

VHASR: A Multimodal Speech Recognition System With Vision Hotwords

Anonymous ACL submission

Abstract

The image-based multimodal automatic speech recognition (ASR) model enhances speech recognition performance by incorporating audio-related image. However, some works suggest that introducing image information to model does not help improving ASR performance. In this paper, we propose a novel approach effectively utilizing audio-related image information and set up VHASR, a multimodal speech recognition system that uses vision as hotwords to strengthen the model’s speech recognition capability. Our system utilizes a dual-stream architecture, which firstly transcribes the text on the two streams separately, and then combines the outputs. We evaluate the proposed model on four datasets: Flickr8k, ADE20k, COCO, and OpenImages. The experimental results show that VHASR can effectively utilize key information in images to enhance the model’s speech recognition ability. Its performance not only surpasses unimodal ASR, but also achieves SOTA among existing image-based multimodal ASR.

1 Introduction

The unimodal ASR (Chan et al., 2015; Radford et al., 2023) only takes audio as input and produces corresponding transcription. In order to further reduce transcription errors, additional information related to the speech can be input, which can be in textual or visual modality. The ASR model that utilizes audio-related information from various modalities is referred to as multimodal ASR.

Common textual cues include hotwords, which are terms in certain professional fields or words that are easily confused with other homonyms. There have been many studies on how to freely customize hotwords and improve the recall of hotwords (Han et al., 2021; Shi et al., 2024). It is also possible to use captions as textual information (Moriya and Jones, 2018; Han et al., 2023).

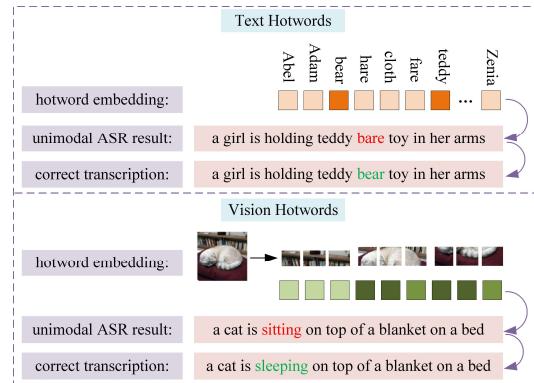


Figure 1: Comparison between text hotwords and the vision hotwords proposed in this paper. Text hotwords are a set of custom keywords that are prone to errors, while image hotwords refer to patches of an image. The hotword with a darker rectangle indicates that it is more relevant to transcription.

Visual cues can be in the form of video or image. Audio-Visual Speech Recognition (AVSR) enhances the accuracy of speech recognition by capturing lip movement information of characters in video (Ivanko et al., 2023). Image-based multimodal ASR extracts visual feature from image associated with speech to correct transcription errors. We abbreviate image-based multimodal ASR as IBSR. Because the lip movement information of video’s role is closely linked to his speech, it influences nearly every word in the transcribed text. In contrast, IBSR only impacts a subset of the words as the image is only associated with specific audio clips (Oneață and Cucu, 2022). IBSR currently lacks a universal and effective method for utilizing image information, leading to various experimental results in different studies. Some works (Sun et al., 2016; Srinivasan et al., 2020a,c) have a positive effect by incorporating image information, while others (Srinivasan et al., 2020b; Oneață and Cucu, 2022; Han et al., 2023), have the opposite effect.

In this paper, we propose a novel approach ef-

064 fectively utilizing audio-related image information
 065 and set up VHASR, a multimodal speech recogni-
 066 tion system that utilizes vision hotwords to enhance
 067 the model’s speech recognition capability. It cal-
 068 culates the similarity between different modalities
 069 to improve the effectiveness of cross-modal fusion.
 070 Drawing inspiration from text hotwords, we utilize
 071 Vision Transformer (ViT) to partition images into
 072 multiple visual tokens and consider each visual to-
 073 ken as an vision hotword. Our system adopts a
 074 dual-stream architecture. One stream is the ASR
 075 stream, which receives audio information and pro-
 076 duces transcribed text. The other stream is the vi-
 077 sion hotwords (VH) stream, which receives vision
 078 hotwords and audio hidden features, and generates
 079 corresponding text. In the VH stream, we calculate
 080 the similarity between audio and vision hotwords to
 081 reduce the weight of vision hotwords with low sim-
 082 ilarity. This process helps to extract fine-grained
 083 image information. When inferring, VHASR first
 084 transcribes the text separately from the ASR stream
 085 and the VH stream, and then merges the outputs.
 086 We ensure the high accuracy of the merged output
 087 by comparing the similarity of different modalities.
 088 Specifically, we first calculate the audio-image sim-
 089 ilarity to discard the VH stream if the similarity is
 090 low. Then, we calculate the image-text token simi-
 091 larity to compare the ASR stream and VH stream
 092 outputs by tokens. Finally, tokens with higher sim-
 093 ilarity are selected for the merged output.

