{mot}} = \lambda_{\text{rot}} \times L_{\text{rot}} + \lambda_{\text{pos}} \times L_{\text{pos}} + \lambda_{\text{speed}} \times L_{\text{speed}}, \quad (4)$$

where  $\lambda_{\text{rot}}$ ,  $\lambda_{\text{pos}}$ ,  $\lambda_{\text{speed}}$  are weights. We detail each term in the following.

Angular distance, i.e., geodesic distance, between the predicted rotation and the GT is adopted as the rotation loss. Mathematically,

$$L_{\text{rot}} = \frac{1}{J \times T} \sum_{j=1}^J \sum_{t=1}^T \cos^{-1} \frac{\text{Tr}(R_t^j (\hat{R}_t^j)^{-1}) - 1}{2}. \quad (5)$$

Position loss is the  $L_1$  distance between the predicted and target joint positions as follows:

$$L_{\text{pos}} = \frac{1}{J \times T} \sum_{j=1}^J \sum_{t=1}^T \|\hat{p}_t^j - p_t^j\|_1. \quad (6)$$

Speed loss is introduced to help the model learn the complicated motion dynamics. In our work, the joint speed  $v_t^j$  is defined as  $v_t^j = p_{t+1}^j - p_t^j$ . We optimize the predicted and target joint speed as follows:

$$L_{\text{speed}} = \frac{1}{J \times (T-1)} \sum_{j=1}^J \sum_{t=1}^{T-1} \|\hat{v}_t^j - v_t^j\|_1. \quad (7)$$

Our model can be trained with 2D motion data or 3D motion data. When modeling the 2D human motion, our method directly predicts the joint position. When modeling the 3D human motion, our method predicts the joint rotation and calculates the 3D joint positions with forward kinematics (FK). Concretely, the FK equation takes in as input the joint rotation matrix about its parent joint and the relative translation to its parent joint (i.e. bone length) and outputs joint positions as follows:

$$p_t^j = p_t^{\text{parent}(j)} + R_t^j s^j, \quad (8)$$

where  $R_t^j$  represents the rotation matrix of joint  $j$  in frame  $t$ ,  $p_t^j$  represents the position of joint  $j$  in frame  $t$ ,  $s^j$  represents the relative translation of joint  $j$  to its parent, and  $\text{parent}(j)$  represents the parent joint index of the joint  $j$ . We will always use  $j$  and  $t$  to index joints and frames in the following. Our model predicts joint rotation in 6D representation [41], which is a continuous representation that help the optimization of the model. The representation is then converted to rotation matrix  $R_t^j$  by Gram-Schmidt-like process, where  $R_t^j$  is the rotation matrix of joint  $j$  in frame  $t$ .

## 4. Audio2Gestures

The proposed Audio2Gestures algorithm is detailed in this section. We first present our Audio2Gestures network by formulating the multimodal motion generation problem in Sec. 4.1, then we detail the training process in Sec. 4.2.

### 4.1. Network structure

We use a conditional encoder-decoder network to model the correlation between audio  $A$  and motion  $M = [p_1, p_2, \dots, p_T]$ , where  $p_t$  represents the joint positions of frame  $t$ . In Fig. 2, our proposed model is made up of an audio encoder  $f_A$ , a motion encoder  $f_M$ , a mapping net  $f_R$  to produce motion-specific code during inference, and a common decoder  $g$  to generate motions from latent codes. The

