

# T2S-GPT: Dynamic Vector Quantization for Autoregressive Sign Language Production from Text

Aoxiong Yin, Haoyuan Li, Kai Shen, Siliang Tang\*, Yueting Zhuang

Zhejiang University

{yinaoxiong, lihaoyuan, shenkai, siliang, yzhuang}@zju.edu.cn

## Abstract

In this work, we propose a two-stage sign language production (SLP) paradigm that first encodes sign language sequences into discrete codes and then autoregressively generates sign language from text based on the learned codebook. However, existing vector quantization (VQ) methods are fixed-length encodings, overlooking the uneven information density in sign language, which leads to under-encoding of important regions and over-encoding of unimportant regions. To address this issue, we propose a novel dynamic vector quantization (DVA-VAE) model that can dynamically adjust the encoding length based on the information density in sign language to achieve accurate and compact encoding. Then, a GPT-like model learns to generate code sequences and their corresponding durations from spoken language text. Extensive experiments conducted on the PHOENIX14T dataset demonstrate the effectiveness of our proposed method. To promote sign language research, we propose a new large German sign language dataset, PHOENIX-News, which contains 486 hours of sign language videos, audio, and transcription texts. Experimental analysis on PHOENIX-News shows that the performance of our model can be further improved by increasing the size of the training data. Our project homepage is <https://t2sgpt-demo.yinaoxiong.cn>.

## 1 Introduction

Sign language is a visual language with complex grammatical structures and is the primary means of communication for nearly 70 million deaf people worldwide<sup>1</sup>. Research on sign language production (Baltatzis et al., 2023a; Fang et al., 2023; Huang et al., 2021; Hwang et al., 2021, 2022; Saunders et al., 2020a) and sign language translation (Camgoz et al., 2018a; Zhang et al., 2023a;

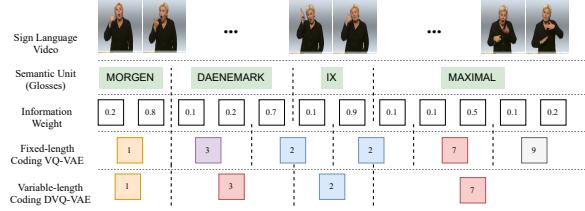


Figure 1: Comparison of fixed-length encoding and variable-length encoding.

Zhou et al., 2021; Yin et al., 2021, 2023) has attracted widespread attention. Sign language production (SLP) is a challenging problem that aims to automatically translate spoken language descriptions into corresponding continuous sign sequences. SLP can help deaf people better access information and communicate with others, thereby facilitating their lives, which has important social significance.

SLP models are expected to learn precise mapping from the spoken language space to the sign language space. Early work used 2D or 3D skeleton poses to represent sign language (Huang et al., 2021; Saunders et al., 2021b, 2020b), while recent work has suggested using 3D human models, such as SMPL-x(Pavlakos et al., 2019), to represent sign language, as it introduces human priors and can better animate (Baltatzis et al., 2023b). To learn the mapping between these two different modal spaces, some work uses autoregressive models (Saunders et al., 2021a,b, 2020b), non-autoregressive models (Huang et al., 2021; Hwang et al., 2021, 2022), or diffusion models (Baltatzis et al., 2023b) to learn the direct mapping from spoken language text to sign language skeleton poses. (Xie et al., 2023) proposed to learn the discrete representation of sign language through VQ-VAE (van den Oord et al., 2017) and then learn the mapping from text to discrete representation through a discrete diffusion model. However, we found that existing sign language discrete representation methods are fixed-

<sup>\*</sup>Corresponding author.

<sup>1</sup>According to World Federation of the Deaf <https://wfdeaf.org/our-work/>

length encodings, as shown in [Figure 1](#), which overlooks the uneven information density in sign language. In addition, many existing works rely on expert-annotated intermediate representations, i.e. glosses, which limit the scalability of the model.

In this work, we are inspired by recent advances from learning the discrete representation for generation ([Huang et al., 2023; Zhang et al., 2023b; Williams et al., 2020; van den Oord et al., 2017; Ao et al., 2022](#)). Specifically, we investigate a two-stage framework based on Dynamic Vector Quantized Variational Autoencoders (DVQ-VAE) and Generative Pre-trained Transformer (GPT) ([Radford et al., 2018](#)) for text-to-sign language production. In the first stage, as shown in [Figure 1](#), DVQ-VAE will learn the weights of each frame and the boundaries of the basic semantic units. Then, the weighted latent vectors are mapped to discrete code indices. Further quantitative analysis of the uneven information density in sign language is provided in [section 3](#). To encourage models to perform variable-length encoding and compress sequence lengths, we propose a novel budget loss. Additionally, to preserve the semantic information of the reconstructed sign language sequences, we also introduce a translation auxiliary loss. In the second stage, a GPT-like model is learned to generate code index sequences from spoken language text. Furthermore, since the duration of quantized code in a sequence can also vary dynamically, we further propose a duration transformer to predict the duration of the next code based on the previous code’s duration and the current code.

The experimental results on the widely used SLP dataset PHOENIX14T ([Camgoz et al., 2018b](#)) demonstrate that our proposed method achieves superior back translation performance compared to previous approaches. Furthermore, throughout the entire development process of image generation and text generation, the scale of the dataset has played a crucial role. A large amount of high-quality corpus is also very important for SLP tasks. In this paper, we present the largest known German Sign Language dataset, PHOENIX-News, which consists of 486 hours of sign language videos, audio, and transcription texts. The native expression, clear hand details, and extensive coverage of our large-scale dataset make it suitable for a variety of sign language research tasks, such as sign language translation and sign language production. Based on this dataset, we further explore the impact of training data size on SLP tasks. Empirical analy-

sis shows that the performance of our model can be further improved by increasing the size of the training data.

Our main contributions are summarized as follows:

- We analyse the uneven information density in sign language. Additionally, we propose for the first time an information density based variable length coding method suitable for sign language.
- We propose a two-stage SLP framework consisting of two components: 1) DVQ-VAE to dynamically assign variable-length codes to sequences based on their different information densities through a novel *adaptive downsampling module* and *budget loss*. 2) A novel *T2M-GPT model* to predict variable-length codes and their corresponding durations.
- Extensive experiments on the challenging PHOENIX14T dataset show the effectiveness of our proposed method.
- We propose the largest known German sign language dataset, PHOENIX-News, which can be used for a variety of sign language research tasks.