094 We evaluate the proposed model on four datasets:
 095 Flickr8k, ADE20k, COCO, and OpenImages. The
 096 experimental results show that VHASR can effec-
 097 tively utilize critical information in images to im-
 098 prove the model’s ASR performance. Its perfor-
 099 mance is not only better than ordinary unimodal
 100 ASR models but also surpasses existing IBSR mod-
 101 els. The contributions of this paper are as follows:

- 102 (1) We demonstrate that through our idea of vi-
 103 sion hotwords, injecting audio-related image
 104 into the ASR model can help the model cor-
 105 rect transcription errors.
- 106 (2) We propose VHASR, by utilizing a dual-
 107 stream architecture and calculating the cross-
 108 modal similarity, it promote effective utiliza-
 109 tion of visual information in vision hotwords.
- 110 (3) The proposed model achieves SOTA on four
 111 datasets: Flickr8k, ADE20k, COCO, and
 112 OpenImages.

2 Related Work 113

Image-based multimodal ASR. Sun et al. (2016) 114
 115 introduces a multimodal speech recognition sce-
 116 nario based on RNN. This approach utilizes im-
 117 ages to assist the language model in decoding the
 118 most probable words and rescores the top hypothe-
 119 ses. Caglayan et al. (2019) proposes a novel multi-
 120 modal grounding method implemented by LSTM
 121 for sequence-to-sequence ASR. To utilize the vi-
 122 sual modality, they first project the visual vector
 123 into the speech feature space, and then use the vi-
 124 sual vector as the initial hidden and cell state for
 125 all LSTM layers. Srinivasan et al. (2020b) presents
 126 a model for multimodal ASR that integrates visual
 127 feature from object proposals. It calculates each
 128 modality’s attention distributions separately and
 129 combines attentions using a hierarchical attention
 130 mechanism in the decoder. Oneață and Cucu (2022)
 131 combines speech and visual embeddings using two
 132 fusion approaches. One approach fuses along the
 133 embedding dimension, and another fuses along the
 134 sequence dimension. They find that the first method
 135 performs better. Han et al. (2023) proposes a novel
 136 multimodal ASR model called ViLaS, which is
 137 based on the continuous integrate-and-fire (CIF)
 138 mechanism. It can integrate image and caption in-
 139 formation simultaneously or separately to facilitate
 140 speech recognition. Chang et al. (2023) proposes
 141 a multimodal ASR system for embodied agents.
 142 Their model is based on Transformer, where the vi-
 143 sual feature vector is concatenated to the decoder’s
 144 input word embedding at every timestep of genera-
 145 tion.

Function of image information. Srinivasan et al. 146
 147 (2020a) conducts the experiment called audio cor-
 148 ruption, in which they mask the words related to
 149 nouns and places with silence and white noise,
 150 respectively. The study demonstrates that visual
 151 representations help in recovering words that are
 152 masked in the input acoustic signal. Srinivasan
 153 et al. (2020c) thinks the previous work has only
 154 masked a fixed set of words in the audio, which is
 155 an unrealistic setting. So, they propose a method
 156 called RandWordMask, where masking can occur
 157 for any word segment to improve the audio cor-
 158 ruption experiment. Kumar et al. (2023) proposes
 159 two effective ASR error correction methods, which
 160 use gated fusion and image captions as prompts,
 161 respectively. Both methods demonstrate that vi-
 162 sual information helps restoring incorrect words
 163 in transcription. In short, image information helps