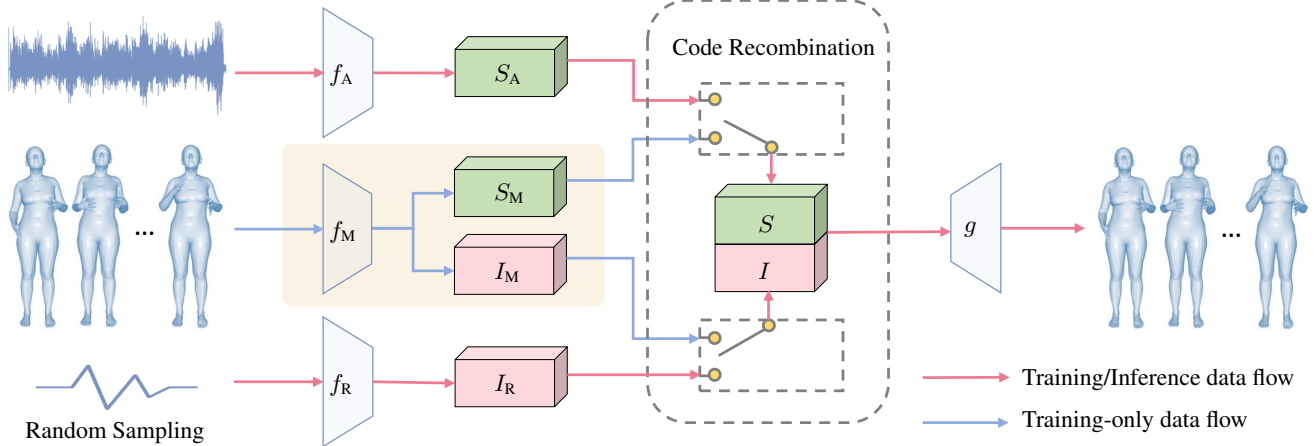


Figure 2. Our method explicitly models the audio-motion mapping by splitting the latent code into shared and motion-specific codes. The decoder generates different motions by recombining the shared and motion-specific codes extracted from different sources. The data flow in blue is only used at the training stage because we do not have motion data during inference.

latent code has been explicitly *split* into two parts (code  $S$  and  $I$ ) to account for the frequently occurred one-to-many mapping between the *same* (technically, very similar) audio and many different possible motions. The mapping network is introduced to facilitate sampling motion-specific codes. Under this formulation, given the same audio input (resulting in the same shared code  $S_A$ ), varied motions produce different motion-specific code  $I_M$  through motion encoder  $f_M$ , resulting in different *full* latent codes ( $S_A \oplus I_M$ ) so that the network can better model the one-to-many mapping ( $M = g(S_A, I_M)$ ).

During inference, shared feature  $S_A$  is extracted with  $f_A$  from the given audio  $A$ . Motion-specific feature  $I_R$  is generated with  $f_R$  from a randomly sampled signal. Both  $S_A$  and  $I_R$  are fed into the decoder  $g$  to produce the final motion  $M$ , i.e  $M = g(S_A, I_R)$ .

During training, given a paired audio-motion data  $A$  and  $M$ , their features  $S_A, S_M, I_M$  are firstly extracted by the encoders. Concretely,  $S_A = f_A(A)$  and  $(S_M, I_M) = f_M(M)$ . The decoder learns to reconstruct the input motion from the extracted features. To be specific, the decoder models the motion space by reconstructing the input motion by  $\hat{M} = g(S_M, I_M)$ . The model is expected to learn the joint embedding of audio and motion by guiding the decoder generate the same target motion from shared codes extracted from different source. But in practice, we notice the decoder will ignore the shared codes  $S$  and reconstruct the motion only from  $I_M$ . This is unwanted since the final motion is solely determined by the motion-specific features, being completely not correlated with the control signal (audio). Thus, another data flow ( $\hat{M}_{S_A I_R} = g(S_A, I_R)$ ) is introduced so that the decoder has to utilize the information contained in the shared code extracted from audio to reconstruct the target. The  $I_R$  is generated from the mapping net  $f_R$ , whose input is a random signal from a Gaussian distribution. The mean and variance of the distribution is cal-

culated from the  $I_M$  of the target motion per channel. We experimentally find using a mapping network  $f_R$  is helpful to improve the realism of the generated motions, which is mainly caused by the mapping network helps align the sampled feature with the motion-specific feature.

## 4.2. Latent code learning

To better learn the split audio-motion shared and motion-specific latent codes, five types of losses are introduced (Fig. 3). *Alignment constraint* and *relaxed motion loss* are introduced to learn the joint embedding (i.e., shared code) of the audio and motion. *Bicycle constraints* and *diversity loss* are introduced to model the multimodality of the motions. *KL divergence* has been described in Sec. 3 and thus omitted. The details are as follows.

**Shared code alignment.** The shared code of paired audio and motion is expected to be the same so that we can safely use audio-extracted shared code during inference and generate realistic and audio-related motions. We align the shared code of audio and motion by the alignment constraint:

$$L_{AC} = \|S_A - S_M\|_1. \quad (9)$$

**Degeneration avoidance.** As we described in Sec. 4.1, the model easily results in the degenerated network, which means the shared code is completely ignored and has no effect on the generated motion. Our solution to alleviating such degeneration is introducing an extra motion reconstruction with audio extracted shared code  $S_A$  and random motion-specific code  $I_R$ . Ideally, the generated motion  $\hat{M}_{S_A I_R}$  resembles its GT from some aspects but is not the same as its GT. In our case, We assume the generated poses are similar in the 3D world space. Thus we propose *relaxed motion loss*, which calculates the position loss and penalizes the model only when the distance is larger than a certain threshold  $\rho$ :

$$L_S = \frac{1}{J} \sum_{i=1}^J \max(\|\hat{p}_i - p_i\|_1 - \rho, 0). \quad (10)$$

