
Vector Quantized Diffusion Model with CodeUnet for Text-to-Sign Pose Sequences Generation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Sign Language Production (SLP) aims to translate spoken languages into sign
2 sequences automatically. The core process of SLP is to transform sign gloss
3 sequences into their corresponding sign pose sequences (G2P). Most existing G2P
4 models usually perform this conditional long-range generation in an autoregressive
5 manner, which inevitably leads to an accumulation of errors. To address this issue,
6 we propose a vector quantized diffusion method for conditional pose sequences
7 generation, called PoseVQ-Diffusion, which is an iterative non-autoregressive
8 method. Specifically, we first introduce a vector quantized variational autoencoder
9 (Pose-VQVAE) model to represent a pose sequence as a sequence of latent codes.
10 Then we model the latent discrete space by an extension of the recently developed
11 diffusion architecture. To better leverage the spatial-temporal information, we
12 introduce a novel architecture, namely CodeUnet, to generate higher quality pose
13 sequence in the discrete space. Moreover, taking advantage of the learned codes,
14 we develop a novel sequential k-nearest-neighbours method to predict the variable
15 lengths of pose sequences for corresponding gloss sequences. Consequently,
16 compared with the autoregressive G2P models, our model has a faster sampling
17 speed and produces significantly better results. Compared with previous non-
18 autoregressive G2P methods, PoseVQ-Diffusion improves the predicted results with
19 iterative refinements, thus achieving state-of-the-art results on the SLP evaluation
20 benchmark.

21 1 Introduction

22 Sign Language Production (SLP), as an essential task for the Deaf community, aims to provide
23 continuously sign videos for spoken language sentences. Since sign languages are distinct linguistic
24 systems [1] which differ from natural languages, sign languages have different word orders from their
25 corresponding natural languages. Therefore, directly learning the alignment mapping between them
26 is challenging. To tackle this issue, previous works first translate spoken languages into glosses¹, then
27 generate the sign pose sequences based on the gloss sequences (G2P) [2, 3], and finally optionally
28 use the sign pose sequence to generate the photo-realistic sign video [4]. Accordingly, G2P is the
29 heart procedure of this task, and it is the focus of this paper.

30 Existing approaches for G2P can be categorized into autoregressive [2, 3] and non-autoregressive [5]
31 methods depending on their decoding strategies. Autoregressive models [2, 3] generate the next pose
32 frame depending on previous frames relying on the teacher forcing strategy [6]. In inference, the
33 recurrent decoding is likely to lead to prediction error propagation over time due to the exposure
34 bias [7]. To break the bottleneck of autoregression, non-autoregressive methods are proposed to

¹Sign glosses are spoken language words that match the meaning of signs and, linguistically, manifest as minimal lexical items.

35 induce the decoder to generate all target predictions simultaneously [8, 9]. Huang *et al.* [5] proposed
36 a non-autoregressive G2P model to generate sign pose sequence parallelly in a one-shot decoding
37 scheme, and used an External Aligner (EA) for sequence alignment learning.

38 Motivated by the recent developed Discrete Denoising Diffusion Probabilistic Model (D3PMs) [10,
39 11, 12] which achieved impressive results for language generation and vector quantized image
40 generation. We propose a Pose Vector Quantized Diffusion (PoseVQ-Diffusion) model to learn the
41 conditional pose sequence generation in the latent discrete space instead of the continuous coordinate
42 space. It is also a non-autoregressive method that performs parallel refinement on the generated
43 results with iterations and therefore shows expressive generative capacity.

44 We will elaborate our approach in three steps. Firstly, we utilize a vector quantized variational
45 autoencoder (VQ-VAE) to represent the pose sequence as sequential latent codes. Different from
46 image VQ-VAE [13, 14], we devise a specific architecture, Pose-VQVAE, to learn the meaningful
47 codebook by reconstructing the pose sequence. Specifically, we divide a sign skeleton into three local
48 point patches representing pose, right hand and left hand separately. Furthermore, a Tokenizer with a
49 vector quantized variational autoencoder is designed to learn discrete point tokens containing local
50 semantic information.

51 Next, we extend the standard vector quantized diffusion methods [11, 12] to model the sequential
52 alignments between sign glosses and quantized codes. The discrete diffusion model samples the data
53 distribution by reversing a forward diffusion process that gradually corrupts the input via a fixed
54 Markov chain. Its corruption process by adding noise data (*e.g.*, [MASK] token) draws our attention
55 to the mask-based generative model, Mask-Predict [9], which is proved to be a variant of diffusion
56 model [11]. In this paper, we explore two variants of the diffusion model for our quantized pose
57 sequence generation. Expanding further, to better leverage the spatial and temporal information of the
58 quantized pose sequences, we introduce a new architecture CodeUnet. In contrast to Unet [15] which
59 is a “fully convolution network” for image data, CodeUnet is a “fully transformer network” designed
60 for discrete tokens. As a result of iterative refinements and better spatial-temporal modelling, our
61 model achieves a higher quality of the conditional pose sequence generation.

62 Lastly, the length prediction of the non-autoregressive G2P models is challenging since the corre-
63 sponding lengths of different sign glosses are different and variable. In this paper, we propose a novel
64 clustering method for this typical sequential data that local adjacent frames should belong to a cluster.
65 Specifically, taking advantage of the meaningful learned codes in the first stage, we firstly apply the
66 k-nearest-neighbor based density peaks clustering algorithm [16, 17] to locate the peaks with higher
67 local density. Secondly, we design a heuristic algorithm to find the boundary between two peaks
68 according to their semantic distance with the two peak codes. Finally, we leverage the length of each
69 gloss as the additional supervised information and predict the length of the gloss sequence in the
70 inference.

71 Our model significantly improves the generation quality on the challenging RWTH-PHOENIX-
72 WEATHER-2014T [18] dataset. The evaluation of conditional sequential generation is evaluated
73 using a back-translated model. Extensive experiments show that our model increases the WER score
74 from 82.01% [5] to 78.21% on generated pose sequence to gloss sequence, and BLEU score from
75 6.66 [5] to 7.42 on generated pose sequence to spoken language.