## 2 Related Work

### 2.1 Sign Language Production

Sign language production (SLP) has been an active area of research for nearly two decades ([Cox et al., 2002; McDonald et al., 2016](#)). Early approaches focused on mapping text to glosses using neural models. ([Stoll et al., 2020](#)) proposed a seq2seq architecture for SLP, which mapped text input to glosses. To generate 2D joint locations, they utilized an empirical lookup table paradigm. Then, ([Saunders et al., 2020b](#)) proposed a progressive transformer to directly learn the mapping between annotations and skeleton pose sequences. ([Saunders et al., 2020a](#)) proposed to improve the quality of skeleton pose generation through adversarial training. In addition, several approaches have been proposed to enhance the generation quality through the utilization of mixture density networks ([Saunders et al., 2021a](#)), Mixture-of-Experts ([Saunders et al., 2021b](#)), dictionary representations ([Saunders et al., 2022](#)), and diffusion models ([Baltatzis et al., 2023a](#)). Several studies have proposed the use of non-autoregressive

Table 1: Summary statistics for different sign language datasets.

Dataset	Language	Attribute				Statistics			Source
		Transcription	Pose	Speech	Document-level	Duration(h)	Vocab	Signers	
BOBSL (Albanie et al., 2021)	BSL	✓	✓	✓	✓	1447	77k	39	TV
How2Sign(Duarte et al., 2021)	ASL	✓	✓	✓	✗	79	16k	11	Lab
OpenASL(Shi et al., 2022)	ASL	✓	✗	✗	✗	288	33k	220	Web
YouTube-ASL(Uthus et al., 2023)	ASL	✓	✗	✗	✓	984	60k	>2519	Web
CSL-Daily(Zhou et al., 2021)	CSL	✓	✗	✗	✗	23	2k	10	Lab
SWISSTXT(Camgöz et al., 2021)	DSGS	✓	✓	✓	✓	88	-	-	TV
VRT-Raw(Camgöz et al., 2021)	VGT	✓	✓	✓	✓	100	-	-	TV
KETI(Ko et al., 2019)	KVK	✓	✓	✗	✗	29	419	14	Lab
SP-10(Yin et al., 2022)	various	✓	✓	✗	✗	14	17k	79	Web
AfriSign(Gueuwou et al., 2023)	various	✓	✗	✗	✗	152	20k	-	Web
PHOENIX2014T(Camgoz et al., 2018b)	DGS	✓	✗	✗	✗	11	3k	9	TV
Public DGS Corpus(Hanke et al., 2020)	DGS	✓	✗	✗	✗	50	-	-	TV
PHOENIX-News (ours)	DGS	✓	✓	✓	✓	486	190k	11	TV

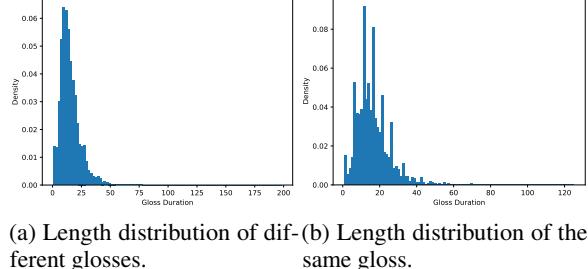
models to generate sign language, thereby improving generation speed (Huang et al., 2021; Hwang et al., 2021, 2022). Additionally, researchers have explored the generation of photo-realistic sign language videos using Generative Adversarial Networks (GANs) (Saunders et al., 2022) or diffusion models (Fang et al., 2023; Xie et al., 2024). Recent studies have shown that using 3D human models, such as SPML-x(Pavlakos et al., 2019), is a better choice for sign language understanding (Lee et al., 2023) and production tasks (Stoll et al., 2022). (Inan et al., 2022) found that representing the intensification level of glosses connected with the duration of a sign.

## 2.2 Vector Quantization for SLP

Vector Quantized Variational Autoencoders (VQ-VAE) proposed by (van den Oord et al., 2017) is an autoencoder structure that aims to learn a discrete representation of data. Recently, VQ-VAE has been used for the SLP task, such as (Saunders et al., 2021a) using a modified VQ-GAN for isolated word sign language video generation. Recently, VQ-VAE has been applied to the SLP task. For instance, (Xie et al., 2024) utilized a modified VQ-GAN (Esser et al., 2021) to generate isolated sign language videos. (Xie et al., 2023) employed VQ-VAE to generate sign pose sequences from gloss sequences. However, existing methods rely on fixed-length encodings and overlook the unequal distribution of information in sign language. To address this issue, we propose a pioneering approach: a variable-length dynamic vector quantization method specifically designed for sign language.

## 2.3 Sign Language Dataset

High-quality sign language datasets are crucial for the SLP task. Table 1 summarizes the publicly available datasets used for sign language research. The PHOENIX14T (Camgoz et al., 2018b) dataset is the most commonly used dataset for SLP tasks, but it has limited data. As an important supplement, we propose PHOENIX-News, which contains 486 hours of sign language data. To the best of our knowledge, this is the largest German sign language dataset to date.



(a) Length distribution of different glosses. (b) Length distribution of the same gloss.

## 3 Analyzing Information Density in Sign Language

The most commonly used discrete representation for sign language is glosses, which are the basic semantic units in sign language and are annotated by sign language experts. We first counted the length distribution of glosses in the PHOENIX14T dataset, as shown in Figure 2a. It can be seen that the length distribution of glosses is uneven, with most glosses having a length between 0 and 50, but some glosses have a length of more than 50. This indicates that uneven information density does exist in sign language. We then counted the length distribution of the most frequently occurring glosses (REGEN) in different contexts, as shown in Figure 2b. It can be seen that even the same gloss has differ-

ent lengths in different contexts. These analysis results inspire us to design a dynamic vector quantization method as described in subsection 4.2 and a duration transformer to predict duration based on context as described in subsection 4.3.

## 4 Method

Our overall two-stage framework is depicted in Figure 3, which consists of two stages: DVQ-VAE and T2M-GPT. In the following, we will first briefly revisit the formulation of VQ and then describe our proposed method in detail.