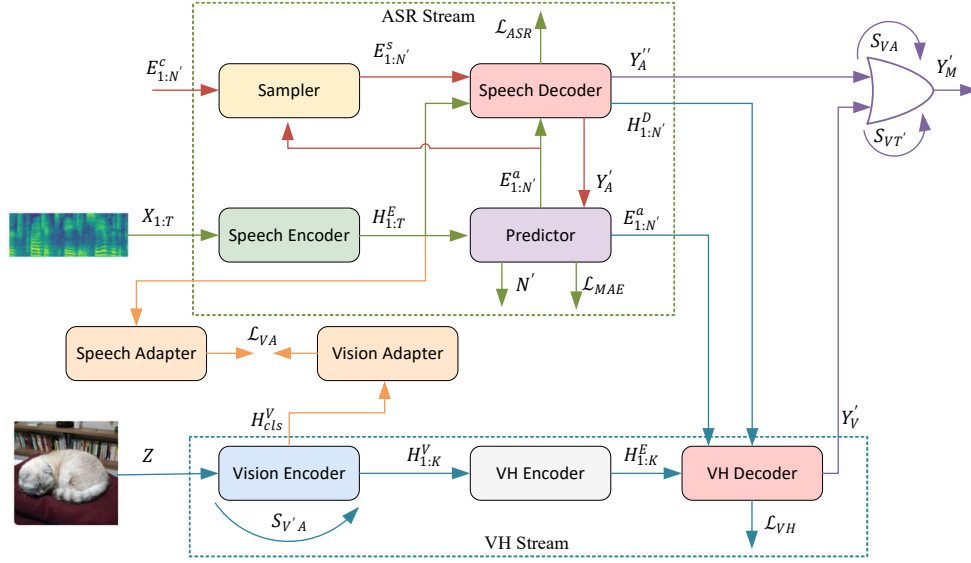


Figure 2: The structure of our proposed model, VHASR. The green dashed box contains the modules of the ASR stream, while the blue dashed box contains the modules of the VH stream. The data flow in the ASR part is indicated by green and red lines. It only passes through the red lines during ASR model’s second pass of training. The VH stream’s data flow is denoted by blue lines. The data flow for calculating audio-image similarity is represented by yellow lines. The purple lines illustrate the data flow when merging two streams.

to recover incorrect words in transcription that are caused by masked acoustic signals or ASR model’s error.

3 VHASR

3.1 ASR Stream

Follow Gao et al. (2022), we adopt this parallel Transformer (Vaswani et al., 2017) for non-autoregressive end-to-end speech recognition as the basic framework of our ASR stream. As shown in green dashed box of Figure 2, the adopted framework consists of four parts: speech encoder, predictor, sampler, and decoder. The framework adopts two-pass training and one-pass inference.

3.1.1 Acoustic Representation Learning

Let X be a speech sequence with T frames, $X = \{x_1, x_2, x_3, \dots, x_T\}$. Y is a sequence of tokens, and its length is N . Each token is in the vocabulary V , $Y = \{y_1, y_2, y_3, \dots, y_N \mid y_i \in V\}$.

The speech encoder adopts the SAN-M (Gao et al., 2020) structure, which is a special Transformer Layer that combines self-attention mechanism with deep feed-forward sequential memory networks (DFSMN). It converts the input $X_{1:T}$ to the hidden representation $H_{1:T}^E$.

$$H_{1:T}^E = \text{SpeechEncoder}(X_{1:T})$$

The predictor is a two-layer Deep Neural Networks (DNN) model that aligns speech and text based on CIF. It is used to predict the length of sentences N' and extract acoustic representation $E_{1:N'}^a$ from the speech encoder’s hidden representation $H_{1:T}^E$.

$$N', E_{1:N'}^a = \text{Predictor}(H_{1:T}^E)$$

The sampler does not contain learnable parameters and is only applied when training. It strengthens acoustic representation to semantic representation by incorporating text feature, aiming to better train the context modeling ability of the speech decoder. In the first pass of training, the sampler selects the text vectors $E_{1:N'}^c$ based on the number of different tokens between Y'_A and Y and their positions. Then, according to the position of incorrect tokens, it mixes $E_{1:N'}^c$ into $E_{1:N'}^a$ to obtain semantic feature $E_{1:N'}^s$. The semantic feature is used for the second pass of training. The sampler’s formula is as follows, λ is the sampler ratio, and $\lambda \in (0, 1)$.

$$E_{1:N'}^s = \text{Sampler}(E_{1:N'}^a, E_{1:N'}^c, \lceil \lambda \sum_{i=1}^{N'} (y'_i \neq y_i) \rceil)$$

3.1.2 Decoding Process

The speech decoder adopts the bidirectional SAN-M structure. In the first pass of training, the hid-

den representation $H_{1:T}^E$ obtained by the speech encoder and the acoustic representation $E_{1:N'}^a$ generated by the predictor are input to the speech decoder to obtain the initial decoding result Y'_A .