76 2 Related Works

77 **Sign Language Production.** Most sign language works focus on sign language recognition (SLR)
78 and translation (SLT) [18, 19, 20, 21, 22, 23], aiming to translate the video-based sign language
79 into text-based sequences. And few attempts have been made for the more challenging task of sign
80 language production (SLP) [24, 25]. Stoll *et al.* proposed the first deep SLP model, which adopts the
81 three-step pipeline. In the core process for G2P, they learn the mapping between the sign glosses and
82 the skeleton poses via a look-up table. After that, B. Saunders *et al.* [3] proposed the progressive
83 transformer to learn the mapping with an encoder-decoder architecture and generate the sign pose
84 in an autoregressive manner in the inference. Further, B. Saunders *et al.* [4] proposed a Mixture
85 Density Network (MDN) to generate the pose sequences condition on the sign glosses and utilize a
86 GAN-based method [26] to produce the photo-realistic sign language video. **B. Saunders *et al.* [27]
87 translates the spoken language to sign language representation with an autoregressive Transformer
88 network, and uses the gloss information to provide the additional supervision. Then they propose a**

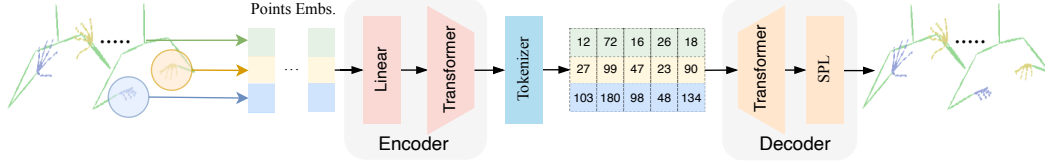


Figure 1: The architecture of the first stage model Pose-VQVAE for learning the discrete latent codes.

89 **Mixture of Motion Primitives(MoMP) architecture to combine distinct motion primitives to produce**
 90 **a continuous sign language sequence.**

91 Different from these autoregressive methods, Huang *et al.* [5] proposed a non-autoregressive model
 92 to parallelly generate the sign pose sequence avoiding the error accumulation problem. They apply
 93 the monotonic alignment search [28] to generate the alignment lengths of each gloss. Our model
 94 also explores a non-autoregressive method with a diffusion strategy, and the adopted diffusion model
 95 architecture provides us with a chance to refine the results with multiple iterations.

96 **Denosing Diffusion Probabilistic Models.** Diffusion generative models have achieved outstanding
 97 results on continuous data, such as image generation [29, 30, 31, 32, 33] and speech synthesis [34, 35,
 98 36]. However, most previous works focus on Gaussian diffusion processes that operate in continuous
 99 state spaces. The discrete diffusion model is first introduced in [37], and it is applied to text generation
 100 in Argmax Flow [10]. To improve and extend the discrete diffusion model, D3PM [11] use a structured
 101 categorical corruption process to shape data generation and embed structure in the forward process.
 102 VQ-Diffusion [12] apply the discrete diffusion model to conditional vector quantized image synthesis
 103 with a mask-and-replace diffusion strategy.

104 3 The Proposed Method

105 The overall objective of this work is to extend the discrete diffusion model for conditional sign pose
 106 sequence generation. The proposed PoseVQ-Diffusion model consists of three key components, Pose-
 107 VQVAE to learn the latent codes, a discrete diffusion model with CodeUnet to model the discrete
 108 codes generation, and a sequential-KNN algorithm on the length prediction for this non-autoregressive
 109 method.

110 3.1 Pose VQ-VAE

111 In this section, we introduce how to tokenize the points of a sign pose skeleton into a set of discrete
 112 tokens. A naive approach is to treat per point as one token. However, such a points-wise reconstruction
 113 model tends to tremendous computational cost due to the quadratic complexity of self-attention
 114 in Transformers. On the other hand, since the details of hand points are essential for sign pose
 115 understanding, treating all the points into one token leads to remarkable inferior reconstruction
 116 performance. To achieve a better trade-off between quality and speed, we propose a simple yet
 117 efficient implementation that groups the points of a sign skeleton into three local patches, representing
 118 pose, right hand and left hand separately. Figure 1 illustrates the framework of our proposed Pose-
 119 VQVAE model with the following submodules.

120 **Encoder.** Given a sign pose sequence of N frames $\mathbf{s} = (s_1, s_2, \dots, s_n, \dots, s_N) \in \mathbb{R}^{N \times J \times K}$, where
 121 $\{x_n^j\}_{j=1}^J$ presents a single sign skeleton containing J joints and K denotes the feature dimension
 122 for human joint data. We separate these points into three local patches, $\mathbf{s}_p \in \mathbb{R}^{N \times (J_p \times K)}$, $\mathbf{s}_r \in$
 123 $\mathbb{R}^{N \times (J_r \times K)}$, $\mathbf{s}_l \in \mathbb{R}^{N \times (J_l \times K)}$ for pose, right hand and left hand respectively, where $J = J_p + J_r + J_l$.
 124 In the encoder module $E(e|\mathbf{s})$, we first transform these three points sequences into feature sequences
 125 by simple three linear layers and concatenate them together. Then we apply a spatial-temporal
 126 Transformer network to learn the long-range interactions within the sequential point features. Finally,
 127 we arrive at the encoded features $\{e_n \in \mathbb{R}^{3 \times h}\}_{n=1}^N$.

128 **Point Tokenizer.** Similar to image VQ-VAE [14], we take the encoded features as inputs and convert
 129 them into discrete tokens. Specifically, we perform the nearest neighbors method $\mathcal{Q}(z|e)$ to quantize
 130 the point feature to the quantized features $\{z_n \in \mathbb{R}^{3 \times h}\}_{n=1}^N$. The quantized features are maintained
 131 by a codebook whose size is V .

