

WavBriVL: Robust Audio Representation and Generation of Audio Driven Diffusion Models

Anonymous EMNLP submission

Abstract

Multimodal large models have been recognized for their advantages in various performance and downstream tasks. The development of these models is crucial towards achieving general artificial intelligence in the future. In this paper, we propose a novel audio representation learning method called WavBriVL, which is based on Bridging-Vision-and-Language (BriVL). WavBriVL embeds audio, image, and text into a shared space, enabling the realization of various multimodal applications. Our approach addresses major challenges in robust audio representation learning and effectively captures the correlation between audio and image. Additionally, we demonstrate the qualitative evaluation of the generated images from WavBriVL, which serves to highlight the potential of our approach in creating images from audio. Overall, our experimental results demonstrate the efficacy of WavBriVL in downstream tasks and its ability to generate appropriate images from audio. The proposed approach has the potential for various applications such as speech recognition, music signal processing, and captioning systems. We would like to highlight that WavBriVL is the first universal method for generating images from audio-driven diffusion models.

1 Introduction

Sound and vision affect people’s core cognition in many areas, such as feeling, information processing and communication. Sound and vision are closely related. However, most of the existing methods only have a single cognitive ability, and some only study text-vision, text-voice, etc. Recent studies have shown that leveraging large-scale Internet data for self-supervised pre-training of models offers better results than relying on high-quality or manually labeled data sets (Pan et al., 2022), such as the recently amazing chatGPT¹. Moreover, mul-

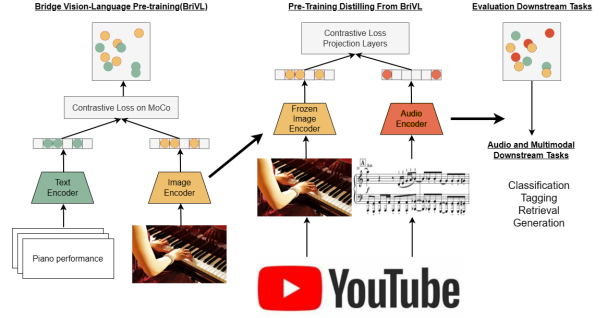


Fig. 1: Bridge Vision-Language Pre-training (BriVL), and our two-stage approaches including pre-training and evaluation.

iple studies demonstrate the effectiveness of multimodal models over single or bimodal models in several fields and tasks (Chen et al., 2022a), such as Microsoft’s latest BEiT3 model (Wang et al., 2022).

Data volume is the basic element for training large-scale language models. Since BERT of Devlin et al. (2018) (perhaps even earlier (Ma and Zhang, 2015)), the pre-training model of NLP has been benefiting from large-scale corpora. According to theory of OpenAI (Kaplan et al., 2020), the language model gradually reflects a scaling law (the rule that the model capacity increases with the model volume). In supervised learning, manual annotation of large amounts of data is very expensive, so self-supervised learning is valued for large model training. In order to expand the boundary of the research field and break the limitation of the lack of relevant resources (Hsu et al., 2021), we explore a new multimodal self-monitoring model based on the latest excellent work: **Bridging-Vision-and-Language** (Fei et al., 2022). It’s a new effort similar to OpenAI CLIP (Radford et al., 2021) and Google ALIGN (Jia et al., 2021). BriVL² model has excellent effect on image and text retrieval tasks, surpassing other common

¹<https://chat.openai.com/>

²<https://github.com/BAAI-WuDao/BriVL>

multimodal pre-training models in the same period.

In this work, we propose WavBriVL, an audio-visual correspondence model that extracts training from the BriVL model. The principle of WavBriVL is to freeze the BriVL visual model, run video on the visual stream of the model, and train a new model to predict BriVL embedding independently from the audio stream. Our method is very simple, which can train, expand, and output images. In addition, WavBriVL embeddings originate from BriVL, which means they align with text. In theory, this makes audio guided image repair (Zhao et al., 2022), audio subtitles and cross mode text/audio to audio retrieval be true. Our approach addresses the unique challenges presented by audio input to improve the overall performance and applicability of multimodal models in various tasks, such as media captioning and recommendation systems. And in this work, we conducted detailed testing on these possible tasks, WavBriVL generates audio subtitles, and cross-mode text/audio to audio retrieval, which surpasses previous methods’ outcomes, proving its significance in future multimodal models’ development.

Finally, we use WavBriVL to guide the generation of model DALL·E³ (Ramesh et al., 2021) output images, and intuitively verify that the embedded space is meaningful. Experimental results show that this method can effectively generate appropriate images from audio. Our approach enables the generation of high-quality images solely from the audio input by leveraging the shared embeddings in the BriVL framework. This is a significant contribution to the field of multimodal learning, as prior methods mainly focused on generating images from text or image inputs, rather than audio inputs. The novel contribution of our work, being the first method that audio-driven diffusion model to generate images, thus demonstrates the potential of the approach in advancing the field of multimodal learning. In addition, compared with other fully supervised models, WavBriVL theoretically requires less data to obtain competitive performance in downstream tasks, that is, it performs pre-training more effectively than competitive methods, because it does not need to completely re learn the visual model, only needs to train the audio model. It is a reproducible and potential application model, and we will provide more code information after publication.