$$Y'_A = \text{SpeechDecoder}(H_{1:T}^E, N', E_{1:N'}^a)$$

In the second pass of training, the hidden representation $H_{1:T}^E$ and the semantic representation $E_{1:N'}^s$ obtained by the sampler are input to the speech decoder to obtain the second decoding result Y''_A

$$Y''_A = \text{SpeechDecoder}(H_{1:T}^E, N', E_{1:N'}^s)$$

During the first pass, no gradient backpropagation is performed, and Y'_A is only used to determine the sampling number of the sampler. Y''_A obtained in the second pass is used to calculate the ASR loss. In inference, the model directly takes Y'_A as output and does not calculate Y''_A .

3.2 Vision Hotwords Stream

3.2.1 Vision Representation Learning

In the VH stream, we need to extract visual feature from images by the vision encoder firstly. A naive idea is to extract the feature from the entire image. Because most of the information in the image is unrelated to the audio, especially the background of the image. The introduction of irrelevant information may cause the visual feature to become noise. Therefore, we should consider a strategy to extract fine-grained image information.

The vision encoder is essentially ViT (Dosovitskiy et al., 2020). ViT uses Transformer to extract visual feature. It follows the application of the Transformer in natural language processing by initially dividing the image into multiple patches, considering each patch as a token, embedding the positional information, and then feeding visual tokens into the Transformer. The feature outputted by ViT are the feature of each visual token. If the downstream task of ViT is classification, a trainable CLS token can be added in front of the visual token. The score on the CLS token can then be utilized for classification. It would be a good choice if we utilize each visual tokens' feature instead of entire image's feature. At the token granularity level, we can diminish the impact of tokens unrelated to audio and amplify the influence of tokens related to audio.

So, our strategy is to calculate the feature of each visual token and then adjust the weight of visual tokens. For the ASR model with text hotwords, it is often necessary to consider how to capture involved hotwords and exclude unrelated hotwords when there are many customized hotwords. This is similar to our consideration, so we call each visual token an vision hotword. Let Z be the input image. First, utilize the vision encoder to transform it into token-level visual feature $H_{0:K}^V$, where K represents the number of vision hotwords. The initial feature of $H_{0:K}^V$, corresponds to the feature of the CLS token, while others are vision hotwords' features.

$$H_{0:K}^V = \text{VisionEncoder}(Z)$$

$$H_{CLS}^V = H_{0:K}^V[0]; H_{1:K}^V = H_{0:K}^V[1:K]$$

We determine the correlation between each vision hotword and audio by calculating their cosine similarity. Specifically, the first step is to input $H_{1:K}^V$ into the vision adapter, which is composed of a linear layer, to obtain $H_{1:K}^{V'}$. Then, input the acoustic feature $H_{1:T}^E$ output by the speech encoder into the speech adapter, which is composed of a Transformer layer, to obtain $H_{1:T}^{E'}$.

$$H_{1:K}^{V'} = \text{VisionAdapter}(H_{1:K}^V)$$

$$H_{1:T}^{E'} = \text{SpeechAdapter}(H_{1:T}^E)$$

Then, calculate cosine similarity between vision hotwords and audio, denoted as $S_{V'A}$.

$$S_{V'A} = \cos(H_{1:K}^{V'}, H_{1:T}^{E'})$$

Finally, we adjust the weight of $H_{1:K}^V$ by $S_{V'A}$.

$$H_{1:K}^V = H_{1:K}^V \times S_{V'A}$$

In order to enhance the effectiveness of similarity-based weight adjustment, an additional loss needs to be introduced to train the adapters. We utilize the acoustic feature and the CLS token's feature of the image to calculate the image-audio contrastive loss \mathcal{L}_{VA} to optimize the adapters. The reason for using image-audio contrastive loss instead of vision hotwords-audio contrastive loss is that the former has a coarser granularity, making it easier to converge. Moreover, during inference, we need to use image-audio similarity for decoding optimization, which will be explained at length in Section 3.3. Figure 3 illustrates in detail our optimization of visual representation by calculating the