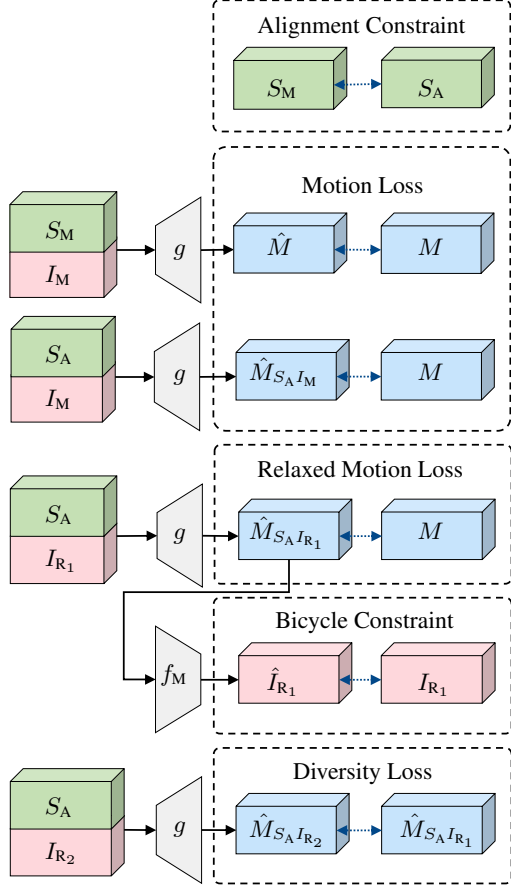


Figure 3. The training details of our model. Our model is trained with alignment constraint, motion reconstruction losses, relaxed motion loss, bicycle constraint, diversity loss and KL divergence. The alignment constraints and motion reconstruction loss help the model learn the audio-motion joint embedding. The relaxed motion loss avoids the degeneration of the shared code. The bicycle constraints and the diversity loss help reduce the mode-collapse problem and guide the model to generate multimodal motions. The KL divergence is omitted in the figure for the sake of brevity.

**Motion-specific code alignment.** Although the motion-specific code could be sampled from Gaussian distribution directly, we noticed that the realism and diversity of the generated motions are not good. The problem is caused by the misalignment of the Gaussian distribution and the motion-specific code distribution. Thus, the mapping net is introduced to map the signal sampled from Gaussian space to the motion-specific embedding. At the training stage, we calculate the mean and variance for every channel and every sample of the  $I_M$ . The sampled features will be fed into a mapping network, which is also a variational autoencoder, before concatenating them with different shared codes to generate motions.

**Motion-specific code reconstruction.** Although the model could model the multimodal distribution of audio-motion pair by splitting the motion code into audio-motion shared

one and motion-specific ones, it is not guaranteed the decoder can sample multimodal motions. For example, suppose the mapping net only maps the sampled signal to a single mode of the multimodal distribution. In that case, the decoder still could only generate unimodal motions, which is also known as the mode-collapse problem. The bicycle constraint [42] ( $M \rightarrow I \rightarrow M$  and  $I \rightarrow M \rightarrow I$ ) is introduced to avoid the mode-collapse problem, which encourages a bijection between the motion and the motion-specific code. Since the motion reconstruction loss has already been introduced, an extra reconstruction loss of the motion-specific code is added as supplement:

$$L_{\text{cyc}} = \|\hat{I}_R - I_R\|_1. \quad (11)$$

**Motion diversification.** To further encourage multimodality of the generated motion, diversity loss [30, 7] is introduced. Maximizing the multimodality loss encourages the mapping network to explore the meaningful motion-specific code space. We follow the setting in [7] and directly maximize the joint position distance between two sampled motions since it is more stable than the original one [30]:

$$L_{\text{DS}} = -L_{\text{pos}}(\hat{M}_{S_M I_{R_1}}, M). \quad (12)$$

## 5. Experiments

In this section, we first introduce the datasets, evaluation metrics and implementation details separately in Sec. 5.1-5.3. Then we show the performance of our algorithm and compare it with three state-of-the-art methods 5.4. Finally, we analyze the influence of each module of our model on the performance by ablation studies 5.5. More results are presented in our project page<sup>1</sup>.