### 4.1 Preliminary

Vector quantization (VQ) (van den Oord et al., 2017) represents a technique for learning a codebook to encode sign language sequences into discrete code representations. Given a sign language sequence  $X = [x_1, x_2, \dots, x_T]$  with  $x_t \in \mathbb{R}^d$ , where  $T$  is the number of frames and  $d$  is the dimension of the sign language, we aim to recover the sign language sequence through an autoencoder and a learnable codebook containing  $K$  codes  $C = \{c_k\}_{k=1}^K$  with  $c_k \in \mathbb{R}^{d_c}$ , where  $d_c$  is the dimension of codes. The sign language sequence  $X$  is first encoded by the encoder  $E$  into a sequence of latent vectors  $Z = [z_1, z_2, \dots, z_{T/l}]$ , and  $z_t \in \mathbb{R}^{d_c}$ , where  $l$  represents the temporal downsampling rate of the encoder  $E$ . For fixed-length encoding,  $l$  is fixed, while for variable-length encoding,  $l$  is dynamically changing. For  $i$ -th latent feature  $z_i$ , the quantization through  $C$  is to find the most similar element in  $C$ , which can be properly written as:

$$\hat{z}_i = \arg \min_{c_k \in C} \|z_i - c_k\|_2 \quad (1)$$

### 4.2 Stage 1: Dynamic Vector Quantization VAE (DVQ-VAE)

Existing methods employ a fixed downsampling rate  $l$  for fixed-length encoding, neglecting the uneven information density in sign language. This oversight introduces redundancy in the learned codebook, leading to a decrease in both generation quality and speed. To address this issue, we propose DVQ-VAE, which consists of a dynamic encoder and a dynamic decoder.

**Dynamic Encoder.** As shown in Figure 3, the sign language sequence  $X$  first passes through a sign language embedding layer. Then, after adding

positional encoding information, it is input into a Transformer Encoder to obtain a sequence of latent vectors  $H = [h_1, h_2, \dots, h_T]$ , with  $h_t \in \mathbb{R}^{d_h}$ , where  $d_h$  is the dimension of the latent vectors. We formulate these operations as:

$$\begin{aligned} X'_t &= \text{relu}(\text{LN}(W_1 X_t + B_1)) + f_{\text{pos}}(t) \\ H &= \text{TransformerEncoder}(X') \end{aligned} \quad (2)$$

where LN denotes Layer Normalization (Ba et al., 2016),  $f_{\text{pos}}$  denotes the positional encoding function, and  $W_1 \in \mathbb{R}^{d_h \times d}$  and  $B_1 \in \mathbb{R}^{d_h}$  are learnable parameters.

The dynamic encoder then contains an **information-based adaptive downsampling module**, which adaptively adjusts the downsampling rate by considering the information weight of each frame. Specifically, we input the latent vector sequence  $H$  into a multi-layer perceptron (MLP) to obtain the information weight of each frame  $I = [i_1, i_2, \dots, i_T]$ , where  $i_t \in [0, 1]$ . We then segment the latent vector sequence  $H$  according to the information weight threshold  $O$  (we set it to 1.0) for semantic unit, and then perform weighted averaging within the segment to obtain the downsampled latent vector sequence  $Z = [z_1, z_2, \dots, z_{T/l}]$ . The downsampling process of the entire module can be formulated as:

$$I = \sigma(W_3(\text{relu}(W_2 H + B_2) + H) + B_3) \quad (3)$$

$$S = \text{cumsum}(I) // O \quad (4)$$

$$Z_t = \sum_{j=1}^T H_j \cdot I_j \cdot F_j, \quad D_t = \sum_{j=1}^T \cdot F_j \quad (5)$$

$$\text{where } F_j = \begin{cases} 1, & \text{if } S_j = t - 1 \\ 0, & \text{otherwise} \end{cases}$$

Equation 3 represents the operation of the MLP, where  $\sigma$  denotes the sigmoid activation function. Equation 4 represents the process of segmenting the latent vector sequence  $H$  according to the information weight threshold  $O$ , where  $S = [s_1, s_2, \dots, s_T]$  and  $s_t \in [0, \text{sum}(I) // O]$ , representing the position markers of the segments.  $\text{cumsum}$  denotes the cumulative sum function, and  $//$  denotes the integer division. Equation 5 represents the process of weighted downsampling, where  $Z_t$  denotes the downsampled latent vector and  $D_t$  denotes the duration of the current latent vector.

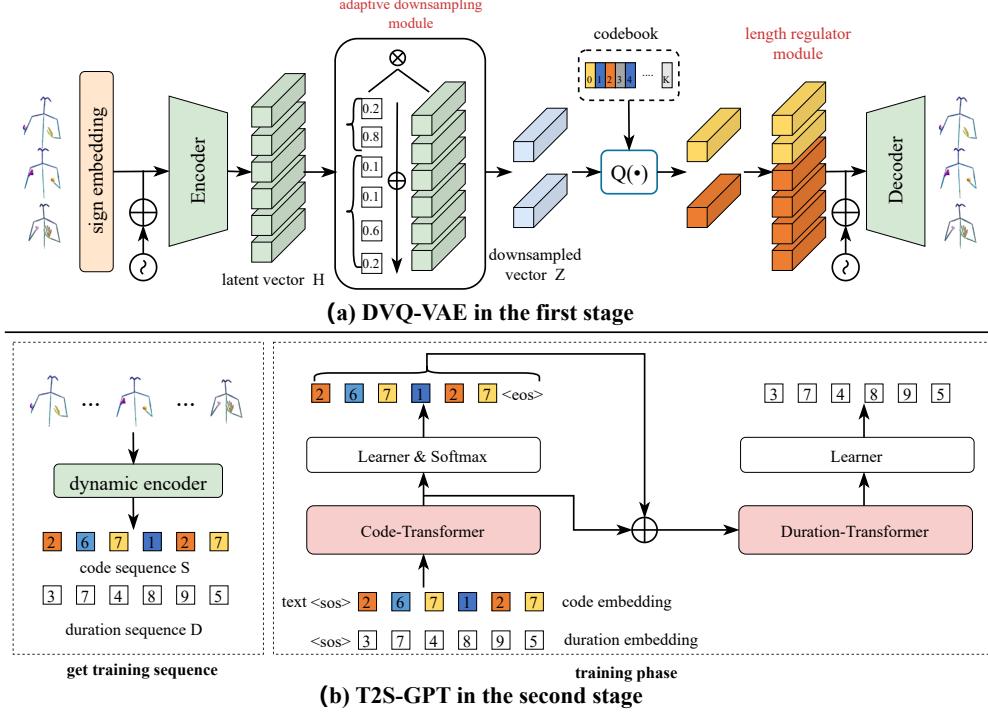


Figure 3: The overview of our proposed two-stage framework.