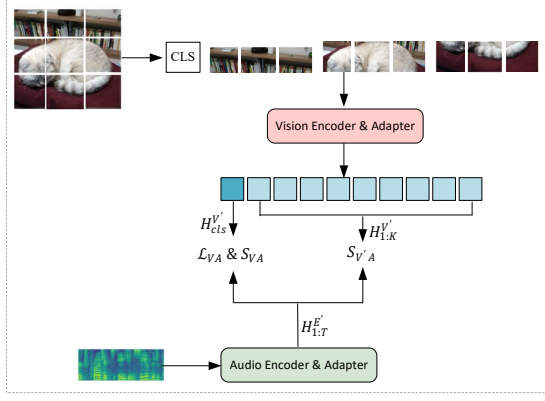


Figure 3: Using vision hotword-audio similitude and image-audio similitude to learn fine visual representation.

similitude between vision hotwords and audio, as well as the similitude between image and audio.

$$H_{CLS}^{V'} = \text{VisionAdapter}(H_{CLS}^V)$$

$$\mathcal{L}_{VA} = \text{ContrastiveLoss}(H_{CLS}^{V'}, H_{1:T}^E)$$

3.2.2 Decoding Process

The blue line in Figure 2 illustrates the data flow of the VH module. After extracting the fine visual representation of $H_{1:K}^V$, we further refine it using an LSTM-based VH encoder to obtain $H_{1:K}^E$.

$$H_{1:K}^E = \text{VHEncoder}(H_{1:K}^V)$$

The next step is to use a text decoder to obtain the probability distribution of each token. Obviously, if we only use $H_{1:K}^E$ which just contains image information as input, it will result in a significant deviation in the probability distribution of tokens, and the VH stream's outcome will be completely inconsistent with the correct transcription. So, we need to incorporate certain hidden features of the ASR stream to modify the output of the VH stream. Drawing lessons from the idea of Shi et al. (2024), we integrate the acoustic feature vector $E_{1:N'}^a$ outputted by the predictor and the hidden feature $H_{1:N'}^D$ outputted by the speech decoder with $H_{1:K}^E$ separately to derive $E_{1:N'}^{a'}$ and $H_{1:N'}^{D'}$, which have been influenced by image information. The VH decoder adopts the same bidirectional SAN-M architecture as the speech decoder.

$$E_{1:N'}^{a'} = \text{VHEncoder}(E_{1:N'}^a, H_{1:K}^E)$$

$$H_{1:N'}^{D'} = \text{VHEncoder}(H_{1:N'}^D, H_{1:K}^E)$$

The final input to the VH output layer is the average of $E_{1:N'}^{a'}$ and $H_{1:N'}^{D'}$.

$$Y_V' = \arg \max_{y_i \in V} (W_{1:V}^V \frac{(E_{1:N'}^{a'} + H_{1:N'}^{D'})}{2} + b_{1:V}^V)$$

3.3 Dual-stream Merging

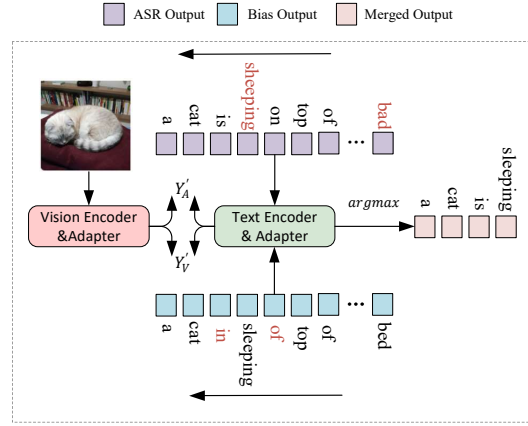


Figure 4: The specific process of decoding optimization.

In this section, we will discuss how to merge the outputs of the ASR stream and the VH stream. A straightforward approach is to add the probability distributions of tokens from two modules by assigning a specific weight, denoted as M_1 . The formula for M_1 is as follows, where p_A , p_V , and p_M are the tokens' probability distributions of the ASR stream, VH stream, and merged result. α is the proportion of p_A , and $\alpha \in (0, 1)$.

$$p_M = \alpha \text{Softmax}(p_A) + (1 - \alpha) \text{Softmax}(p_V)$$

$$Y_{M_1}' = \arg \max_{y_i \in V} (p_M)$$

The M_1 has low flexibility, making it difficult to achieve good results in practice. Figure 4 illustrates a merging method based on image-token similarity, referred to as M_2 . The vision encoder and adapter are used to calculate the visual feature of the image, $H_{CLS}^{V'}$, and the text encoder and adapter are used to calculate the feature of each token, $H_{1:N'}^{T'}$. The formula for $H_{CLS}^{V'}$ has been provided in Section 3.2.1, and the formula for $H_{1:N'}^{T'}$ is as follows. The text encoder consists of Transformer layers, the text adapter consists of a linear layer, and *Embedding* is an additional embedding layer.