### 5.1. Datasets

**Trinity dataset.** Trinity Gesture Dataset [10] is a large-scale speech to gesture synthesis dataset. This dataset records a male native English speaker talking many different topics, such as movies and daily activities. The dataset contains 23 sequences of paired audio-motion data, 244 minutes in total. The audio of the dataset is recorded at 44kHz. The motion data, consisting of 56 joints, are recorded at 60 frame per second (FPS) or 120 FPS using Vicon motion capture system.

**S2G-Ellen dataset.** The S2G-Ellen dataset, which is a subset of the Speech2Gesture dataset [12], contains positions of 49 2D upper body joint estimated from 504 YouTube videos, including 406 training sequences (469513 frames), 46 validation sequences (46027 frames), and 52 test sequences (59922 frames). The joints, which is estimated using OpenPose [5], include neck, shoulders, elbows, wrists, and hands.

<sup>1</sup><https://jingli513.github.io/audio2gestures>



## 5.2. Evaluation metrics

### 5.2.1 Quantitative metrics

**Realism.** Following Ginosar *et al.*'s [12] suggestion, the  $L_1$  distance of joint position in Eq. 13 and the percentage of correct 3D keypoints (PCK) in Eq. 14 are adopted to evaluate the realism of the generated motion. Specifically,  $L_1$  distance is calculated by averaging the corresponding joint's position error of all joints between prediction  $\hat{M}$  and GT  $M$ :

$$L_1 = \frac{1}{T \times J} \sum_{t=1}^T \sum_{j=1}^J \|\hat{M}_t^j - M_t^j\|_1. \quad (13)$$

The PCK metric calculates the percentage of correctly predicted keypoints, where a predicted keypoint is thought as correct if its distance to its target is smaller than a threshold  $\delta$ :

$$\text{PCK} = \frac{1}{T \times J} \sum_{t=1}^T \sum_{j=1}^J \mathbf{1}[\|p_t^j - \hat{p}_t^j\|_2 < \delta], \quad (14)$$

where  $\mathbf{1}$  is the indicator function and  $p_t^j$  indicates joint  $j$ 's position of frame  $t$ . As in [12], the  $\delta$  is set to 0.2 in our experiments.

**Diversity.** Diversity measures how many different poses/motions have been generated within a long motion. For example, RNN-based methods easily get stuck to some static motion as the generated motion becomes longer and longer. And static motions, which are undesired apparently, should get low diversity scores. We first split the generated motions into equal-lengthed non-overlapping motion clips (50 frames per clip in our experiments) and we calculate diversity as the averaged  $L_1$  distance of the motion clips. Formally, it is defined as:

$$\text{Diversity} = \frac{1}{N \times \lfloor N/2 \rfloor} \sum_{a=1}^N \sum_{a_2=a_1+1}^N \|\hat{M}_{a_1} - \hat{M}_{a_2}\|_1, \quad (15)$$

where the  $\hat{M}_{a_1}$  and  $\hat{M}_{a_2}$  represent clips from the same motion sequence,  $N$  represents the count of the motion clips, which is  $\frac{T}{50}$  in our experiments. Please note that jitter motion and invalid poses can also result in high diversity score. So higher diversity is preferred only if the generated motion is natural.

**Multimodality.** Multimodality measures how many different motions could be sampled (through multiple runs) for a given audio clip. Note that multimodality calculates motion difference across different motions while diversity calculates (short) motion clip difference within the same (long) motion. We measure the multimodality by generating motions for an audio  $N$  times, which is 20 in our experiments, and then calculate the average  $L_1$  distance of the motions.

$$\text{Multimodality} = \frac{1}{N \times \lfloor N/2 \rfloor} \sum_{a=1}^N \sum_{b=a+1}^N \|\hat{M}_a - \hat{M}_b\|_1, \quad (16)$$

where the  $\hat{M}_a$  and  $\hat{M}_b$  represent sampled motions generated through different runs for the given audio. Similar to diversity, invalid motion will also result in abnormally high multimodality score.