**Dynamic Decoder.** The goal of the dynamic decoder is to reconstruct the original sign language sequence  $X$  based on the quantized latent vector sequence  $\hat{Z}$  and the duration information  $D = [d_1, d_2, \dots, d_{T/l}]$ . We use a **length regulator module** to address the issue of mismatched lengths between the vector sequence  $\hat{Z}$  and the original sign language sequence  $X$  during dynamic decoding.

$$\hat{X} = \text{LR}(\hat{Z}, D) \quad (6)$$

where  $\hat{X}$  denotes the extended sequence, and LR denotes the length regulator module. For example, if  $\hat{Z} = [\hat{z}_1, \hat{z}_2, \hat{z}_3]$  and  $D = [1, 2, 3]$ , then  $\hat{X} = [\hat{z}_1, \hat{z}_2, \hat{z}_2, \hat{z}_3, \hat{z}_3, \hat{z}_3]$ . We then input the extended sequence  $\hat{X}$  into a Transformer-based decoder to obtain the reconstructed sign language sequence  $X_{re}$ .

**Training of DVQ-VAE.** The optimization goal of the original VQ-VAE (van den Oord et al., 2017)  $\mathcal{L}_{\text{vq}}$  contains three components: a reconstruction loss  $\mathcal{L}_{\text{re}}$ , an embedding loss  $\mathcal{L}_{\text{embed}}$ , and a commitment loss  $\mathcal{L}_{\text{commit}}$ .

$$\mathcal{L}_{\text{vq}} = \mathcal{L}_{\text{re}} + \underbrace{\|Z - sg[\hat{Z}]\|_2}_{\mathcal{L}_{\text{embed}}} + \lambda_1 \underbrace{\|sg[Z] - \hat{Z}\|_2}_{\mathcal{L}_{\text{commit}}} \quad (7)$$

where  $\lambda_1$  is a hyper-parameter for the commitment loss and  $sg$  is the stop-gradient operator. In our work, the calculation formula for the reconstruction loss is as follows:

$$\mathcal{L}_{\text{re}} = \mathcal{L}_1^{\text{smooth}}(X, X_{\text{re}}) + \mathcal{L}_1^{\text{smooth}}(V(X), V(X_{\text{re}})) \quad (8)$$

where  $V(\cdot)$  denotes the calculation of velocity, for example,  $V(X) = [v_1, v_2, \dots, v_{T-1}]$ , where  $v_i = x_{i+1} - x_i$ . In addition to the original optimization goal, we introduce two new loss functions: **budget loss**  $\mathcal{L}_{\text{budget}}$  and **sign language translation auxiliary loss**  $\mathcal{L}_{\text{slt}}$ . Without using the budget loss, the model tends to use more codes to represent the sign language sequence, resulting in longer sequence lengths. To encourage the model to use a higher downsampling rate  $l$ , we define the budget loss as:

$$\mathcal{L}_{\text{budget}} = \mathbb{E}[\max(0, (\text{sum}(I) - T/R))] \quad (9)$$

Since the length of the downsampled sequence is  $\text{sum}(I)/O$ , the budget loss can be interpreted as the expectation of the length of the downsampled sequence. Where  $T$  denotes the length of the original sign language sequence, and  $R$  denotes the expected downsampling rate. The goal of the sign

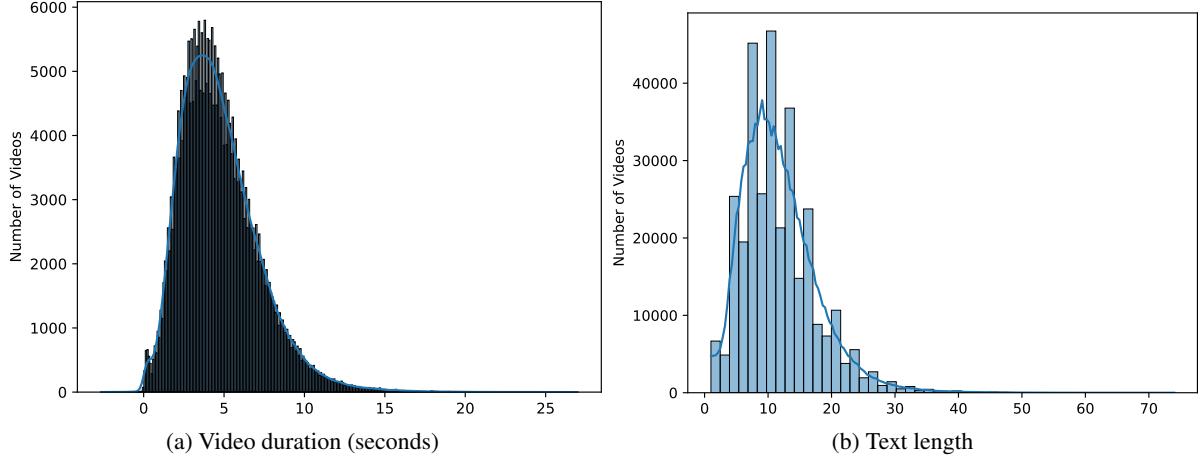


Figure 4: Distribution of text length and video duration in the PHOENIX-News dataset.

language translation auxiliary loss is to preserve the semantic information of the reconstructed sign language sequence, and its calculation formula is as follows:

$$\mathcal{L}_{\text{slt}} = \mathbb{E}[-\log P(Y|X_{re})] \quad (10)$$

where  $Y$  denotes the spoken language text corresponding to the sign language sequence. The final loss for DVQ-VAE is defined as:

$$\mathcal{L} = \mathcal{L}_{vq} + \lambda_2 \mathcal{L}_{\text{budget}} + \lambda_3 \mathcal{L}_{\text{slt}} \quad (11)$$

We also use two common training recipes (Razavi et al., 2019), exponential moving average (EMA) and code book restart, to improve the utilization of the codebook.