$$H_{1:N'}^T = \text{TextEncoder}(\text{Embedding}(Y_V'))$$

$$H_{1:N'}^{T'} = \text{TextAdapter}(H_{1:N'}^T)$$

Based on $H_{CLS}^{V'}$ and $H_{1:N'}^{T'}$, the cosine similarity of the image and tokens, $S_{VT'}^V$, can be calculated.

$$S_{VT'}^V = \cos(H_{CLS}^{V'}, H_{1:N'}^{T'})$$

When calculating Y_{M_2}' , we first calculate the text feature of the ASR stream output Y_A' and the VH stream output Y_V' , respectively, namely $H_{1:N'}^{T'_A}$ and $H_{1:N'}^{T'_V}$. Then calculate their cosine similarities with $H_{CLS}^{V'}$ separately, namely $S_{VT'}^{A}$ and $S_{VT'}^V$. Finally, a token by token comparison of the dual-stream is conducted according to $S_{VT'}^A$ and $S_{VT'}^V$. Specifically, the value of these two similarities at any position represents the similarity score between the token at that position and the image. At the same position, Y_A' and Y_V' may obtain different tokens. We determine which token to choose as the final result by judging the value of $S_{VT'}^A$ and $S_{VT'}^V$ at that position. If $S_{VT'}^A > S_{VT'}^V$, we take the token on Y_A' , and vice versa. After completing N' comparisons, Y_{M_2}' can be obtained.

In Section 3.2.1, to achieve an fine-grained visual representation, we additionally introduce speech and vision adapters in VHASR to compute the similarity between vision hotwords and audio. Then, to train the adapter, we calculate contrastive loss between the image and audio. In the inference stage, we can further utilize the trained adapter to optimize M_2 by calculating image-audio similarity. Specifically, we calculate the image-audio similarity S_{VA} for a batch of data. If the audio of a piece of data does not match its own image, it is considered that the correlation between this image and audio is low. Therefore, for this data, the output of the VH stream is discarded, and the output of the ASR stream is directly used as the final output. We introduce a novel merging method called M_3 . It involves initially filtering data with low image and audio correlation using S_{VA} , followed by dual-stream merging as outlined in M_2 . We will conduct a detailed comparative experiment on these three merging methods in Section 4.

4 Experiment

4.1 Configuration

Table 1 shows all the datasets used in this paper, with Flickr8k, ADE20k, COCO, and OpenImages

used for training and testing, and SpokenCOCO used for pre-training. Flickr8k is from Harwath and Glass (2015) and SpokenCOCO is from Hsu et al. (2021). ADE20k, COCO and OpenImages are from Local Narratives proposed by (Harwath et al., 2016). In order to shorten the experimental period, we filter data with audio exceeding 40s in ADE20k, and with more than 40 tokens or an audio duration of more than 20 seconds in COCO and OpenImages. We use word error rate (WER) as an evaluation metric to evaluate the speech recognition performance of ASR stream, VH stream, M_1 , M_2 , and M_3 .

our baseline is 220M English Paraformer. In Flickr8k, we compare our model with Acoustic-LM-RNN proposed by Sun et al. (2016), model utilizing object features as visual information (abbreviated as Multimodal (object) in the paper) from Srinivasan et al. (2020a), Weighted-DF in Srinivasan et al. (2020c), MAG proposed by Srinivasan et al. (2020b), model fusing the two modalities along the sequence dimension (abbreviated as Multimodal (emb) in the paper) from Oneață and Cucu (2022) and ViLaS in Han et al. (2023).

The modules in CLIP-Base (Radford et al., 2021) is utilized to construct the vision encoder and vision adapter for the VH stream, as well as the vision encoder and text encoder for M_2 . The vision module of the VH stream freeze parameters during training, and the M_2 's modules do not require training. The 220M English Paraformer is chosen as the foundational framework for ASR stream, initialized with the same parameters as the baseline. λ of sampler is set to 0.75 and α of M_1 is set to 0.5. The experimental environment is constructed using Funasr (Gao et al., 2023) and ModelScope. We trained the model until convergence, with a maximum of 80 training epochs. We consistently utilize the Adam optimizer with a learning rate of $5e-5$, and the training is conducted on GeForce RTX 3090.