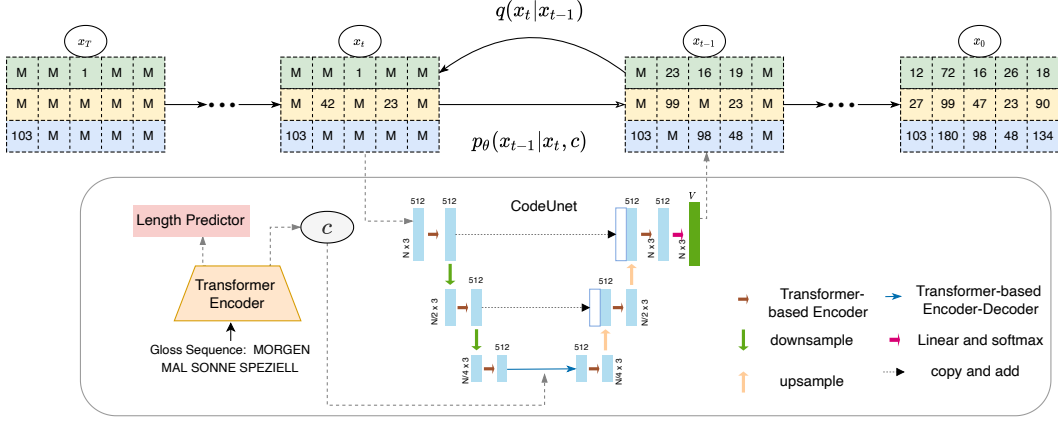


Figure 2: Our approach uses a discrete diffusion model to represent the Vector-Quantized sign pose sequence allowing non-autoregressive pose sequence generation. Specifically, after compressing the sign pose sequences to meaningful discrete codes, each code is randomly masked or replaced and a CodeUnet model is trained to restore the original data.

132 **Decoder.** The decoder $D(\tilde{s}|z)$ receives the quantized features as inputs and also applies spatial-
 133 temporal Transformer to get the output features $\{o_n \in \mathbb{R}^{3 \times h}\}_{n=1}^N$. Finally, we separate the output
 134 feature for three sub-skeleton and utilize a structured prediction layer (SPL) [38] $\mathcal{P}(\tilde{s}|o)$ to reconstruct
 135 the corresponding sub-skeleton $\tilde{s}_p \in \mathbb{R}^{N \times (J_r \times K)}$, $\tilde{s}_l \in \mathbb{R}^{N \times (J_r \times K)}$, $\tilde{s}_r \in \mathbb{R}^{N \times (J_r \times K)}$. We adopt
 136 the SPL to rebuild the skeleton from feature because it explicitly model the spatial structure of the
 137 human skeleton and the spatial dependencies between joints. The hierarchy chains of the pose, right
 138 hand and left hand skeleton are given in supplemental material.

139 **Training.** The encoder $E(e|s)$, tokenizer $Q(z|e)$ and decoder $D(\tilde{s}|z)$ can be trained end-to-end via
 140 the following loss function:

$$\mathcal{L}_{\text{Pose-VQVAE}} = \|s_p - \tilde{s}_p\| + \|s_r - \tilde{s}_r\| + \|s_l - \tilde{s}_l\| + \|sg[e] - z\| + \beta \|sg[z] - e\|, \quad (1)$$

141 where $sg[\cdot]$ stands for stop-gradient operation. **In practice, we replace the forth term with exponential**
 142 **moving averages (EMA) to update the codebook.**

143 3.2 Discrete Diffusion Model with CodeUnet

144 To allow conditional sampling, a discrete diffusion model is trained on the latent codes obtained from
 145 the Pose-VQVAE model. Figure 2 shows the architecture of our proposed PoseVQ-Diffusion, which
 146 aims to model the latent space in a non-autoregressive manner. We will subsequently introduce the
 147 diffusion process, reverse denoising process and the parametered model CodeUnet.

148 **Diffusion Process.** Given a sequence of latent codes $x_0 \in \mathbb{R}^{N \times 3}$ obtained from the vector quantized
 149 model, where $x_0^{(i,j)} \in \{1, 2, \dots, V\}$ at location (i, j) represents the index within the codebook. The
 150 diffusion process aims to corrupt the original data x_0 via a fixed Markov chain $p(x_t|x_{t-1})$ by adding
 151 small amount noise continuously. After a fixed T timesteps, it produces a sequence of increasingly
 152 noisy data x_1, \dots, x_T with the same dimensions as x_0 , and x_T becomes a pure noise sample.

153 For the scalar discrete variables with V categories $x_t^{(i,j)} \in [1, V]$, the forward transition probabilities
 154 from x_{t-1} to x_t can be represented by matrices $[Q_t]_{mn} = q(x_t = m|x_{t-1} = n) \in \mathbb{R}^{V \times V}$. Note that
 155 we omit the superscripts (i, j) to avoid confusion. Then the forward diffusion process can be written
 156 as:

$$q(x_t|x_{t-1}) = \mathbf{x}_t^T Q_t \mathbf{x}_{t-1}, \quad (2)$$

157 where $\mathbf{x}_t \in \mathbb{R}^{V \times 1}$ is the one-hot version of x_t and $Q_t \mathbf{x}_{t-1}$ is the categorical distribution for x_t . A
 158 nice property of the above Markov diffusion process is that we can sample x_t as any timestep directly
 159 from x_0 as:

$$q(x_t|x_0) = \mathbf{x}_t^T \bar{Q}_t \mathbf{x}_0, \text{ with } \bar{Q}_t = Q_t \dots Q_1. \quad (3)$$

160 D3PM [11] formulate the transition matrix $Q_t \in \mathbb{R}^{V \times V}$ by introducing a small number of uniform
161 noises to the categorical distribution. As formulated as the first matrix in Eq. (4) with $\alpha \in [0, 1]$ and
162 $\beta_t = (1 - \alpha_t)/V$. It can be interpreted as each token having a probability of $\alpha_t + \beta_t$ to remain the
163 previous value and a probability of β_t to be the value from the whole V categories. Based on D3PM,
164 VQ-Diffusion [12] propose a mask-and-replace diffusion strategy that not only replaces the previous
165 value but also insert [MASK] token to explicitly figure out the tokens that have been replaced. We
166 extend this mask-and-replace strategy to our quantized pose sequence modelling. Since the length of
167 pose sequences may be different in a minibatch, we have to add two special tokens, [MASK] and
168 [PAD] tokens, so each token has $V + 2$ states. The mask-and-replace diffusion process can be defined
169 as follows: each token has a probability of α_t to be unchanged, $V\beta_t$ to be uniformly resampled and
170 $\gamma_t = 1 - \alpha_t - V\beta_t$ to be replaced with [MASK] token. **Note that [MASK] and [PAD] tokens always**
171 **keep its own state. The difference is that [PAD] is used to represent the padding part in the initial**
172 **sequence x_0 , and [MASK] is used to replace the original token of the code sequence in the diffusion**
173 **process. Moreover, in the revised denoising process, the [MASK] positions are required to predict**
174 **but [PAD] positions need to be ignored.** The transition matrix $Q_t \in \mathbb{R}^{(V+2) \times (V+2)}$ is formulated as
175 the second matrix of the following:

$$Q_t = \begin{bmatrix} \alpha_t + \beta_t & \beta_t & \cdots & \beta_t \\ \beta_t & \alpha_t + \beta_t & \cdots & \beta_t \\ \vdots & \vdots & \ddots & \vdots \\ \beta_t & \beta_t & \cdots & \alpha_t + \beta_t \end{bmatrix}; Q_t = \begin{bmatrix} \alpha_t + \beta_t & \beta_t & \cdots & \beta_t & 0 & 0 \\ \beta_t & \alpha_t + \beta_t & \cdots & \beta_t & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \beta_t & \beta_t & \cdots & \alpha_t + \beta_t & 0 & 0 \\ \gamma_t & \gamma_t & \cdots & \gamma_t & 1 & 0 \\ 0 & 0 & \cdots & 0 & 0 & 1 \end{bmatrix}. \quad (4)$$

176 Finally, the categorical distribution of \mathbf{x}_t can be derived as following using reparameterization trick:

$$\begin{aligned} \text{when } x_0 \neq V + 2, \quad \bar{Q}_t \mathbf{x}_0 &= \begin{cases} \bar{\alpha}_t + \bar{\beta}_t, & x_t = x_0 \\ \bar{\beta}_t, & x_t \neq x_0 \text{ and } x_t \leq V \\ \bar{\gamma}_t, & x_t = V + 1 \\ 0, & x_t = V + 2 \end{cases} \\ \text{when } x_0 = V + 2, \quad \bar{Q}_t \mathbf{x}_0 &= \begin{cases} 0, & x_t \neq V + 2 \\ 1, & x_t = V + 2 \end{cases} \end{aligned} \quad (5)$$

177 where $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, $\bar{\gamma}_t = 1 - \prod_{i=1}^t (1 - \gamma_i)$, and $\bar{\beta}_t = (1 - \bar{\alpha}_t - \bar{\gamma}_t)/V$. Therefore, we can directly
178 sample x_t within the computation cost $O(V)$.

179 **Reverse Denoising Process.** The reverse denoising process aims to recreate the real sample from a
180 full noise input by gradually sampling from $q(x_{t-1}|x_t)$. However, it is intractable to estimate the
181 conditional probability $q(x_{t-1}|x_t)$ since it needs to use the whole dataset. Fortunately, the conditional
182 probability is tractable when conditioned on x_0 using Bayes' rule:

$$q(x_{t-1}|x_t, x_0) = \frac{q(x_t|x_{t-1}, x_0)q(x_{t-1}|x_0)}{q(x_t|x_0)} = \frac{(\mathbf{x}_t^T Q_t \mathbf{x}_{t-1})(\mathbf{x}_{t-1}^T \bar{Q}_{t-1} \mathbf{x}_0)}{\mathbf{x}_t^T \bar{Q}_t \mathbf{x}_0}, \quad (6)$$

183 thus we train a denoising model $p_\theta(x_{t-1}|x_t, c)$ to approximate the tractable distribution
184 $q(x_{t-1}|x_t, x_0)$, where c is the conditional feature of gloss sequence. And the model is trained
185 to minimize the variational lower bound [37]:

$$\begin{aligned} \mathcal{L}_{vb} = & \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(x_T|x_0) \parallel p_\theta(x_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t, c))}_{L_{t-1}} \right. \\ & \left. - \underbrace{\log p_\theta(x_0|x_1, c)}_{L_0} \right]. \end{aligned} \quad (7)$$

186 **Reparameterization Trick on Model Learning.** Compared with directly predicting $p_\theta(x_{t-1}|x_t, c)$,
187 recent works [30, 10, 12] find that predict the data x_0 gives better quality at every reverse step. Thus,

188 we let our denoising model to predict the distribution $p_\theta(\tilde{x}_0|x_t, c)$. With a reparameterization trick,
 189 the conditional reverse distribution can be formulated as:

$$p_\theta(x_{t-1}|x_t, c) = \sum_{\tilde{x}_0=1}^V q(x_{t-1}|x_t, \tilde{x}_0)p_\theta(\tilde{x}_0|x_t, y). \quad (8)$$

190 Under this x_0 -parameterization trick, we introduce an auxiliary denoising objective to encourage
 191 good predictions of the data x_0 at each time step [11]. The final loss function is combined with the
 192 negative variational lower bound and the auxiliary loss:

$$\mathcal{L}_{\text{ddm}} = \mathcal{L}_{\text{vb}} - \lambda \log p_\theta(x_0|x_t, c), \quad (9)$$

193 where λ is a coefficient for the auxiliary loss term.

194 **CodeUnet for Model Learning.** Most image diffusion models [29, 30, 39] adopt the Unet [15] as
 195 their architectures since it is effective for data with spatial structure. However, directly applying the
 196 Unet in discrete sequence generation, *e.g.*, text generation [11] and quantized image synthesis [12],
 197 will bring information leakage problem, since the convolution layer over adjacent tokens may provide
 198 shortcuts for the mask-based prediction [40]. Therefore, Austin *et al.* [11] and Gu *et al.* [12] use the
 199 token-wise Transformer framework to learn the distribution $p_\theta(\tilde{x}_0|x_t, c)$. In this work, to incorporate
 200 the advantages of Unet and Transformer networks, we propose a novel architecture CodeUnet to
 201 learn the spatial-temporal interaction for our quantized pose sequence generation.

202 As shown in Figure 2, the CodeUnet consists of a contracting path (left side), an expansive path (right
 203 side) and a middle module. The middle module is an encoder-decoder Transformer framework. The
 204 encoder consists of 6 Transformer blocks. It takes the gloss sentence as input and obtains a conditional
 205 feature sequence. The decoder has two blocks. Each block has a self-attention, a cross-attention, a
 206 feed-forward network and an Adaptive Layer Normalization [41, 12](AdaLN). The AdaLN operator
 207 is devised to incorporate timestep t information as $\text{AdaLN}(h, t) = \alpha_t \text{LayerNorm}(h) + b_t$, where
 208 h is the intermediate activations, α_t and β_t are obtained from a linear projection of the timestep
 209 embedding.