³<https://github.com/lucidrains/DALLE-pytorch>

2 Method And Tasks

Bridging-Vision-and-Language (BriVL) is a model trained on 650 million text image weak semantic datasets. They designed a cross modal comparison learning algorithm based on the monomodal comparison learning method MoCo (He et al., 2020), and maintained the negative sample queue in different training batches through a mechanism called Memory Bank, so as to obtain a large number of negative samples for use in the comparison learning method. As shown in the left part of Figure 1, its core idea is to realize the general artificial intelligence model (AGI) by simulating the multi-mode processing idea of human brain. It also shows the SOTA results in such scenes as image annotation, image zero sample classification, and input features of other downstream multimodal tasks. Even the guidance generation model has excellent performance.

As shown in Figure 1, WavBriVL replaces the text encoder with the audio encoder by freezing the visual model of BriVL, runs the image through it, and trains the new model to predict that only the matching image embedded content is obtained from the audio. We refer to the exclusive multilayer perceptron of BriVL, which can not only enhance performance but also prepare for possible downstream tasks. After the audio encoder is trained, we freeze it and use it in the WavBriVL image generation task as a qualitative evaluation of our experimental results.

2.1 Dataset for WavBriVL performance test

We select diverse set of data ranging from various number of clips, number of categories, and perform diverse tasks including classification, retrieval, and generation. For evaluation, we use relevant metrics detailed in Table 1 for each task. BriVL needs more than 100 A100 graphics cards to train for 10 days, so we don’t consider retraining it. Our training and performance testing are based on the pre-trained model.

2.2 Dataset for diffusion model

We used VGG-Sound (Chen et al., 2020a) and AudioSet (Gemmeke et al., 2017) video datasets. VGG-Sound consists of short clips of audio sound, which includes 310 video classes and 200,000 audio that span a large number of challenging acoustic environments and noise characteristics of practical applications. All videos are non man-made and

more corresponding. We randomly select one image from each sample video, cut them into squares, and sample them down to 64×64 . The audio sampling rate is 16,000Hz. We use it to train the model, which helps to increase the applicability of the model. AudioSet collected 10 second clips from 2.1 million videos. We randomly selected the audio of 18 videos for our image generation task.

2.3 Feature extraction processing methods

For image and audio encoders, we use EfficientNet-B7 (Tan and Le, 2019) as the CNN in the image encoder, and the backbone WavLM (Chen et al., 2022b) as the basic transformer in the audio encoder. The self concerned block is composed of 4 Transformer encoder layers and MLP block respectively, with two fully connected layers and one ReLU activation layer. For all models, we use grid search to find the best hyperparameter. For other hyperparameters (such as batch size, training steps, learning rate, etc.), we directly use the suggested values in the original papers. Note that for per-instance perturbation, we adopt the appropriate quantity compared to the original epochs.

Image Encoder. In the input image, the method of BriVL using random grayscale for the input image and random color jitter for data enhancement is followed. For all videos in the dataset, we use 720P resolution and separate images (if not, use 480P). All images are cropped down to 360×360 pixels. We use Transformer to capture patch features, and use the average pooling layer to fuse and extract. To better capture the relationship of image patch features, BriVL’s team⁴ deploys a self-attention (SA) block containing multiple Transformer encoder layers. Every Transformer encoder layer consists of a multi-head attention (MHA) layer and a feed forward network (FFN) layer (Fei et al., 2022):

$$\mathbf{S}' = \text{LayerNorm}(\mathbf{S} + \text{MHA}(\mathbf{S})) \quad (1)$$

$$\mathbf{S} = \text{LayerNorm}(\mathbf{S}' + \text{FFN}(\mathbf{S}')) \quad (2)$$

Then, they use the average pooling layer to fuse the extracted patch features:

$$\mathbf{r}^{(i)} = \frac{1}{N_p} \sum_{j=1}^{N_p} \mathbf{S}_j \in \mathbb{R}^c \quad (3)$$

where \mathbf{S}_j is the j -th column of \mathbf{S} . A two-layer MLP block with a ReLU activation layer is adopted

to project $\mathbf{r}^{(i)}$ to the joint cross-modal embedding space, resulting in the final d -dimensional image embedding $\mathbf{z}^{(i)} \in \mathbb{R}^d$.

Audio Encoder. For audio input, we first convert the original audio waveform (1D) into a spectrum (2D) as the input of WavLM, and pool the entire 512 dimensional audio sequence to output an embedding. The WavLM embedding is computed by the weighted average of outputs from all transformer layers. The WavLM⁵ model inspired by HuBERT contains two main networks as follows: a CNN encoder and a Transformer with L blocks. During training, some frames of the CNN encoder output \mathbf{x} are masked randomly and fed to the Transformer as input. The Transformer is optimized to predict the discrete target sequence \mathbf{z} , in which each $z_t \in [C]$ is a C -class categorical variable. The distribution over the classes is parameterized with

$$p(c|\mathbf{h}_t) = \frac{\exp(\text{sim}(\mathbf{W}^P \mathbf{h}_t^L, \mathbf{e}_c)/\tau)}{\sum_{c'=1}^C \exp(\text{sim}(\mathbf{W}^P \mathbf{h}_t^L, \mathbf{e}_{c'})/\tau)} \quad (4)$$

where \mathbf{W}^P is a projection matrix, \mathbf{h}_t^L is the output hidden state for step t , \mathbf{e}_c is the embedding for class c , $\text{sim}(a, b)$ means the cosine similarity between a and b , and $\tau = 0.1$ scales the logit (Chen et al., 2022b). The WavLM embedding is calculated by the weighted average of all transformer layer outputs of WavLM, where the weights are learned during fine tuning. In the process of fine-tuning, we either update or freeze the parameters of WavLM.