### 4.3 Stage 2: Text-to-Sign GPT (T2S-GPT)

**Code-Transformer.** With a learned DVQ-VAE, a sign language sequence  $X$  can be mapped to a sequence of indices  $S = [s_1, s_2, \dots, s_{T/l}, End]$ , which are indices from the learned codebook. Note that a special *End* token is added to indicate the stop of the sign language code sequence. By projecting  $S$  back to their corresponding codebook entries, we obtain  $\hat{Z} = [\hat{z}_1, \hat{z}_2, \dots, \hat{z}_{T/l}]$ , where  $\hat{z}_i = c_{s_i}$ . The generation of the sign language code sequence  $S$  can be formalized as an autoregressive next index prediction problem: given the previous  $i - 1$  indices, i.e.,  $S_{<i}$ , and the text condition  $Y$ , our goal is to predict the distribution of the possible next index  $p(S_i | Y, S_{<i})$ , which can be solved by a transformer, as shown in Figure 3. The negative log-likelihood (NLL) loss for code autoregressive training is:

$$\mathcal{L}_{\text{code}} = \mathbb{E}[-\log p(S_i | Y, S_{<i}, D_{<i})] \quad (12)$$

We introduce a duration embedding layer to embed the duration information  $D$  into the transformer.

**Duration-Transformer.** As mentioned in subsection 4.2 and Equation 6, to decode  $X_{re}$ , we need not only  $\hat{Z}$ , but also the duration information  $D$ . Therefore, we design a duration-transformer to predict the duration of the next code based on the previous code’s duration and the current codes. As shown in Figure 3, the duration-transformer takes the sum of Code-Transformer’s output hidden vector  $H_{\text{code}}$  and an extra code embedding as input:

$$H_{\text{dur}} = H_{\text{code}}[N_y : N_y + l - 1] + f_{\text{code}}(S[\leq l]) \quad (13)$$

where  $H_{\text{dur}}$  denotes the input of the duration-transformer, and  $N_y$  denotes the length of the condition text. The design idea behind this is that when predicting the next code’s duration, the model should not only be aware of previous steps’ codes and their duration information but also should be aware of current code information. The optimization goal of the duration-transformer is to minimize the difference between the predicted duration and the real duration. The calculation formula is as follows:

$$\mathcal{L}_{\text{dur}} = \mathbb{E}[\|D_i - \hat{D}_i\|_2] \quad (14)$$

In inference, we round the output of the duration-transformer to obtain the duration. The final optimization goal is:

$$\mathcal{L} = \mathcal{L}_{\text{code}} + \mathcal{L}_{\text{dur}} \quad (15)$$

## 5 The Proposed PHOENIX-News Dataset

As shown in Table 1, PHOENIX-News aims to provide the community with a new large-scale

Table 2: Quantitative results for text to sign language task on PHOENIX14T test set.

Methods	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	BLEU-4
GT	39.17	37.75	24.92	18.25	14.34
PT( <a href="#">Saunders et al., 2020b</a> )	20.58	17.47	7.76	5.50	4.38
NAT-EA( <a href="#">Huang et al., 2021</a> )	26.81	27.00	14.12	9.20	6.67
T2M-GPT ( <a href="#">Zhang et al., 2023b</a> )	29.19	28.32	16.05	10.77	8.01
MDM ( <a href="#">Tevet et al., 2022</a> )	30.37	27.59	15.83	10.29	7.55
<b>T2S-GPT (ours)</b>	<b>34.65</b>	<b>33.16</b>	<b>21.09</b>	<b>15.26</b>	<b>11.87</b>

document-level sign language dataset, which contains 486 hours of sign language videos, audio, and transcription texts. We collected daily news programs in German Sign Language from the German public television station PHOENIX from 2013 to 2023. We then used whisper ([Radford et al., 2023](#)) to transcribe the program’s speech into text. Finally, we performed preprocessing steps such as domain cropping, sign language pose estimation, and sign language text alignment to obtain the final dataset. Since there are new sign language news programs every day, PHOENIX-News will be a continuously updated project. Therefore, we did not divide the dataset into training and test sets, but used the existing dataset for testing. Each video in the dataset has an average duration of 4.7 seconds, and the average text length is 11 words. We show the distribution of video duration and text length in [Figure 4a](#) and [Figure 4b](#).

## 6 Experiments

### 6.1 Experimental Setup

We will introduce our experimental setup in this section, including the dataset and evaluation metrics. We will provide the implementation details of the model in the [Appendix C](#).

**Dataset and Evaluation Metrics.** We evaluate our proposed T2S-GPT model on the PHOENIX14T ([Camgoz et al., 2018b](#)) dataset, which is the most commonly used dataset for SLP tasks and has been used as a benchmark in many previous SLP works ([Huang et al., 2021; Saunders et al., 2021a, 2020b; Xie et al., 2023](#)). The PHOENIX14T dataset contains 7,096 training samples (with 2,887 words in German spoken language translations), 519 validation samples, and 642 test samples. We use the pose parameter  $\theta$  in the SMPL-X model to represent the sign language pose, and the rotation 6D representation ([Zhou et al., 2019](#)) is used to represent the rotation in the pose. Follow-

ing the most widely used setting in SLP ([Saunders et al., 2020b](#)), we use the back translation metric to evaluate the generation quality. Since previous works did not publicly release the weights of their SLT models used to calculate the back translation metric, following the previous setting, we train an SLT model using the code from ([Camgoz et al., 2020](#)). We provide the details of the sign language representation in the [Appendix B](#). We provide the training details of the SLT model in the [Appendix D](#).

### 6.2 Comparisons with State-of-the-Art Methods

We compare our T2S-GPT model with several other models, including the state-of-the-art text-to-sign model and the state-of-the-art text-to-motion model.

**Comparison methods.** 1) Progression Transformer (PT) ([Saunders et al., 2020b](#)) directly predicts the sign language pose sequence in an autoregressive manner. 2) NAT-EA ([Huang et al., 2021](#)) generates the sign language pose sequence in a non-autoregressive manner. 3) T2M-GPT ([Zhang et al., 2023b](#)) is a state-of-the-art autoregressive text-to-motion model, and its prediction target is the discrete representation of sign language processed by VQ-VAE. 4) MDM ([Tevet et al., 2022](#)) uses a diffusion model to generate motion sequences based on text in a non-autoregressive manner. Both T2M-GPT and MDM use CLIP ([Radford et al., 2021](#)) to extract text features as condition signals, but the original CLIP does not support German well. To make a fair comparison, we use a multilingual CLIP model([Reimers and Gurevych, 2019](#))<sup>2</sup>.