210 Both contracting path and the expansive path are hierarchical structures and each level has two
 211 Transformer encoder blocks. For downsampling in contracting path, given the feature of quantized
 212 pose sequence, *e.g.*, $h \in \mathbb{R}^{N \times 3 \times d_{\text{model}}}$, where d_{model} is the feature dimension, we first sample uniformly
 213 with stride 2 in the temporal dimension and remain constant in the spatial dimension. Then we set the
 214 downsampled feature as query $Q \in \mathbb{R}^{N/2 \times 3 \times d_{\text{model}}}$, and keep key K and value V unchanged for the
 215 following attention network. In the upsampling of expansive path, we directly repeat the feature 2
 216 times as a query, but the key and value remains for the following attention network:

$$\forall n = 1, \dots, N, Q_n^{\text{up}} = h_{n//2}, K^{\text{up}} = V^{\text{up}} = h, \quad (10)$$

217 where $\cdot//\cdot$ denotes floor division. Finally, a linear layer and a softmax layer are applied to make the
 218 prediction.

219 3.3 Length Prediction with Sequential-KNN

220 In this section, inspired by [17] which merges tokens with similar semantic meanings from different
 221 locations, we propose a novel clustering algorithm to get the lengths for corresponding glosses.
 222 Specifically, given a token sequence which is obtained from the Pose-VQVAE model, we compute
 223 the local density ρ of each token according to its k-nearest-neighbors:

$$\rho_i = \exp\left(-\frac{1}{k} \sum_{z_j \in \text{KNN}(z_i)} \|z_i - z_j\|_2^2\right), \text{ where } |i - j| \leq l \quad (11)$$

224 where i, j is the position in the sequence, l is a predefined hyperparameter indicating that we only
 225 consider the local region since the adjacent tokens are more likely to belong to a gloss. z_i and z_j are
 226 the latent feature for i^{th} and j^{th} tokens. $\text{KNN}(z_i)$ represents the k-nearest neighbors for i^{th} token.

227 We assign $\{p_1, \dots, p_M\}$ positions with a higher local density as the peaks, where M is the length
 228 of the gloss sequence. Then between two adjacent peaks, for example p_1 and p_2 , we sequentially
 229 iterate from p_1 to p_2 , and find the first position that is farther from z_{p_1} and closer to z_{p_2} , which is

230 the boundary we determined. After finding these boundaries, we get the lengths of the contiguous
 231 pose sequence for its corresponding glosses. As shown in Figure 2, we define the obtained lengths
 232 as $\{L_1, \dots, L_M\}$, and the Transformer encoder for gloss sequence is trained under the supervised
 233 information of lengths. For each gloss word, we predict a number from $[1, P]$, where P is the
 234 maximum length of the target pose sequence. Mathematically, we formulate the classification loss of
 235 length prediction as:

$$\mathcal{L}_{\text{len}} = \frac{\delta}{M} \sum_i^M \sum_j^P (-L_i = j) \log p(L_i|c). \quad (12)$$

236 In the training of the discrete diffusion mode, \mathcal{L}_{len} is trained together with a coefficient δ . In the
 237 inference, we predict the length of glosses, and their summation is the length of target pose sequence.

238 In summary, we arrive at our proposed two-stage approach, PoseVQ-Diffusion, with the first-stage
 239 Pose-VQVAE model and the second-stage discrete diffusion model with a length predictor. **The**
 240 **whole training and inference algorithm is shown in supplementary material.**

241 4 Experiments

242 4.1 Experiment Setups and Implementation Details

243 **Dataset.** We evaluate our G2P model on RWTH-PHOENIX-Weather 2014T (RPWT) dataset [18]. It
 244 is the *only* publicly available SLP dataset with parallel sign language videos, gloss annotations, and
 245 spoken language translations. This corpus contains 7096 training samples (with 1066 different sign
 246 glosses in gloss annotations and 2887 words in German spoken language translations), 519 validation
 247 samples, and 642 test samples.

248 **Evaluation Criteria.** Following the widely-used setting in SLP [3], we adopt the back-translation
 249 method for evaluation. Specifically, we utilize the state-of-the-art SLT [19] model to translate the
 250 generated sign pose sequence back to gloss sequence and spoken language, where its input is modified
 251 as pose sequence. Specifically, we compute BLEU [42] and Word Error Rate (WER) between the
 252 back-translated spoken language translations and gloss recognition results with ground truth spoken
 253 language and gloss sequence.

254 **Data Processing.** Since the RWTH-PHOENIX-Weather 2014T (RPWT) dataset doesn't contain pose
 255 information, we generate the pose sequence as the ground truth. Following B. Saunders *et al.* [3],
 256 we extract 2D joint points from sign video using OpenPose [43] and lift the 2D joints to 3D with
 257 a skeletal model estimation improvement method [44]. Finally, similar to [24], we apply skeleton
 258 normalization to remove the skeleton size difference between different signers.

259 **Model Setting.** The Pose-VQVAE consists of an Encoder, a Tokenizer, and a Decoder. The Encoder
 260 contains a linear layer to transform pose points to hidden feature with dimension set as 256, a 3-layer
 261 Transformer module with divided space-time attention [45]. The Tokenizer maintains a codebook
 262 with a size set as 2048. The Decoder contains the same 3-layer Transformer module as the encoder
 263 and an SPL layer to predict the structural sign skeleton. For the discrete diffusion model, we set
 264 the timestep T as 100. All Transformer blocks of CodeUnet have $d_{\text{model}}=512$ and $N_{\text{depth}}=2$. The
 265 size of local region l in Equation 11, is set as 16 which is the average length of a gloss. And the
 266 number of nearest neighbors k is set as 16. **We optimize our network using AdamW [46] with**
 267 **$\beta_1 = 0.9$ and $\beta_2 = 0.96$. The learning rate is set to 0.0004 after 5000 iterations of warmup.** We train
 268 the model on 8 NVIDIA Tesla V100 GPUs. We include all hyperparameters setting and the details of
 269 implementation in the supplementary material.