2.4 Training process

We continue to use a similar cross modal comparative loss in BriVL. It is defined based on MoCo (He et al., 2020), which provides a mechanism of building dynamic sample queues for contrastive learning. Since the two negative queues used in our BriVL decouple the queue size from the mini-batch size, we can have a much larger negative sample size than the mini-batch size (thus GPU-resource-saving). Loss function with cross projection defined as $CXLoss = L(f(\text{Image}), \text{Audio}) + L(\text{Image}, g(\text{Audio}))$ (f, g : projection functions and L : contrastive loss).

For all models, we use grid search to find the best hyperparameter. For other hyperparameters (such as batch size, training steps, learning rate,

⁴<https://github.com/BAAI-WuDao/BriVL>

⁵<https://github.com/microsoft/unilm/tree/master/wavlm>

etc.), we directly use the suggested values in the original papers. Note that for per-instance perturbation, we adopt the appropriate quantity compared to the original epochs. In this paper, we utilize several key parameters to achieve our experimental results. The topk parameter is set to 1, which indicates that we only consider the top-scoring prediction for each input instance. The queue_size parameter is set to 9600, which controls the number of instances that can be processed in parallel. We use a momentum value of 0.99 to stabilize the learning process and prevent oscillations during training. The temperature parameter is set to 0.07, which scales the logits output of the model to control the softness of the predicted probability distribution. Finally, we use a grid_size of 4 to divide the input image into a grid of smaller sub-regions for object detection tasks.

3 Related Works

Our motivation comes from the relevant work proposed in the first half of this year (2022): we can see that BriVL has demonstrated better performance than CLIP (Radford et al., 2021) in many aspects, and Microsoft’s new WavLM (Chen et al., 2022b) is also better than the previous Wav2Vec (Baevski et al., 2020) in most cases. We guess that the combination of these two new works will also be better than Wav2CLIP⁷. More importantly, there is currently a lack of groundbreaking work on audio guided diffusion models to generate images, which is a very meaningful attempt.

3.1 Audio dependent multimodal models

There have been many multimodal works that have taken audio into account before, and some have replaced text with audio as the main object for matching with images (Ilharco et al., 2019; Chrupala, 2022). In addition to AudioCLIP (Guzhov et al., 2021) and other similar but actually different work, the most similar to us is Wav2CLIP (Wu et al., 2022). For CLIP, the BriVL we use has the following differences and advantages: Firstly, BriVL has more weak semantic relevance, so our model is more imaginative. For example, here are two groups of graphs in Figure 2 generated by using CLIP and BriVL respectively using GAN for comparison and understanding in the field of text-guided generation. Secondly, for our network

⁷<https://github.com/descriptinc/lyrebird-wav2clip>

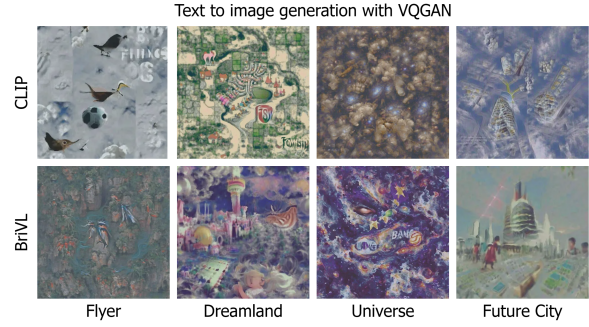


Fig. 2: Examples of CLIP (top) and BriVL (bottom) to image generation from text, BriVL’s labels in x-axis are translated.

architecture, because there is not necessarily a fine-grained area match between the image and audio, we lost the time-consuming target detector and adopted a simple and more efficient dual tower architecture, so we can encode the image and audio input through two independent encoders. Thirdly, BriVL designed a cross modal comparative learning algorithm based on the single modal comparative learning method MoCo (He et al., 2020), which has different advantages than CLIP.

3.2 Audio driven image generation

For many years, people have been trying to give AI people multimodal perception and thinking, and one of the main ideas is to simulate people’s impressions of different external inputs, namely image generation. The pursuit of applications and methods for generating different images is the direction of researchers’ efforts. With the emergence of different generation models, such as Goodfellow introduced GAN in 2014, there has been a lot of excellent work in the field of GAN-based image generation (Karras et al., 2017; Cudeiro et al., 2019; Yi et al., 2020; Zhang et al., 2021a; Song et al., 2022; Zhang et al., 2021b,c; Wu et al., 2021; Lahiri et al., 2021; Richard et al., 2021; Thies et al., 2020; Wen et al., 2020; Song et al., 2021; Chen et al., 2020b). Then, from single mode to multi-mode, from text guidance about 15 years later to audio guidance (Qiu and Kataoka, 2018) 20 years later (of course, there are more and earlier attempts and exceptions), several impressive works appeared (Xu et al., 2018; Zhu et al., 2021; Hessel et al., 2021; Saharia et al., 2022b,a). At a time when diffusion models have achieved success in many fields, exploring based on this work is meaningful.