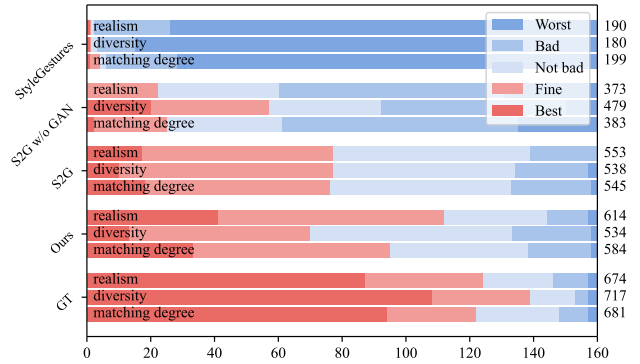


Figure 4. User study results comparing our method against the state-of-the-art methods. ‘‘S2G’’ is short for Speech2Gesture [12]. The horizontal axis represents the number of samples rated by the participants. In total, 160 comparisons have been rated (40 participants, 4 comparisons each questionnaire). The average score (higher is better) for each method is listed on the right. Bars with different colors indicate the count of the corresponding ranking of each algorithm. The video results are in our project page.

### 5.2.2 User studies

To evaluate the results qualitatively, we conduct user studies to analyze the visual quality of the generated motions. Our questionnaire contains four 20-second long videos. The motion clips shown in one video is generated by various methods from the same audio clip. The participants are asked to rate the motion clips from the following three aspects respectively:

1. Realism: which one is more realistic?
2. Diversity: which motion has more details?
3. Matching degree: which motion matches the audio better?

The results of the questionnaires are shown in Fig. 4. We show the count of different ranking in the figure. The average score of different metrics for each algorithm is listed after the corresponding bar. The scores assigned to each ratings are  $\{5,4,3,2,1\}$  for  $\{\text{best, fine, not bad, bad, worst}\}$  respectively.

## 5.3. Implementation details

**Data processing.** We detail the data processing of Trinity dataset and S2G-Allen dataset here.

(1) Trinity dataset. The audio data are resampled to 16kHz for extracting log-mel spectrogram [37] feature using librosa [32]. More concretely, the hop size is set to SR/FR where SR is the sample rate of the audio and FR is the frame rate of the motion so that the resulting audio feature have the same length as the input motion. In our case, the resulting hop size is 533 since SR is 16000 and FR is 30. The dimension of the log-mel spectrogram is 64.

The motion data are downsampled to 30 FPS and then retargeted to the SMPL-X [34] model. SMPL-X is an expressive articulated human model consisting of 54 joints

Dataset	Method	$L_1$ ↓	PCK ↑	Diversity ↑	Multimodality ↑
Trinity	S2G w/o GAN [12]	7.71	0.82	5.99	-
	S2G [12]	24.68	0.39	2.46	-
	StyleGestures [3]	18.97 (18.07)	0.34 (0.34)	2.34 (3.79)	<b>7.55</b>
	Ours	<b>7.84 (7.65)</b>	<b>0.82 (0.83)</b>	<b>6.32 (6.52)</b>	4.11
S2G-Ellen	S2G w/o GAN [12]	0.74	0.37	0.61	-
	S2G [12]	1.08	0.23	0.89	-
	Ours	0.94 (0.92)	0.33 (0.34)	0.84 (0.85)	0.77

Table 1. Quantitative results on Trinity dataset and S2G-Ellen dataset. ↑ means the higher is better and ↓ means the lower is better. For methods supporting sampling, we run 20 tests and report their *average* score and the *best* score (in parentheses). Speech2Gesture (“S2G” in the table) could not generate multimodality motions.

Method	$L_1$ ↓	PCK ↑	Diversity ↑	Multimodality ↑
baseline	8.22	0.80	6.20	-
+ split	8.69 (8.30)	0.77 (0.78)	5.83 (6.02)	<b>5.90</b>
+ mapping net	8.06 (7.91)	0.80 (0.81)	5.86 (6.05)	3.44
+ bicycle constraint	7.94 ( <b>7.63</b> )	0.80 (0.82)	6.31 (6.46)	3.68
+ diversity loss	<b>7.84 (7.65)</b>	<b>0.82 (0.83)</b>	<b>6.32 (6.52)</b>	4.11

Table 2. Ablation study results on the Trinity dataset. Note that every line adds a new component compared to its previous line. For methods supporting sampling, we run 20 tests and report their *average* score and the *best* score (in parentheses).

(21 body joints, 30 hand joints, 3 face joints, respectively), which has been widely used in 3D pose estimation and prediction [20, 34, 40, 25]. The joint rotation is in 6D rotation representation [41] in our experiments, which is a smooth representation and could help the model approximate the target easier. Note that the finger motions are removed due to unignorable noise.