270 4.2 Comparisons with State-of-the-Art Methods

271 **Competing Methods.** We compare our PoseVQ-Diffusion with previous state-of-the-art G2P
 272 models. **Progressive Transformer (PTR)** [3] is the first SLP model to tackle the G2P problem in an
 273 autoregressive manner. Since they use the ground-truth first sign pose frame and timing information,
 274 their reported results are not comparable to ours. Thus we adopt the results reported by Huang *et*
 275 *al.* [5]. **NAT-EA** [5] propose a non-autoregressive method to directly predict the target pose sequence

Table 1: Quantitative results for G2P task on RWTH-PHOENIX-Weather 2014T test dataset. † indicates the results is provided by Huang et al. [5]. **Note that smaller WER is better, higher BLEU is better, and lower DTW-MJE is better. The closer all the results are to the GT, the better.**

Method	WER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	DTW-MJE
PTR† [3]	94.65	11.45	7.08	5.08	4.04	0.191
MoMP† [27]	92.41	13.17	8.24	6.25	4.75	0.188
NAT-AT [5]	88.15	14.26	9.93	7.11	5.53	0.177
NAT-EA [5]	82.01	15.12	10.45	7.99	6.66	0.146
PoseVQ-AR (Ours)	85.27	14.26	10.02	7.57	5.94	0.172
PoseVQ-MP (Ours)	79.38	15.43	10.69	8.26	6.98	0.146
PoseVQ-Diffusion (Ours)	78.21	16.03	11.32	9.17	7.42	0.122
GT	50.23	23.47	15.86	12.03	10.47	0.0

276 with the External Aligner (EA) to learn alignments between glosses and pose sequence. **NAT-AT** is
 277 the NAT model without EA that uses the decoder-to-encoder attention to learn the alignments.

278 **Quantitative Comparison.** The comparison between our PoseVQ-Diffusion and the competing
 279 methods is shown in Tabel 1. The row of **PoseVQ-AR** refs to the vector quantized model with
 280 an autoregressive decoder. The row of **PoseVQ-MP** refs to the vector quantized model with the
 281 Mask-Predict [9] strategy, which is also a variant of discrete diffusion model [11]. **PoseVQ-Diffusion**
 282 refs to the vector quantized model with mask-and-replace diffusion strategy. As indicated in Table 1,
 283 both diffusion-based models outperform the state-of-the-art G2P models with relative improvements
 284 on WER score by 4.6% (82.01 \rightarrow 78.21) and on BLEU-4 by 11.4% (6.66 \rightarrow 7.42). This shows
 285 the effectiveness of the iterative mask-based non-autoregressive method on the vector quantized
 286 pose sequence. In addition, the Mask-Predict strategy is a mask only strategy that is similar to
 287 PoseVQ-Diffusion with $\bar{\gamma}_T = 1$. Therefore, PoseVQ-Diffusion achieves better performance than
 288 PoseVQ-MP. This reflects the mask-and-replace strategy is superior to the mask only strategy.

289 4.3 Model Analysis and Discussions

290 We also investigate the effects of different components and design choices of our proposed model.

291 **Analysis of The Design of Pose-VQVAE.** As shown in the first three rows of Table 2a, we study
 292 the design of our Pose-VQVAE model. Pose-VQVAE-joint means we compress all points into one
 293 token, Pose-VQVAE-separate means the points are separated into three local patches according to the
 294 structure of a sign skeleton. Empirically, Pose-VQVAE-separate achieves much better reconstruction
 295 (MSE) performance. This indicates that compressing all skeleton points into one token embedding is
 296 not advisable, leading to information loss. The second row of Figure 3 shows the sample of sign pose
 297 sequences reconstructed by Pose-VQVAE-separate.

298 **CodeUnet vs. Transformer.** For a fair comparison, we replace our CodeUnet with Transformer
 299 network with keeping other settings the same. As shown in the last three rows of Table 2a, the
 300 diffusion-based model with our CodeUnet achieves better performance on the back-translate evalua-
 301 tion. This phenomenon suggests that the hierarchical structure of CodeUnet makes it particularly
 302 effective for data with spatial structure. Moreover, in our experiments with the same batch size,
 303 CodeUnet coverages faster than Transformer. Having said that, due to sign pose sequences being
 304 temporally redundant, the compression of CodeUnet in the time dimension makes it more efficient in
 305 training.

306 **Number of Timesteps.** We compare the performance of the model with different numbers of training
 307 steps. As shown in the left two columns of Table 2b, we find that the results get better when the
 308 training step size is increased from 10 to 100. As it increased further, the results seemed to saturate.
 309 Therefore, we set the training step to 100 to trade off performance and speed.

310 **Length Candidates.** Length prediction is essential for a non-autoregressive generation. Our approach
 311 proposes a Sequential-KNN algorithm to learn the lengths for corresponding glosses and then treat
 312 the length prediction as a classification problem. As shown in the right two columns in Table 2b,
 313 we study the performance with different length candidates and compare it with the reference (gold)
 314 target sequence length. The results show that multiple candidates can increase performance, but too
 315 many candidates can even degrade performance.



Figure 3: G2P qualitative results. We show some examples of predicted sign pose sequences compared with our reconstruction model and previous G2P model [3]. For readability, we sampled every 5 frames for a total of 16 frames. See our supplementary material for more results.

Table 2: Analysis into the effects of different designs and hyperparameters for our propose model.

(a) Ablation on design of reconstruction and prediction model.

Reconstruction Model	MSE (\downarrow)
Pose-VQVAE-joint	0.0242
Pose-VQVAE-seperate	0.0139
Prediction Model	WER (\downarrow)
Transformer	80.36
CodeUnet	78.21

(b) Ablation on the hyperparameters of training steps and length candidates.

Training Steps	WER (\downarrow)	Length Candidates	WER (\downarrow)
10	81.06	1	79.45
50	79.31	2	78.69
100	78.21	3	78.21
150	78.17	4	78.74
200	78.15	Gold	77.26

316 5 Conclusion

317 In this paper, we presented a novel paradigm for conditional sign pose sequence generation through an
 318 iterative non-autoregressive method. Specifically, we first devise a specific architecture Pose-VQVAE
 319 to learn discrete codes by reconstruction. Then we extend the discrete diffusion model to model the
 320 sequential alignments between sign glosses and quantized codes. And a “fully transformer” network
 321 CodeUnet is proposed for the spatial-temporal information in discrete space. Finally, we propose a
 322 sequential-KNN algorithm to learn the length of corresponding glosses and then predict the length as
 323 a classification task. Compared with previous state-of-the-art autoregressive and non-autoregressive
 324 methods, extensive experiments demonstrate the effectiveness of our proposed PoseVQ-Diffusion
 325 framework.