Dataset	Task	Clip (Split)	Class Metric
ESC-50 (Piczak, 2015)	MC/ZS	2k (5 folds)	50 ACC
UrbanSound8K (Salamon et al., 2014)	MC/ZS	8k (10 folds)	10 ACC
VGGSound (Chen et al., 2020a)	MC/ZS	185k	309 mAP
DESED (Turpault et al., 2019)	AR	2.5k (valid)	10 F1
VGGSound (Chen et al., 2020a)	CMR	15k (test)	309 MRR
Clotho (Drossos et al., 2020)	AC	5k (evaluation)	COCO ⁶

Table 1: Downstream tasks, including 1. classification: multi-class (MC), zero-shot (ZS), 2. retrieval: audio (AR) and cross-modal retrieval (CMR), and 3. audio captioning (AC) task, with various of clips, classes, and common metrics.

4 Task 1: WavBriVL Performance Test

We begin by discussing the training, development, and evaluation process of the WavBriVL model. We use publicly available datasets of varying sizes and tasks, including classification, retrieval, and audio captioning tasks. We compare WavBriVL with some widely used as strong benchmarks in this field, and evaluate its performance in these tasks. Additionally, we investigate the effect of sound volume on the generated images. We hypothesize that the volume of sounds can influence the generated images. Hence, we explore the influence of sound volume on image features extracted from the sound using the sound correlation model. We also perform quantitative image analysis to evaluate the performance of WavBriVL compared to previous work, such as S2I and Pedersoli et al. We test model with five categories from VEGAS (Zhou et al., 2018) and compare its performance with other methods in terms of generating visually plausible images.

4.1 Training, development, and evaluation

We selected publicly available audio classification data of different sizes, which are generally used for evaluation (Cramer et al., 2019), and also included some audio tasks/data, as shown in table 1, including classification, retrieval and audio captioning. ESC-50 (Piczak, 2015) is a simple data set with only 2 thousand samples, while UrbanSound8K (Salamon et al., 2014) is a large environmental data set with 10 categories. VGGSound (Chen et al., 2020a) is a huge set of audio and video materials as we said before. DESED is used again as an audio extraction (AR) job because DESED can perform sound extraction at the fragment level. Finally, Clotho (Drossos et al., 2020) is a unique set of audio subtitles.

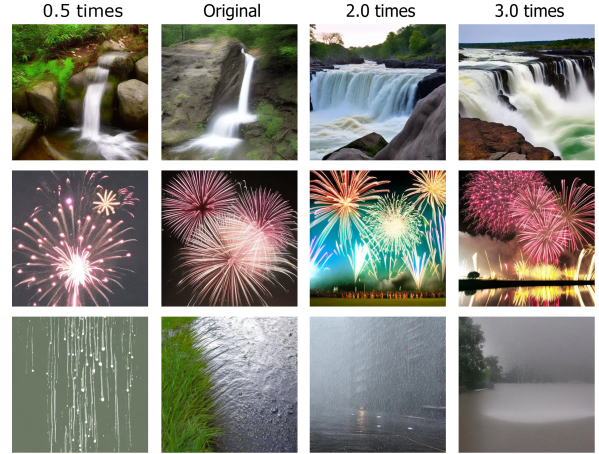


Fig. 3: Generated images by inputting different volumes of sounds. The numbers in the table is the relative loudness to the original sound.

For multi-class (MC) classification problems, an MLP-based classifier is employed, with a corresponding number of classes as output. In DESED, we use the way of simulating WavBriVL and sed_eval⁸ to realize audio retrieval (AR). At the same time, we also explore the performance of ours when dealing with multimodal tasks, and how to transfer zero samples to other modalities.

4.2 Sound volume

To establish the reliability of our method’s capability to learn the connection between sound and vision, we analyzed the influence of sound volume on generated images. To achieve this, we adjusted the sound volume levels during testing and extracted features for the corresponding sound files. These modified sound features were then input into our pre-trained generator, which was trained on a standard volume scale. The final three sets of images can prove our hypothesis that the magnitude of

⁸https://github.com/TUT-ARG/sed_eval

⁸<https://github.com/tylin/coco-caption>

Model	Classification				Retrieval	
	ESC-50	UrbanSound8K	VGGSound	DESED (AR)	VGGSound (CMR)	
	ACC	ACC	mAP	F1	A→I (MRR)	I→A (MRR)
Supervise	0.5200	0.6179	0.4331			
OpenL3	0.733	0.7588	0.3487	0.1170	0.0169	0.0162
Wav2CLIP	0.8595	0.8101	0.4663	0.3955	0.0566	0.0678
WavBriVL	0.9117	0.8832	0.4741	0.3720	0.0611	0.0608
SOTA	0.959	0.8949	0.544			
WavBriVL (ZS)	0.412	0.4024	0.1001			

Table 2: In the subsequent classification and acquisition work, there will be supervised training, other audio representation modes, OpenL3, and the latest SOTA (Guzhov et al., 2021; Kazakos et al., 2021). ZS is based on WavBriVL as a zero sample size model, some of which are derived from the original literature.

Method		VEGAS (5 classes)		
		R@1	FID (↓)	IS (↑)
(A)	Pedersoli et al.	23.10	118.68	1.19
(B)	S2I	39.19	114.84	1.45
(C)	S2V	77.58	34.68	4.01
(D)	Ours	81.31	31.48	5.42

Table 3: **Comparison to the baseline: Pedersoli et al. (2022) and existing sound-to-image/video method: S2I and S2V (Fanzeres and Nadeu, 2021; Sung-Bin et al., 2023).** Our method outperforms the others both qualitatively and quantitatively in the VEGAS dataset.

different volume levels is usually positively correlated with the effects and meanings displayed in the images.