(2) S2G-Ellen dataset. Following [12], the data are split into 64-frame long clips (4.2 seconds). Audio features are extracted in the same way as the Trinity dataset. The body joints are represented in a local coordinate frame relative to its root. Namely, the origin of the coordinate is the root joint.

**Network.** Every encoder, decoder and mapping net consists of four residual blocks [15], including 1D convolution and ReLU non-linearity [2]. The residual block is similar to [4] except several modifications. To be specific, the casual convolutions whose kernels see only the history are replaced with normal symmetric 1D convolutions seeing both the history and the future. Both the shared code and motion-specific code are set to 16 dimensions.

**Training.** At the training stage, we randomly crop a 4.2-second segment of the audio and motion data, which is 64 frames for the S2G dataset (15 FPS) and 128 frames for the Trinity dataset (30 FPS). The model weights are initialized with the Xavier method [13] and trained 180K steps using the Adam [21] optimizer. The batch size is 32 and the learning rate is  $10^{-4}$ . The  $\lambda_{rot}$ ,  $\lambda_{pos}$ ,  $\lambda_{speed}$  are set as 1, 1, 5 respectively, and  $\rho$  is set as 0.02 in our experiments. Our model is implemented with PyTorch [33].

#### 5.4. Comparison with state-of-the-art methods

We compare our method with two recent representative state-of-the-art methods, including one LSTM-based

method named StyleGestures [3] and one CNN-based method named Speech2Gesture [12] on Trinity dataset. StyleGestures adapts normalizing flows [24, 22, 16] to speech-driven gesture synthesis. We train StyleGestures using the code released by the authors. The training data of the StyleGestures are processed in the same way as the authors indicate<sup>2</sup>. Speech2Gesture, originally designed to map speech to 2D human keypoints, consists of an audio encoder and a motion decoder. Its final output layer has been adjusted to predict 3D joint rotations and is trained with the same losses as our method.

Quantitative experimental results are listed in Tab. 1 and user study results in Fig. 4. Both results show that our method outperforms previous state-of-the-art algorithms on the realism and diversity metrics, demonstrating that it is beneficial to explicitly model the one-to-many mapping between audio and motion in the network structure.

While StyleGestures supports generating different motions for the same audio by sampling, the quality of its generated motions is not very appealing. Also, its diversity score is the lowest, because LSTM output easily gets stuck into some poses, resulting in long static motion afterwards. The algorithm is not good at generating long-term sequences due to the error accumulation problem of the LSTM. The authors test their algorithm on 400 frames (13 seconds) length sequences. However, obviously deteriorated motions are generated when evaluating their algorithm to generate 5000-frame (166 seconds) long motions.

As for Speech2Gesture, the generated motions show similar realism with ours but obtain lower diversity score (Tab. 1) than our method. But Speech2Gesture does not support generating multimodal motions. Also note that

<sup>2</sup>The motions generated by StyleGestures are 20 FPS and have a different skeleton from our method. We upsample the predicted motion to 30 FPS and retarget it to SMPL-X skeleton with MotionBuilder.

Speech2Gesture with GAN generates many invalid poses and gets the worst performance. We have trained the model several times changing the learning rate range from 0.0001 to 0.01, and report the best performance here. The bad performance may be caused by the unstable of the training process of the generative adversarial network.

### 5.5. Ablation study

To gain more insights into the proposed components of our model, we test some variants of our model on the 3D Trinity dataset (Tab. 2). We run every variant 20 times and report the *averaged* performance and the *best* performance to avoid the influence of randomness. Note that the randomness of our model comes from two different parts, the randomness introduced by the variational autoencoder and by the motion-specific feature sampling.

We start with a “baseline” model, which excludes the mapping net and the split code. It is trained only with the motion reconstruction losses (Eq. 4) and shared code constraint (Eq. 9). The averaged scores “avg  $L_1$ ”, “avg PCK” and “avg Diversity” of the model equal to the best scores “min  $L_1$ ”, “max PCK” and “max Diversity” on  $L_1$ , which indicates that the randomness of the VAE model have almost no affect on generating multimodal motions.

The next setting is termed as “+split”, which splits the output of the motion encoder into shared and motion-specific codes and introduces the relaxed motion loss (Eq. 10). This modification explicitly enables the network to handle the one-to-many mapping, but it harms the realism (see “avg  $L_1$ ”, “min  $L_1$ ”, “avg PCK” and “max PCK”) and diversity. As we can see, both the  $L_1$  and the PCK metrics are worse than “baseline”. The abnormal results is mainly caused by the misalignment between the sampled signal and the motion-specific feature. We analyze the difference between the sampled signals with the motion-specific feature, and find that there is a big difference in their statistical characteristics, such as the mean and variance of derivatives.