**Quantitative Comparison.** We report the back translation metrics (including BLEU ([Papineni et al., 2002](#)) scores and ROUGE-L ([Lin and](#)

<sup>2</sup><https://huggingface.co/sentence-transformers/clip-ViT-B-32-multilingual-v1>

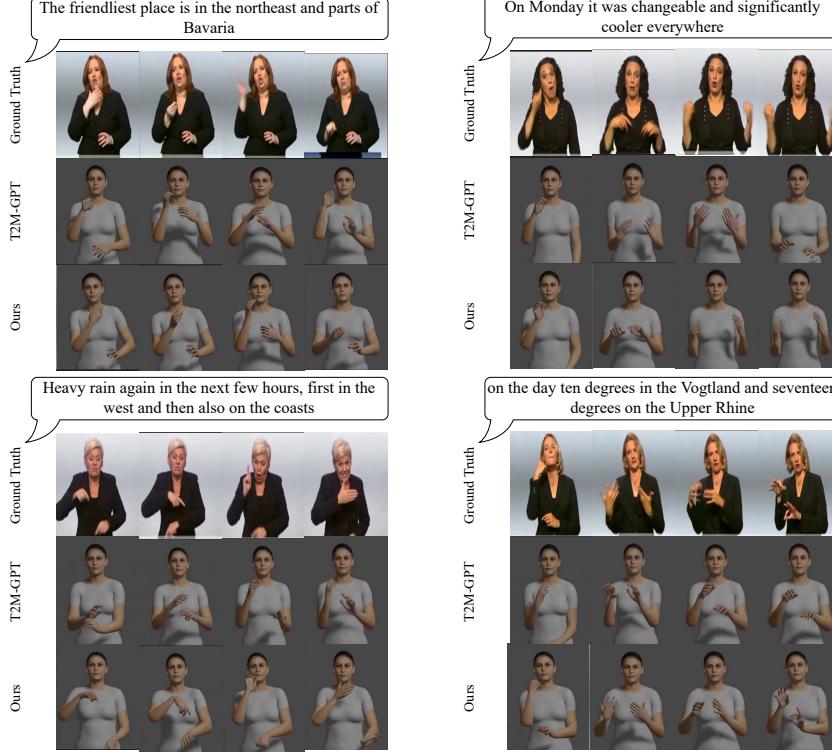


Figure 5: Qualitative results of our T2S-GPT model.

Och, 2004) scores) obtained by all models on the PHOENIX14T dataset in Table 2. We only use the PHOENIX14T dataset to train all models, and the training settings are consistent with those in the original papers. As shown in Table 2, our T2S-GPT model achieves the best results on all metrics. Specifically, our T2S-GPT model achieves a score of 11.87 on BLEU-4, which is 3.86 points higher than the state-of-the-art T2M-GPT model. On ROUGE-L, our T2S-GPT model achieves a score of 34.65, which is 4.28 points higher than the state-of-the-art MDM model. These results indicate that our T2S-GPT model can generate higher-quality sign language.

**Qualitative Results.** We qualitatively compare our T2S-GPT method with other methods and the ground truth sign language pose sequence on the PHOENIX14T test set, as shown in Figure 5. As shown in Figure 5, compared with other methods, the sign language generated by our T2S-GPT method is closer to the ground truth. Note that we provide a video demonstration of our method on the anonymous project homepage <https://t2sgpt-demo.yinaoxiong.cn/>, which can better convey the temporal information.

Table 3: Results of ablation experiments on the PHOENIX14T dataset.

Model	R	B1	B2	B3	B4
T2S-GPT	<b>34.65</b>	<b>33.16</b>	<b>21.09</b>	<b>15.26</b>	<b>11.87</b>
w/o DVQ-VAE	30.80	27.77	16.01	10.96	8.39
w/o Duration-Transformer	31.99	30.05	18.25	12.43	9.39

### 6.3 Ablation and Analysis

**Analysis on DVQ-VAE.** As shown in the second row of Table 3, when we replace DVQ-VAE with the VQ-VAE proposed by (Zhang et al., 2023b) with a downsampling rate of 4, we find that the back translation metrics of the SLP model have decreased significantly. This indicates that our DVQ-VAE model can obtain more compact and higher-quality discrete representations of sign language.

**Analysis on Duration-Transformer.** As shown in the third row of Table 3, when we replace the duration-transformer with a simple fully connected layer, we find that the back translation metrics of the SLP model have decreased significantly.

**Impact of dataset size.** To study whether our proposed T2S-GPT model is scalable, we train the T2S-GPT model by adding different proportions of the PHOENIX-News dataset. The experimental results are shown in , where we find that as

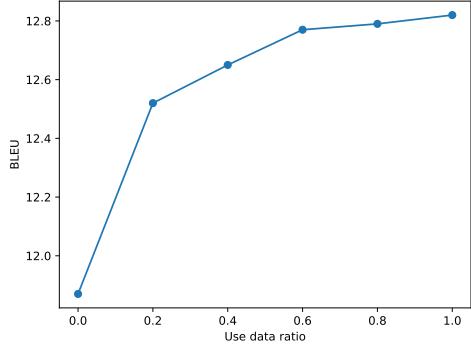


Figure 6: Impact of dataset size on the performance of T2S-GPT.

the dataset size increases, the performance of the T2S-GPT model also continues to improve. This indicates that our T2S-GPT model is scalable.

## 7 Conclusion

In this work, we propose a two-stage text-to-sign model T2S-GPT, which consists of a dynamic vector quantization VAE (DVQ-VAE) and a GPT-like autoregressive generation model. Our method achieves better performance than the previous state-of-the-art text-to-sign model. In addition, we have collected a new large-scale document-level sign language dataset PHOENIX-News, and the experimental results show that a larger dataset can still bring additional improvements to our method.

## 8 Limitations and Potential Risks

Although using a 3D human body model as the sign language representation introduces prior information about human body shape, it does not constrain the rotation motion of the joints themselves. The model’s predictions occasionally produce some abnormal cases that do not conform to the human joint structure, which may make users feel uncomfortable. At the same time, this is also a manifestation of the model’s generation errors. To address this issue, we plan to introduce more prior information in future work, such as human motion priors and physical constraints on human joint rotation angles. SLP technology itself does not have any obvious potential risks, but since the current SLP technology is still in a relatively early stage, if it is directly applied to practical scenarios, it may mislead users. For example, in weather forecasts, if the model generates sign language with incorrect place names, it may mislead users.

## Acknowledgements

This work was supported by the Key Research and Development Projects in Zhejiang Province (No. 2024C01106), the NSFC (No. 62272411), the National Key Research and Development Project of China (2018AAA0101900), and Research funding from FinVolution Group.