4.3 Quantitative image analysis

We conducted a comparative analysis of our proposed model against publicly available prior works S2I⁹ (Fanzeres and Nadeu, 2021; Sung-Bin et al., 2023) and Pedersoli et al. (2022). It should be noted that while the latter is not primarily designed for sound-to-image conversion, it employs a VQVAE-based model to generate sound-to-depth or segmentation. We trained our model and Pedersoli et al. using the same training setup as S2I, including five categories in VEGAS, to ensure a fair comparison. As shown in Table 3, our proposed model outperforms all other models while generating visually compelling and recognizable images. We assert that this superior performance can be attributed to the combination of visually enriched audio embeddings and a powerful image generator.

⁹<https://github.com/leofanzeres/s2i>

4.4 Comparisons with previous work

First, we monitor the benchmark by training from scratch on each downlink (with random initialization of the encoder weights). Next, we compare WavBriVL with other publicly available OpenL3 (Cramer et al., 2019) pre-trained on different pre-text tasks in OpenL3. OpenL3 multimodal self-monitoring training with AudioSet. It serves as a strong benchmark for different audio tasks, such as audio classification and retrieval. We extract features from OpenL3 (512 dim) and WavBriVL (512 dim) and apply the same training scheme to all downstream classification and retrieval tasks. In the chart, we can see that in the retrieval of classification, we are slightly better than our previous work, with an average increase of about 0.04, and only some deficiencies in AR. But it’s only about 0.02. We approach or slightly outperform our previous work in retrieval tasks.

In summary, our model has good effects in both data sets of audio retrieval classification, for the source of our strengths: In the Classification tasks, on the four datasets, three of us achieved good results close to or exceeding SOTA. one of reason may be related to our data, and the other may be the effect of BriVL. As for the lack of excellent performance in AR tasks, it may be due to the excessive divergence of the BriVL dataset. If we retrain the basic model on a large scale, we may achieve better results. In the Retrieval tasks, such mrr tasks from A to I, from I to A we have also achieved excellent results, which mainly comes from the excellent training effect of the previous two towers model and the pre-training model, the structure of the brief is useful for general with tasks.

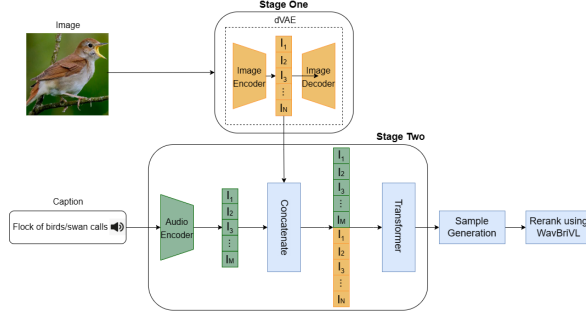


Fig. 4: Working principle of DALL·E+WavBriVL

5 Task2: Speech Generation Picture Based on Diffusion Model

To visually analyze the shared embedded space of WavBriVL and enhance its detectability. We asked WavBriVL to generate images with the help of DALL·E. The purpose of this task is to qualitatively evaluate the effect of our model. The image of the previous model is generated by VQGAN¹⁰, because the old model has not completely tried the effect of the diffusion model. The effect of their use of the diffusion model is worth looking forward to, but they have not yet tried it.

5.1 Processing method

WavBriVL includes an Audio Encoder and an Image Encoder. Its pre training model can accurately judge whether the given audio and image match. Similarly, in order to generate audio guidance image, we only need to match the image generated by DALL·E with the input audio according to whether BriVL "approves" it. If there is no match, feedback will be given to DALL·E (Ramesh et al., 2021) to guide it to generate more accurate images. This process is an iterative process of training DALL·E. In this iteration process, the image quality generated by DALL·E must be getting better and better, and closer to the limit of WavBriVL.

DALL·E image generator was created by OpenAI and it can be able to generate images similar to surrealism directly through text description. VQGAN is the choice of Wav2CLIP, and not very convenient to compare similar work. But, comparisons are still relevant. The goal of DALL·E is to treat the text token and image token as a data sequence and carry out auto-regression through Transformer. Due to the large resolution of the image, if a single pixel is treated as a token, it will lead to a huge amount of computation, so DALL·E introduced a dVAE model to reduce the resolution of image.

¹⁰<https://github.com/nerdyrodent/VQGAN-CLIP>

Options	Positive	Negative	Neither
Wav2CLIP	72 - 78%	9 - 17%	5 - 13%
WavBriVL	75 - 83%	12 - 18%	4 - 7%

Table 4: Human scores on correlation between sounds and images, Wav2CLIP works for comparison

1. In the first stage, first train a dVAE to compress each 256x256 RGB image into a 32x32 image token, and each position has 8192 possible values (that is, the encoder output of the dVAE is the logits with the dimension of 32x32x8192, and then combine the features of the codebook through the logits index. The embedding of the codebook is learnable).
2. In the second stage, the text is encoded with Text Encoder to obtain a maximum of x text tokens. If the number of tokens does not meet the maximum value, the maximum value is padded. Then x text tokens and 1024 image tokens are spliced to obtain 1280 data in length. Finally, the spliced data is input into Transformer for autoregressive training.
3. In the reasoning stage, given a candidate image and an audio, the fused token can be obtained through the transformer, and then the image can be generated by the dVAE decoder. Finally, as shown in Figure 4, the matching score of the audio and the generated image can be calculated through the pre-trained WavBriVL, ultimately achieving the effect of guiding the generation of the most matched image. As in general performance testing, DALL·E and WavBriVL are frozen during the generation process.