Thus, a mapping network (“+mapping net”) is introduced to align the sampled signal with the motion-specific feature automatically. Although the multimodality drops compare to “+split”, this modification helps to improve other metrics of the generated motions a lot. Note that a higher multimodality score only makes sense when the generated motions is natural, as described in Sec. 5.2. The “+mapping net” model also outperforms the baseline model in the  $L_1$  metrics and gets a similar PCK metrics, but the model could generate multimodal motions. We notice that the diversity of the motions generated by “+mapping net” model is not as good as the baseline model. The realism score of the “+mapping net” model is also worse than the baseline model, which may be due to the users prefer the motions with more dynamics. We think the problem may be caused by the mode collapse problem suffered by many

generative methods.

To overcome this problem, two simple yet effective losses – Bicycle constraints and diversity loss – are introduced. Bicycle constraint improves the multimodality of the motions from 3.44 to 3.68. The avg diversity of the motions also increase from 5.86 to 6.31. The diversity loss further improves the motion diversity and multimodality but have little influence on the realism. The final model outperforms the baseline model in all quantitative indicators, which shows that the audio-motion mapping could be better modeled by explicit modeling the one-to-many correlation.

### 5.6. Application

We notice that motion-specific code extracted from a motion strongly controls the final motion output. To be specific, the synthesized motion is almost the same as the original motion used to extract this motion-specific code. This feature is perfect for a type of motion synthesis application where pre-defined motions are provided on the timeline as constraints. For example, if there is a  $n$ -frame long motion clip that we want the avatar to perform from frame  $t$  to  $t + n$ . We could extract its motion-specific code  $I_M$  with the motion encoder and directly replace the sampled motion-specific code  $I_R$  from  $t$  to  $t + n$ . Our model could generate a smooth motion from the edited motion-specific code. Please refer to our project page for the demonstration.

## 6. Conclusion

In this paper, we explicitly model the one-to-many mapping by splitting the latent code into shared code and motion-specific code. This simple solution with our customized training strategy effectively improves the realism, diversity, and multimodality of the generated motion. We also demonstrate an application that the model could insert a specific motion into the generated motion by editing the motion-specific code, with smooth and realistic transitions. Despite the model could generate multimodal motions and provide users the ability to control the output motion, there exist some limitations. For example, the generated motion is not very related to what the person says, future work could be improving the meaning of the generated motion by incorporating word embedding as an additional condition.

## 7. Acknowledgement

This work was supported by the Natural Science Foundation of China (U2013210, 62006060), the Shenzhen Research Council (JCYJ20210324120202006), the Shenzhen Stable Support Plan Fund for Universities (GXWD20201230155427003-20200824125730001), and the Special Research project on COVID-19 Prevention and Control of Guangdong Province (2020KZDZDX1227).