326 6 Broader Impact and Limitations

327 We develop a general paradigm for conditional pose generation in this paper. We do not foresee any
 328 negative ethical/societal impacts at this moment. Although the proposed PoseVQ-Diffusion proves
 329 effective in conditional sign pose sequence generation, we notice several limitations of our approach.
 330 (i) Although our discrete diffusion model-based model has a faster sampling speed in longer sequence
 331 generation than traditional autoregressive models, it is much slower than one-shot non-autoregressive
 332 models. (ii) Our proposed two-stage based models are not end-to-end and thus more difficult to train
 333 than previous methods. We, therefore, plan to resolve the aforementioned issues and further mitigate
 334 the training and inference speed between the one-shot non-autoregressive method.

335 **References**

- 336 [1] Roland Pfau, Martin Salzman, and Markus Steinbach. The syntax of sign language agreement:
337 Common ingredients, but unusual recipe. *Glossa: a journal of general linguistics*, 2018.
- 338 [2] Ben Saunders, Richard Bowden, and Necati Cihan Camgöz. Adversarial training for multi-
339 channel sign language production. In *31st British Machine Vision Conference 2020, BMVC*
340 *2020, Virtual Event, UK, September 7-10, 2020*. BMVA Press, 2020.
- 341 [3] Ben Saunders, Necati Cihan Camgöz, and R. Bowden. Progressive transformers for end-to-end
342 sign language production. *ArXiv preprint*, abs/2004.14874, 2020.
- 343 [4] Ben Saunders, Necati Cihan Camgoz, and R. Bowden. Everybody sign now: Translating spoken
344 language to photo realistic sign language video. *ArXiv preprint*, abs/2011.09846, 2020.
- 345 [5] Wencan Huang, Wenwen Pan, Zhou Zhao, and Qi Tian. Towards fast and high-quality sign
346 language production. *Proceedings of the 29th ACM International Conference on Multimedia*,
347 2021.
- 348 [6] Ronald J. Williams and David Zipser. A learning algorithm for continually running fully
349 recurrent neural networks. *Neural Computation*, 1:270–280, 1989.
- 350 [7] Florian Schmidt. Generalization in generation: A closer look at exposure bias. In *Proceedings*
351 *of the 3rd Workshop on Neural Generation and Translation*, pages 157–167, Hong Kong, 2019.
352 Association for Computational Linguistics.
- 353 [8] Jiatao Gu, James Bradbury, Caiming Xiong, Victor O. K. Li, and Richard Socher. Non-
354 autoregressive neural machine translation. In *6th International Conference on Learning Rep-*
355 *resentations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track*
356 *Proceedings*. OpenReview.net, 2018.
- 357 [9] Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-predict: Parallel
358 decoding of conditional masked language models. In *Proceedings of the 2019 Conference on*
359 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
360 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 6112–6121, Hong Kong, China,
361 2019. Association for Computational Linguistics.
- 362 [10] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forr’e, and Max Welling. Argmax
363 flows and multinomial diffusion: Towards non-autoregressive language models. *ArXiv preprint*,
364 abs/2102.05379, 2021.
- 365 [11] Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg.
366 Structured denoising diffusion models in discrete state-spaces. In *NeurIPS*, 2021.
- 367 [12] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and
368 Baining Guo. Vector quantized diffusion model for text-to-image synthesis. *ArXiv preprint*,
369 abs/2111.14822, 2021.
- 370 [13] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution
371 image synthesis. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*
372 *(CVPR)*, pages 12868–12878, 2021.
- 373 [14] Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation
374 learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus,
375 S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing*
376 *Systems 30: Annual Conference on Neural Information Processing Systems 2017, December*
377 *4-9, 2017, Long Beach, CA, USA*, pages 6306–6315, 2017.
- 378 [15] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for
379 biomedical image segmentation. In *MICCAI*, 2015.
- 380 [16] Mingjing Du, Shifei Ding, and Hongjie Jia. Study on density peaks clustering based on k-nearest
381 neighbors and principal component analysis. *Knowl. Based Syst.*, 99:135–145, 2016.