5.2 Correlation between sounds and images

This section aims to investigate whether the proposed method generates graphs that are also relevant to humans. In Figure 5, we demonstrate that our method can generate more eye-catching images; However, simply proving authenticity is not enough to prove the deep connection between sound and image. To demonstrate the connection between the two, we conducted a test similar to previous work (Ilharco et al., 2019; Wan et al., 2019). Participants were presented with two images, each with different sound categories as input and the image closest to the given sound. We conducted three tests and obtained a series of option values. By collecting participants' options, we aim to evaluate the effectiveness of the model in generating images related to different sound categories.

The experimental results are shown in Table 4,

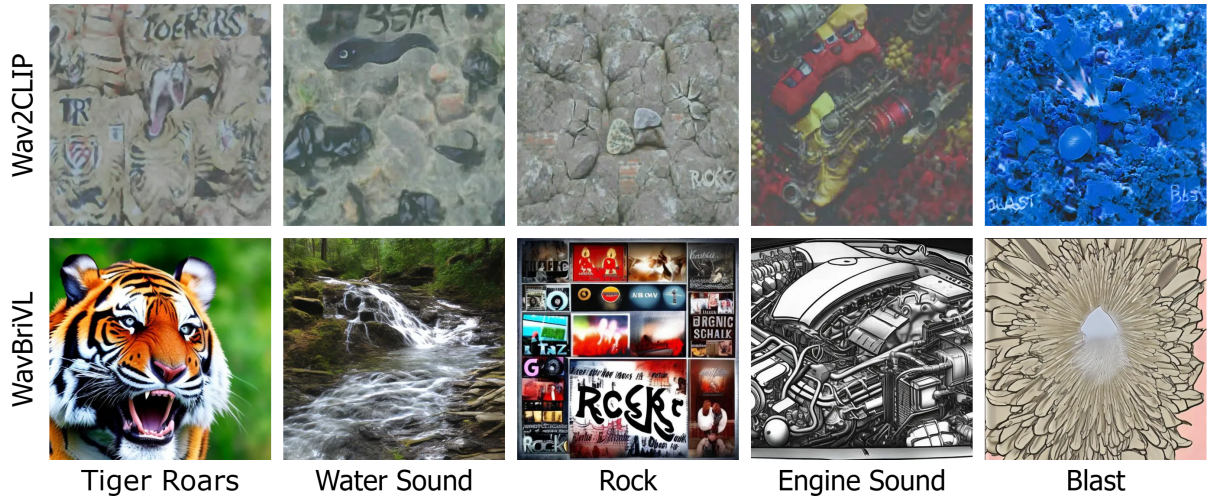


Fig. 5: Images generated from five-piece audio in AudioSet (Gemmeke et al., 2017). Top: Wav2CLIP, Bottom: WavBriVL - corresponding audio input labels in x-axis.

which collected participants' reactions and classified them as positive, negative, or neutral. A positive option indicates that participants have chosen images generated from input sound, while a negative option indicates their preference for images generated from different categories of sound. Participants who believe that neither of these images represents the sound they hear are considered neutral. Our research results indicate that the majority of participants believe that the generated images are related to the input sound, thus verifying our method's ability to generate images related to a given sound.

5.3 Comparison with previous work

In previous work, Wav2CLIP also tried to generate text/audio maps. Here are two sets of pictures for comparison with our work. Figure 2 shows the text output image of CLIP and BriVL. Figure 5 shows another group of pictures generated by Wav2CLIP and WavBriVL using audio.

However, in general, they all generated appropriate images, and they have their own characteristics: for example, in their understanding of "Tiger Roads", WavBriVL is more realistic, and WavCLIP is more abstract. When they faced the input of "Water Sound", our work generated a small stream, and WavCLIP generated symbolic images similar to fish fossils. Even considering the characteristics of the GAN model, this result can further prove the superiority of our work, which also indicates that our exploration and attempt to generate images using a universal audio guided diffusion model is meaningful; For the generation of audio, they exhibit two characteristics of convergence and divergence

between the two models, as we can see, convergence still corresponds to the image. Divergence is reflected in Figure 5 generated by audio, which is more imaginative than Figure 2 generated by text. This is because our BriVL weak semantic text image dataset has strong imagination, and another reason is that audio itself has strong divergence ability, which will enhance the associative ability of audio driven models.

6 Summary & Conclusion

This paper introduces a WavBriVL¹¹ for audio representation. The results show that WavBriVL is able to output general, robust sound representations, and that WavBriVL can be easily transferred to multimodal jobs, such as audio classification, audio retrieval, audio captioning and audio image generation. In future research, we will explore some interpretable machine learning approaches that uses the ability to generate (sound-image) across modalities. Based on learnings from embedded systems, additional speech classification and retrieval efforts are evaluated and compared to more advanced multimodal large models. On this basis, we will try to share the embedding space in multiple modes, so as to achieve the cross mode of image-generated text and image-generated sound. In the future, we will also consider exploring and using Microsoft's latest text-to-speech fusion model, SpeechLM (Zhang et al., 2022), the next release of the Diffusion model (Ho et al., 2020), the Consistency Models (Song et al., 2023) and the NeRF (Mildenhall et al., 2020) as the next version of the work.