## References

- [1] Mixamo. <https://www.mixamo.com>. [Online; accessed 15-March-2021].
- [2] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.
- [3] Simon Alexanderson, Gustav Eje Henter, Taras Kucherenko, and Jonas Beskow. Style-controllable speech-driven gesture synthesis using normalising flows. *Comput. Graph. Forum*, 39(2):487–496, 2020.
- [4] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- [5] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7291–7299, 2017.
- [6] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Adv. Neural Inform. Process. Syst.*, pages 2172–2180, 2016.
- [7] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2020.
- [8] CARL DOERSCH. Tutorial on variational autoencoders. *stat*, 1050:13, 2016.
- [9] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- [10] Ylva Ferstl and Rachel McDonnell. Iva: Investigating the use of recurrent motion modelling for speech gesture generation. In *IVA '18 Proceedings of the 18th International Conference on Intelligent Virtual Agents*, Nov 2018.
- [11] Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *Int. Conf. Comput. Vis.*, pages 4346–4354, 2015.
- [12] S. Ginosar, A. Bar, G. Kohavi, C. Chan, A. Owens, and J. Malik. Learning individual styles of conversational gesture. In *IEEE Conf. Comput. Vis. Pattern Recog.* IEEE, June 2019.
- [13] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- [14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Adv. Neural Inform. Process. Syst.*, pages 2672–2680, 2014.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016.
- [16] Gustav Eje Henter, Simon Alexanderson, and Jonas Beskow. Moglow: Probabilistic and controllable motion synthesis using normalising flows. *ACM Trans. Graph.*, 39(6), Nov. 2020.
- [17] Alejandro Hernandez, Jurgen Gall, and Francesc Moreno-Noguer. Human motion prediction via spatio-temporal inpainting. In *Int. Conf. Comput. Vis.*, pages 7134–7143, 2019.
- [18] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *Eur. Conf. Comput. Vis.*, pages 172–189, 2018.
- [19] Ashesh Jain, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Structural-rnn: Deep learning on spatio-temporal graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5308–5317, 2016.
- [20] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-End Recovery of Human Shape and Pose. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 7122–7131, 2018.
- [21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [22] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Adv. Neural Inform. Process. Syst.*, pages 10215–10224, 2018.
- [23] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *stat*, 1050:1, 2014.
- [24] Ivan Kobyzev, Simon Prince, and Marcus Brubaker. Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [25] Nikos Kolotouros, Georgios Pavlakos, Michael Black, and Kostas Daniilidis. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 2019-October, pages 2252–2261, 2019.
- [26] Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, and Ole Winther. Autoencoding beyond pixels using a learned similarity metric. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, ICML'16*, page 1558–1566. JMLR.org, 2016.
- [27] Sergey Levine, Philipp Krähenbühl, Sebastian Thrun, and Vladlen Koltun. Gesture controllers. In *ACM SIGGRAPH 2010 papers*, pages 1–11. 2010.
- [28] Sergey Levine, Christian Theobalt, and Vladlen Koltun. Real-Time Prosody-Driven Synthesis of Body Language. *ACM Trans. Graph.*, 28(5):1–10, dec 2009.
- [29] Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. Character controllers using motion vaes. *ACM Trans. Graph.*, 39(4), July 2020.
- [30] Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [31] Julieta Martinez, Michael J Black, and Javier Romero. On human motion prediction using recurrent neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2891–2900, 2017.
- [32] Brian McFee, Vincent Lostanlen, Alexandros Metsai, Matt McVicar, Stefan Balke, Carl Thomé, Colin Raffel, Frank

- Zalkow, Ayoub Malek, Dana, Kyungyun Lee, Oriol Nieto, Jack Mason, Dan Ellis, Eric Battenberg, Scott Seyfarth, Ryuichi Yamamoto, Keunwoo Choi, viktorandreevichmorozov, Josh Moore, Rachel Bittner, Shunsuke Hidaka, Ziyao Wei, nullmightybofo, Darío Hereñú, Fabian-Robert Stöter, Pius Friesch, Adam Weiss, Matt Vollrath, and Taewoon Kim. *librosa/librosa*: 0.8.0, July 2020.
- [33] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [34] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [35] Dario Pavlo, David Grangier, and Michael Auli. Quaternet: A quaternion-based recurrent model for human motion. In *Brit. Mach. Vis. Conf.*, 2018.
- [36] Eli Shlizerman, Lucio Dery, Hayden Schoen, and Ira Kemelmacher-Shlizerman. Audio to body dynamics. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7574–7583, 2018.
- [37] Stanley Smith Stevens, John Volkman, and Edwin B Newman. A scale for the measurement of the psychological magnitude pitch. *The Journal of the Acoustical Society of America*, 8(3):185–190, 1937.
- [38] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. Mocogan: Decomposing motion and content for video generation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1526–1535, 2018.
- [39] Xinchun Yan, Akash Rastogi, Ruben Villegas, Kalyan Sunkavalli, Eli Shechtman, Sunil Hadap, Ersin Yumer, and Honglak Lee. Mt-vae: Learning motion transformations to generate multimodal human dynamics. In *Eur. Conf. Comput. Vis.*, pages 265–281, 2018.
- [40] Jason Zhang, Panna Felsen, Angjoo Kanazawa, and Jitendra Malik. Predicting 3D human dynamics from video. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 2019-Octob, pages 7113–7122, 2019.
- [41] Yi Zhou, Connelly Barnes, Lu Jingwan, Yang Jimei, and Li Hao. On the continuity of rotation representations in neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2019.
- [42] Jun-Yan Zhu, Richard Zhang, Deepak Pathak, Trevor Darrell, Alexei A Efros, Oliver Wang, and Eli Shechtman. Toward multimodal image-to-image translation. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 465–476. Curran Associates, Inc., 2017.
- [43] Yizhe Zhu, Martin Renqiang Min, Asim Kadav, and Hans Peter Graf. S3vae: Self-supervised sequential vae for representation disentanglement and data generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6538–6547, 2020.