## References

- Samuel Albanie, GÜl Varol, Liliane Momeni, Hannah Bull, Triantafyllos Afouras, Himel Chowdhury, Neil Fox, Bencie Woll, Rob Cooper, Andrew McParland, et al. 2021. Bbc-oxford british sign language dataset. *arXiv preprint arXiv:2111.03635*.
- Tenglong Ao, Qingzhe Gao, Yuke Lou, Baoquan Chen, and Libin Liu. 2022. Rhythmic gesticulator: Rhythm-aware co-speech gesture synthesis with hierarchical neural embeddings. *ACM Transactions on Graphics (TOG)*, 41(6):1–19.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer Normalization.
- Vasileios Baltatzis, Rolando Alexandros Potamias, Evangelos Ververas, Guanxiong Sun, Jiankang Deng, and Stefanos Zafeiriou. 2023a. Neural sign actors: A diffusion model for 3d sign language production from text. *arXiv preprint arXiv:2312.02702*.
- Vasileios Baltatzis, Rolando Alexandros Potamias, Evangelos Ververas, Guanxiong Sun, Jiankang Deng, and Stefanos Zafeiriou. 2023b. Neural Sign Actors: A diffusion model for 3D sign language production from text.
- Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. 2018a. Neural sign language translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7784–7793.
- Necati Cihan Camgoz, Simon Hadfield, Oscar Koller, Hermann Ney, and Richard Bowden. 2018b. Neural Sign Language Translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7784–7793.
- Necati Cihan Camgoz, Oscar Koller, Simon Hadfield, and Richard Bowden. 2020. Sign Language Transformers: Joint End-to-End Sign Language Recognition and Translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10023–10033.
- Necati Cihan Camgoz, Ben Saunders, Guillaume Ruchette, Marco Giovanelli, Giacomo Inches, Robin Nachtrab-Ribback, and Richard Bowden. 2021. Content4all open research sign language translation datasets. In *2021 16th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2021)*, pages 1–5. IEEE.

- Stephen Cox, Michael Lincoln, Judy Tryggvason, Melanie Nakisa, Mark Wells, Marcus Tutt, and Sanja Abbott. 2002. Tessa, a system to aid communication with deaf people. In *Proceedings of the fifth international ACM conference on Assistive technologies*, pages 205–212.
- Amanda Duarte, Shruti Palaskar, Lucas Ventura, Deepti Ghadiyaram, Kenneth DeHaan, Florian Metze, Jordi Torres, and Xavier Giro-i-Nieto. 2021. **How2Sign: A Large-scale Multimodal Dataset for Continuous American Sign Language**. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2734–2743.
- Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming Transformers for High-Resolution Image Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12873–12883.
- Sen Fang, Chunyu Sui, Xuedong Zhang, and Yapeng Tian. 2023. **SignDiff: Learning Diffusion Models for American Sign Language Production**.
- Shester Gueuwou, Kate Takyi, Mathias Müller, Marco Stanley Nyarko, Richard Adade, and Rose-Mary Owusu Mensah Gyening. 2023. Afrisign: Machine translation for african sign languages. In *4th Workshop on African Natural Language Processing*.
- Thomas Hanke, Marc Schulder, Reiner Konrad, and Elena Jahn. 2020. Extending the Public DGS Corpus in Size and Depth. In *Proceedings of the LREC2020 9th Workshop on the Representation and Processing of Sign Languages: Sign Language Resources in the Service of the Language Community, Technological Challenges and Application Perspectives*, pages 75–82, Marseille, France. European Language Resources Association (ELRA).
- Mengqi Huang, Zhendong Mao, Zhuowei Chen, and Yongdong Zhang. 2023. Towards Accurate Image Coding: Improved Autoregressive Image Generation With Dynamic Vector Quantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22596–22605.
- Wencan Huang, Wenwen Pan, Zhou Zhao, and Qi Tian. 2021. **Towards Fast and High-Quality Sign Language Production**. In *Proceedings of the 29th ACM International Conference on Multimedia, MM ’21*, pages 3172–3181, New York, NY, USA. Association for Computing Machinery.
- Eui Jun Hwang, Jung Ho Kim, Suk Min Cho, and Jong C. Park. 2022. **Non-Autoregressive Sign Language Production via Knowledge Distillation**.
- Eui Jun Hwang, Jung-Ho Kim, and Jong C. Park. 2021. Non-autoregressive sign language production with gaussian space. In *32nd British Machine Vision Conference 2021, BMVC 2021, Online, November 22-25, 2021*, page 197. BMVA Press.
- Mert Inan, Yang Zhong, Sabit Hassan, Lorna Quandt, and Malihe Alikhani. 2022. Modeling intensification for sign language generation: A computational approach. *arXiv preprint arXiv:2203.09679*.
- Sang-Ki Ko, Chang Jo Kim, Hyedong Jung, and Choongsang Cho. 2019. Neural sign language translation based on human keypoint estimation. *Applied sciences*, 9(13):2683.
- Taeryung Lee, Yeonguk Oh, and Kyoung Mu Lee. 2023. Human Part-wise 3D Motion Context Learning for Sign Language Recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20740–20750.
- Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 605–612.
- Ilya Loshchilov and Frank Hutter. 2018. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- John McDonald, Rosalee Wolfe, Jerry Schnepp, Julie Hochgesang, Diana Gorman Jamrozik, Marie Stumbo, Larwan Berke, Melissa Bialek, and Farah Thomas. 2016. An automated technique for real-time production of lifelike animations of american sign language. *Universal Access in the Information Society*, 15:551–566.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. 2019. Expressive Body Capture: 3D Hands, Face, and Body From a Single Image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10975–10985.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2023. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR.
- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.

- Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. 2019. Generating diverse high-fidelity images with vqvae-2. *Advances in neural information processing systems*, 32.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Ben Saunders, Richard Bowden, and Necati Cihan Camgöz. 2020a. Adversarial training for multi-channel sign language production. In *31st British Machine Vision Conference 2020, BMVC 2020, Virtual Event, UK, September 7-10, 2020*. BMVA Press.
- Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. 2020b. Progressive transformers for end-to-end sign language production. In *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XI*, volume 12356 of *Lecture Notes in Computer Science*, pages 687–705. Springer.
- Ben Saunders, Necati Cihan Camgöz, and Richard Bowden. 2021a. Continuous 3D multi-channel sign language production via progressive transformers and mixture density networks. *International Journal of Computer Vision*, 129(7):2113–2135.
- Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. 2021b. Mixed SIGNals: Sign Language Production via a Mixture of Motion Primitives. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1899–1909.
- Ben Saunders, Necati Cihan Camgoz, and Richard Bowden. 2022. Signing at Scale: Learning to Co-Articulate Signs for Large-Scale Photo-Realistic Sign Language Production. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5141–5151.
- Bowen Shi, Diane Brentari, Greg Shakhnarovich, and Karen Livescu. 2022. Open-Domain Sign Language Translation Learned from Online Video.
- Stephanie Stoll, Necati Cihan Camgöz, Simon Hadfield, and Richard Bowden. 2020. Text2Sign: Towards sign language production using neural machine translation and generative adversarial networks. *International Journal of Computer Vision*, 128(4):891–908.
- Stephanie Stoll, Armin Mustafa, and Jean-Yves Guillemaut. 2022. There and Back Again: 3D Sign Language Generation from Text Using Back-Translation. In *2022 International Conference on 3D Vision (3DV)*, pages 187–196.
- Guy Tevet, Sigal Raab, Brian Gordon, Yonatan Shafir, Daniel Cohen-Or, and Amit H. Bermano. 2022. Human Motion Diffusion Model.
- David Uthus, Garrett Tanzer, and Manfred Georg. 2023. YouTube-ASL: A Large-Scale, Open-Domain American Sign Language-English Parallel Corpus.
- Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. 2017. Neural Discrete Representation Learning. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Will Williams, Sam Ringer, Tom Ash, David MacLeod, Jamie Dougherty, and John Hughes. 2020. Hierarchical quantized autoencoders. *Advances in Neural Information Processing Systems*, 33:4524–4535.
- Pan Xie, Taiying Peng, Yao Du, and Qipeng Zhang. 2024. Sign Language Production With Latent Motion Transformer. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3024–3034.
- Pan Xie, Qipeng Zhang, Taiyi Peng, Hao Tang, Yao Du, and Zexian Li. 2023. G2P-DDM: Generating Sign Pose Sequence from Gloss Sequence with Discrete Diffusion Model.
- Aoxiong Yin, Zhou Zhao, Weike Jin, Meng Zhang, Xingshan Zeng, and Xiaofei He. 2022. MLSLT: Towards multilingual sign language translation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 5099–5109. IEEE.
- Aoxiong Yin, Zhou Zhao, Jinglin Liu, Weike Jin, Meng Zhang, Xingshan Zeng, and Xiaofei He. 2021. Simultslt: End-to-end simultaneous sign language translation. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 4118–4127.
- Aoxiong Yin, Tianyun Zhong, Li Tang, Weike Jin, Tao Jin, and Zhou Zhao. 2023. Gloss attention for gloss-free sign language translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2551–2562.
- Biao Zhang, Mathias Müller, and Rico Sennrich. 2023a. Sltnet: A simple unified model for sign language translation. *arXiv preprint arXiv:2305.01778*.
- Jianrong Zhang, Yangsong Zhang, Xiaodong Cun, Shaoli Huang, Yong Zhang, Hongwei Zhao, Hongtao Lu, and Xi Shen. 2023b. T2M-GPT: Generating Human Motion from Textual Descriptions with Discrete Representations.
- Hao Zhou, Wengang Zhou, Weizhen Qi, Junfu Pu, and Houqiang Li. 2021. Improving Sign Language Translation With Monolingual Data by Sign Back-Translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1316–1325.
- Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. 2019. On the Continuity of Rotation Representations in Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753.

## A Example Appendix

## B Sign Language Representation

Inspired by the latest advances in sign language processing (Lee et al., 2023; Stoll et al., 2020), to better represent the complex body movements in sign language, we propose to use the pose parameter  $\vec{\theta} = [\vec{\omega}_0^T, \dots, \vec{\omega}_K^T]^T$  of the SMPL-X human body model (Pavlakos et al., 2019) as the sign language representation, instead of the 3D joint coordinates in Euclidean space used in previous works. Where  $\vec{\omega}_k \in \mathbb{R}^3$  denotes the axis-angle representation of the relative rotation of part  $k$  with respect to its parent in the kinematic tree. However, since the axis-angle form is not a continuous rotation representation, which is not conducive to network learning, we further convert it to the rotation 6D representation (Zhou et al., 2019)  $\vec{o} = [\vec{r}_0^T, \dots, \vec{r}_K^T]^T$ . We ignore the lower body joints outside the visible range. There are three advantages of using this representation: 1) it has rotation and translation invariance; 2) it separates the modeling of human body shape and pose, and the semantics of sign language should only be related to the pose and independent of the shape; 3) the introduction of human body prior avoids generating abnormal results, such as fingers longer than arms.

## C Implementation Details

For DVQ-VAE, we set the dimension of the latent vectors  $d_h$  to 512, the dimension of the code-book  $d_c$  to 512, and the number of codes  $K$  to 1024. The number of transformer layers in the encoder and decoder is set to 6, and the hidden size, number of heads, and feed-forward dimension for each layer are set to 512, 8, and 2048, respectively. The dropout rate is set to 0.1. We use AdamW (Loshchilov and Hutter, 2018) optimizer with  $[\beta_1, \beta_2] = [0.9, 0.99]$ , batch size of 256, and exponential moving constant  $\lambda = 0.99$ . We train for a total of 100K iterations, with an initial learning rate of 2e-4, and then use the cosine learning rate decay strategy during training.  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  in the final loss are set to 1, 0.5, and 1.0, respectively. The  $R$  in  $\mathcal{L}_{budget}$  is set to 12.

For T2S-GPT, the hidden size, number of heads, and feed-forward dimension for each transformer layer are set to 1024, 16, and 4096, respectively. The dropout rate is set to 0.1. The number of transformer layers in the code-Transformer and duration-Transformer is set to 18 and 6, respectively. We

use a batch size of 256 and train for 300K iterations. We optimize the models with the AdamW optimizer, warm up the learning rate for the first 4k updates to a peak of 1e-4, and then linearly decay it to 0. We use a 32GB NVIDIA V100 GPU to train our model.

## D Back Translation Model

To calculate the back translation metric, we train a sign language translation (SLT) model that takes sign language pose sequences as input and outputs the corresponding spoken language text. The SLT model adopts the architecture introduced by (Camgoz et al., 2020). Both the encoder and decoder components of the model are built using transformers. In particular, the hidden size, number of heads, and feed-forward dimension for each layer are configured as 512, 8, and 2048, respectively. Additionally, a dropout rate of 0.4 is applied within the model. The number of transformer layers in the encoder and decoder is set to 3. The training settings are consistent with those in the original paper.