¹¹AnonymousGitHub

Limitations

The limitation of WavBriVL is that the BriVL is trained based on the Chinese text image data set WSCD (the text-image is corresponding (Fei et al., 2022)), while our later training uses the English video data set VGG-Sound (the audio-image is corresponding). However, audio and text do not necessarily correspond strictly. It has no impact on the classification, retrieval, and generation tasks of audio image, but it is not recommended to use them when text image tasks are involved (of course, this is obviously a BriVL task, not our WavBriVL task). When future researchers explore multimodal mutual transformation, it is recommended to find a Chinese video dataset for retraining. The method in this paper is sufficient for generating correlation images. This is also mentioned in the Section 2.1 Dataset chapter.

Ethics Statement

All datasets we train actively exclude harmful, pornographic, and private content, and are only used for research purposes. The participants we recruited, except for some who volunteered, received satisfactory compensation for the rest. The academic tools and human assessment related tests used in this article comply with all regulations or relevant permits.

Biases & Content Acknowledgment Although our ability to generate images through audio is impressive, it should be noted that this model may be influenced by human factors to output content that enhances or exacerbates social biases.

References

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in Neural Information Processing Systems*, 33:12449–12460.

Hang Chen, Hengshun Zhou, Jun Du, Chin-Hui Lee, Jingdong Chen, Shinji Watanabe, Sabato Marco Siniscalchi, Odette Scharenborg, Di-Yuan Liu, Bao-Cai Yin, et al. 2022a. The first multimodal information based speech processing (misp) challenge: Data, tasks, baselines and results. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 9266–9270. IEEE.

Honglie Chen, Weidi Xie, Andrea Vedaldi, and Andrew

Zisserman. 2020a. Vggsound: A large-scale audio-visual dataset. In *ICASSP*, pages 721–725. IEEE.

Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi Xu, and Chenliang Xu. 2020b. Talking-head generation with rhythmic head motion. In *European Conference on Computer Vision*, pages 35–51. Springer.

Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. 2022b. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*.

Grzegorz Chrupała. 2022. Visually grounded models of spoken language: A survey of datasets, architectures and evaluation techniques. *Journal of Artificial Intelligence Research*, 73:673–707.

Jason Cramer, Ho-Hsiang Wu, Justin Salamon, and Juan Pablo Bello. 2019. Look, listen, and learn more: Design choices for deep audio embeddings. In *ICASSP*, pages 3852–3856. IEEE.

Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael J Black. 2019. Capture, learning, and synthesis of 3d speaking styles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10101–10111.

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

Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. 2020. *Clotho: An audio captioning dataset*. In *ICASSP*.

Leonardo A Fanzeres and Climent Nadeu. 2021. Sound-to-imagination: Unsupervised crossmodal translation using deep dense network architecture. *arXiv preprint arXiv:2106.01266*.

Nanyi Fei, Zhiwu Lu, Yizhao Gao, Guoxing Yang, Yuqi Huo, Jingyuan Wen, Haoyu Lu, Ruihua Song, Xin Gao, Tao Xiang, et al. 2022. Towards artificial general intelligence via a multimodal foundation model. *Nature Communications*, 13(1):1–13.

Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP 2017*, New Orleans, LA.

Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. 2021. Audioclip: Extending clip to image, text and audio. *arXiv preprint arXiv:2106.13043*.

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. *Momentum contrast for unsupervised visual representation learning*. In *2020*

698	<i>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 9726–9735.	751
699		752
700	Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7514–7528.	753
701		754
702		755
703		756
704		757
705		758
706	Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. <i>Advances in Neural Information Processing Systems</i> , 33:6840–6851.	759
707		760
708		761
709		762
710	Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. 2021. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 29:3451–3460.	763
711		764
712		765
713		766
714		767
715		768
716	Gabriel Ilharco, Yuan Zhang, and Jason Baldridge. 2019. Large-scale representation learning from visually grounded untranscribed speech. In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 55–65, Hong Kong, China. Association for Computational Linguistics.	769
717		770
718		771
719		772
720		773
721		774
722	Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In <i>International Conference on Machine Learning</i> , pages 4904–4916. PMLR.	775
723		776
724		777
725		778
726		779
727		780
728	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .	781
729		782
730		783
731		784
732		785
733	Tero Karras, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. 2017. Audio-driven facial animation by joint end-to-end learning of pose and emotion. <i>ACM Transactions on Graphics (TOG)</i> , 36(4):1–12.	786
734		787
735		788
736		789
737	Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. 2021. Slow-fast auditory streams for audio recognition. In <i>ICASSP</i> , pages 855–859.	790
738		791
739		792
740	Avisek Lahiri, Vivek Kwatra, Christian Frueh, John Lewis, and Chris Bregler. 2021. Lipsync3d: Data-efficient learning of personalized 3d talking faces from video using pose and lighting normalization. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 2755–2764.	793
741		794
742		795
743		796
744		797
745		798
746		799
747	Long Ma and Yanqing Zhang. 2015. Using word2vec to process big text data. In <i>2015 IEEE International Conference on Big Data (Big Data)</i> , pages 2895–2897. IEEE.	800
748		801
749		802
750		803
		804
		805
	Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In <i>ECCV</i> .	
	Xichen Pan, Peiyu Chen, Yichen Gong, Helong Zhou, Xinbing Wang, and Zhouhan Lin. 2022. Leveraging unimodal self-supervised learning for multimodal audio-visual speech recognition. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4491–4503, Dublin, Ireland. Association for Computational Linguistics.	
	Fabrizio Pedersoli, Dryden Wiebe, Amin Banitalebi, Yong Zhang, and Kwang Moo Yi. 2022. Estimating visual information from audio through manifold learning. <i>arXiv preprint arXiv:2208.02337</i> .	
	Karol J. Piczak. 2015. ESC: Dataset for Environmental Sound Classification. In <i>ACM Multimedia</i> , page 1015. ACM Press.	
	Yue Qiu and Hirokatsu Kataoka. 2018. Image generation associated with music data. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops</i> , pages 2510–2513.	
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, et al. 2021. Learning transferable visual models from natural language supervision. <i>ICML</i> .	
	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In <i>International Conference on Machine Learning</i> , pages 8821–8831. PMLR.	
	Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando De la Torre, and Yaser Sheikh. 2021. Meshtalk: 3d face animation from speech using cross-modality disentanglement. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 1173–1182.	
	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. 2022a. Photorealistic text-to-image diffusion models with deep language understanding. In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 36479–36494. Curran Associates, Inc.	
	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022b. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in Neural Information Processing Systems</i> , 35:36479–36494.	