- 382 [17] Wang Zeng, Sheng Jin, Wentao Liu, Chen Qian, Ping Luo, Ouyang Wanli, and Xiaogang Wang.
383 Not all tokens are equal: Human-centric visual analysis via token clustering transformer. *ArXiv*
384 *preprint*, abs/2204.08680, 2022.
- 385 [18] Necati Cihan Camgöz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden.
386 Neural sign language translation. In *2018 IEEE Conference on Computer Vision and Pattern*
387 *Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 7784–7793. IEEE
388 Computer Society, 2018.
- 389 [19] Necati Cihan Camgöz, Oscar Koller, Simon Hadfield, and Richard Bowden. Sign language
390 transformers: Joint end-to-end sign language recognition and translation. In *2020 IEEE/CVF*
391 *Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June*
392 *13-19, 2020*, pages 10020–10030. IEEE, 2020.
- 393 [20] Necati Cihan Camgöz, Oscar Koller, Simon Hadfield, and R. Bowden. Multi-channel transform-
394 ers for multi-articulatory sign language translation. *ArXiv preprint*, abs/2009.00299, 2020.
- 395 [21] Hao Zhou, Wen gang Zhou, Yun Zhou, and Houqiang Li. Spatial-temporal multi-cue network
396 for sign language recognition and translation. *IEEE Transactions on Multimedia*, 24:768–779,
397 2022.
- 398 [22] Pan Xie, Mengyi Zhao, and Xiaohui Hu. Pisltrc: Position-informed sign language transformer
399 with content-aware convolution. *ArXiv preprint*, abs/2107.12600, 2021.
- 400 [23] Hezhen Hu, Weichao Zhao, Wen gang Zhou, Yuechen Wang, and Houqiang Li. Signbert: Pre-
401 training of hand-model-aware representation for sign language recognition. *2021 IEEE/CVF*
402 *International Conference on Computer Vision (ICCV)*, pages 11067–11076, 2021.
- 403 [24] Stephanie Stoll, Necati Cihan Camgöz, Simon Hadfield, and Richard Bowden. Sign language
404 production using neural machine translation and generative adversarial networks. In *British*
405 *Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*, page 304.
406 BMVA Press, 2018.
- 407 [25] Qinkun Xiao, Mingyong Qin, and Yuting Yin. Skeleton-based chinese sign language recognition
408 and generation for bidirectional communication between deaf and hearing people. *Neural*
409 *networks : the official journal of the International Neural Network Society*, 125:41–55, 2020.
- 410 [26] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A. Efros. Everybody dance now.
411 In *2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea*
412 *(South), October 27 - November 2, 2019*, pages 5932–5941. IEEE, 2019.
- 413 [27] Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. Mixed signals: Sign language
414 production via a mixture of motion primitives. In *2021 IEEE/CVF International Conference on*
415 *Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 1899–1909.
416 IEEE, 2021.
- 417 [28] Jaehyeon Kim, Sungwon Kim, Jungil Kong, and Sungroh Yoon. Glow-tts: A generative
418 flow for text-to-speech via monotonic alignment search. In Hugo Larochelle, Marc’Aurelio
419 Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural*
420 *Information Processing Systems 33: Annual Conference on Neural Information Processing*
421 *Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 422 [29] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*,
423 2021.
- 424 [30] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo
425 Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin,
426 editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural*
427 *Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 428 [31] Jonathan Ho, Chitwan Saharia, William Chan, David Fleet, Mohammad Norouzi, and Tim
429 Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*,
430 23:47:1–47:33, 2022.

- 431 [32] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic
432 models. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International*
433 *Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of
434 *Proceedings of Machine Learning Research*, pages 8162–8171. PMLR, 2021.
- 435 [33] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
436 resolution image synthesis with latent diffusion models. *ArXiv preprint*, abs/2112.10752,
437 2021.
- 438 [34] Myeonghun Jeong, Hyeongju Kim, Sung Jun Cheon, Byoung Jin Choi, and Nam Soo Kim.
439 Diff-tts: A denoising diffusion model for text-to-speech. In *Interspeech*, 2021.
- 440 [35] Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-tts: A diffusion model for text-to-
441 speech via classifier guidance. 2021.
- 442 [36] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A
443 versatile diffusion model for audio synthesis. In *9th International Conference on Learning*
444 *Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- 445 [37] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep
446 unsupervised learning using nonequilibrium thermodynamics. In Francis R. Bach and David M.
447 Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning, ICML*
448 *2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*,
449 pages 2256–2265. JMLR.org, 2015.
- 450 [38] Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. Structured prediction helps 3d human
451 motion modelling. In *2019 IEEE/CVF International Conference on Computer Vision, ICCV*
452 *2019, Seoul, Korea (South), October 27 - November 2, 2019*, pages 7143–7152. IEEE, 2019.
- 453 [39] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and
454 Ben Poole. Score-based generative modeling through stochastic differential equations. In *9th*
455 *International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May*
456 *3-7, 2021*. OpenReview.net, 2021.
- 457 [40] Piotr Nawrot, Szymon Tworowski, Michal Tyrolski, Lukasz Kaiser, Yuhuai Wu, Christian
458 Szegedy, and Henryk Michalewski. Hierarchical transformers are more efficient language
459 models. *ArXiv preprint*, abs/2110.13711, 2021.
- 460 [41] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv preprint*,
461 abs/1607.06450, 2016.
- 462 [42] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
463 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Associa-*
464 *tion for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA, 2002.
465 Association for Computational Linguistics.
- 466 [43] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: Realtime
467 multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis*
468 *and Machine Intelligence*, 43:172–186, 2021.
- 469 [44] Jan Zelinka and Jakub Kanis. Neural sign language synthesis: Words are our glosses. *2020*
470 *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 3384–3392, 2020.
- 471 [45] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need
472 for video understanding? In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th*
473 *International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*,
474 volume 139 of *Proceedings of Machine Learning Research*, pages 813–824. PMLR, 2021.
- 475 [46] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International*
476 *Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*.
477 OpenReview.net, 2019.

478 **Checklist**

- 479 • Do the main claims made in the abstract and introduction accurately reflect the paper’s
480 contributions and scope? [Yes]
- 481 • Have you read the ethics review guidelines and ensured that your paper conforms to them?
482 [Yes]
- 483 • Did you discuss any potential negative societal impacts of your work? [Yes] See Section 6
- 484 • Did you describe the limitations of your work? [Yes] See Section 6
- 485 • Did you state the full set of assumptions of all theoretical results? [Yes]
- 486 • Did you include complete proofs of all theoretical results? [Yes] More proofs are in
487 supplemental material.
- 488 • Did you include the code, data, and instructions needed to reproduce the main experimental
489 results (either in the supplemental material or as a URL)? [Yes] The code and instructions
490 are in the the supplemental material
- 491 • Did you specify all the training details (e.g., data splits, hyperparameters, how they were
492 chosen)? [Yes] More training details are in the supplemental material
- 493 • Did you report error bars (e.g., with respect to the random seed after running experiments
494 multiple times)? [No]
- 495 • Did you include the amount of compute and the type of resources used (e.g., type of GPUs,
496 internal cluster, or cloud provider)? [Yes] See Section 4.1
- 497 • If your work uses existing assets, did you cite the creators? [Yes]
- 498 • Did you mention the license of the assets? [No]
- 499 • Did you include any new assets either in the supplemental material or as a URL? [No]
- 500 • Did you discuss whether and how consent was obtained from people whose data you’re
501 using/curating? [No]
- 502 • Did you discuss whether the data you are using/curating contains personally identifiable
503 information or offensive content? [No]
- 504 • Did you include the full text of instructions given to participants and screenshots, if applica-
505 ble? [No]
- 506 • Did you describe any potential participant risks, with links to Institutional Review Board
507 (IRB) approvals, if applicable? [No]
- 508 • Did you describe any potential participant risks, with links to Institutional Review Board
509 (IRB) approvals, if applicable? [Yes]