- J. Salamon, C. Jacoby, and J. P. Bello. 2014. A dataset and taxonomy for urban sound research. In *ACM Multimedia*, pages 1041–1044, Orlando, FL, USA.
- Linsen Song, Wayne Wu, Chen Qian, Ran He, and Chen Change Loy. 2022. Everybody’s talkin’: Let me talk as you want. *IEEE Transactions on Information Forensics and Security*, 17:585–598.
- Luchuan Song, Bin Liu, Guojun Yin, Xiaoyi Dong, Yufei Zhang, and Jia-Xuan Bai. 2021. Tacr-net: Editing on deep video and voice portraits. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 478–486.
- Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. 2023. [Consistency models](#).
- Kim Sung-Bin, Arda Senocak, Hyunwoo Ha, Andrew Owens, and Tae-Hyun Oh. 2023. [Sound to visual scene generation by audio-to-visual latent alignment](#).
- Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR.
- Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner. 2020. Neural voice puppetry: Audio-driven facial reenactment. In *European conference on computer vision*, pages 716–731. Springer.
- Nicolas Turpault, Romain Serizel, Ankit Parag Shah, and Justin Salamon. 2019. [Sound event detection in domestic environments with weakly labeled data and soundscape synthesis](#). In *DCASE*, New York City, United States.
- Chia-Hung Wan, Shun-Po Chuang, and Hung-Yi Lee. 2019. [Towards audio to scene image synthesis using generative adversarial network](#). In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 496–500.
- Wenhui Wang, Hangbo Bao, Li Dong, Johan Bjorck, Zhiliang Peng, Qiang Liu, Kriti Aggarwal, Owais Khan Mohammed, Saksham Singhal, Subhojit Som, et al. 2022. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*.
- Xin Wen, Miao Wang, Christian Richardt, Ze-Yin Chen, and Shi-Min Hu. 2020. Photorealistic audio-driven video portraits. *IEEE Transactions on Visualization and Computer Graphics*, 26(12):3457–3466.
- Haozhe Wu, Jia Jia, Haoyu Wang, Yishun Dou, Chao Duan, and Qingshan Deng. 2021. Imitating arbitrary talking style for realistic audio-driven talking face synthesis. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1478–1486.
- Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. 2022. Wav2clip: Learning robust audio representations from clip. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4563–4567. IEEE.
- Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1316–1324.
- Ran Yi, Zipeng Ye, Juyong Zhang, Hujun Bao, and Yong-Jin Liu. 2020. Audio-driven talking face video generation with learning-based personalized head pose. *arXiv preprint arXiv:2002.10137*.
- Chenxu Zhang, Saifeng Ni, Zhipeng Fan, Hongbo Li, Ming Zeng, Madhukar Budagavi, and Xiaohu Guo. 2021a. 3d talking face with personalized pose dynamics. *IEEE Transactions on Visualization and Computer Graphics*.
- Chenxu Zhang, Yifan Zhao, Yifei Huang, Ming Zeng, Saifeng Ni, Madhukar Budagavi, and Xiaohu Guo. 2021b. Facial: Synthesizing dynamic talking face with implicit attribute learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3867–3876.
- Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan. 2021c. Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3661–3670.
- Ziqiang Zhang, Sanyuan Chen, Long Zhou, Yu Wu, Shuo Ren, Shujie Liu, Zhuoyuan Yao, Xun Gong, Lirong Dai, Jinyu Li, et al. 2022. Speechlm: Enhanced speech pre-training with unpaired textual data. *arXiv preprint arXiv:2209.15329*.
- Pengcheng Zhao, Yanxiang Chen, Lulu Zhao, Guang Wu, and Xi Zhou. 2022. Generating images from audio under semantic consistency. *Neurocomputing*, 490:93–103.
- Yipin Zhou, Zhaowen Wang, Chen Fang, Trung Bui, and Tamara L Berg. 2018. Visual to sound: Generating natural sound for videos in the wild.
- Hao Zhu, Man-Di Luo, Rui Wang, Ai-Hua Zheng, and Ran He. 2021. Deep audio-visual learning: A survey. *International Journal of Automation and Computing*, 18(3):351–376.