Dataset	Train	Validation	Test
Flickr8k	29,999	4,998	4,998
ADE20k	17,067	1,672	-
COCO	49,109	3,232	-
OpenImages	269,749	27,813	-
SpokenCOCO	592,187	25,035	-

Table 1: Datasets used in experiments.

4.2 Main Result

Table 2 presents the results of the proposed method and baseline on four datasets. For the ASR stream

Dataset	Baseline	VHASR					
	WER (\downarrow)	Pretrain	WER _{ASR} (\downarrow)	WER _{VH} (\downarrow)	WER _{M₁} (\downarrow)	WER _{M₂} (\downarrow)	WER _{M₃} (\downarrow)
Flickr8k	3.91	×	3.82	3.91	3.79	3.56	3.51
		✓	3.68	3.75	3.66	3.37	3.35
ADE20k	10.51	×	10.33	10.52	10.38	9.80	9.60
		✓	10.27	10.37	10.32	9.62	9.53
COCO	10.44	×	10.35	10.34	10.28	9.63	9.61
		✓	10.25	10.36	10.28	9.60	9.59
OpenImages	8.72	×	8.61	8.58	8.58	7.73	7.71
		✓	8.58	8.63	8.59	7.70	7.68

Table 2: Main results of proposed model in four datasets.

and VH stream, the WER of the ASR stream is lower. The VH stream can acquire the ability of transcribing by utilizing the hidden layer’s features of the ASR stream as VH decoder’s input. Among the three merge methods, M_3 has the best results, followed by M_2 , and finally M_1 . This is consistent with our expected results. M_1 has limited flexibility, and the fixed weight proportion is not applicable to all data. By calculating image-token similarity, comparing the results of the ASR stream and VH stream token by token, and resulting in a final output with the highest similarity, M_2 achieves WER that are better than both WER_{ASR} and WER_{VH}. Furthermore, by calculating audio-image similarity in addition and excluding the VH stream with low similarity, M_3 reduces the transcription error compared to M_2 . For the baseline and ASR stream, ASR stream performs better, indicating that joint training of the ASR stream, VH stream, and audio-image pairing improves the unimodal ASR’s performance. For the baseline and M_3 , M_3 outperforms the baseline on all four datasets, demonstrating the effectiveness of our method. In addition, pre-training with large-scale corpora can further strengthen the performance of the model. We use SpokenCOCO, which contains the largest amount of data, to pre-train the proposed model, resulting in improvements in all five metrics of the model across all four datasets.

4.3 Ordinary Multimodal Fusion vs Hotword Level Multimodal Fusion

The comparison results are shown in the Table 3. Without vision information, Vilas (Han et al., 2023) performs better than our VHASR since they have done sufficient pretraining. With vision information, VHASR’s ASR performance has been significantly enhanced and it achieves the lowest WER. Obviously, our experimental results indicate that the incorporation of visual information aids in rec-

tifying tokens for ASR transcription errors and decreasing WER. However, Srinivasan et al. (2020b), Oneață and Cucu (2022) and Han et al. (2023) argue that the speech in Flickr8k is sufficiently clear, making it challenging to enhance transcription performance by incorporating additional information from other modalities.

Model	Word Error Rate (\downarrow)	
	w/o vision	w vision
Acoustic-LM-RNN (Sun et al., 2016)	14.75	13.81 (\downarrow 0.94)
Multimodal (object) (Srinivasan et al., 2020a)	16.40	14.80 (\downarrow 1.60)
Weighted-DF (Srinivasan et al., 2020c)	13.70	13.40 (\downarrow 0.30)
MAG (Srinivasan et al., 2020b)	13.60	13.80 (\uparrow 0.20)
Multimodal (emb) (Oneață and Cucu, 2022)	3.80	4.30 (\uparrow 0.50)
ViLaS (Han et al., 2023)	3.40	3.40 (\downarrow 0)
VHASR	3.91	3.35 (\downarrow 0.56)

Table 3: Comparison results with benchmarks in F8k.

MAG (Srinivasan et al., 2020b) utilizes global visual feature, which may introduce a significant amount of information unrelated to audio and potentially impact the model’s ASR performance. They considered this issue and proposed MAOP, which utilizes multiple fine-grained image feature extracted from object proposals. But in terms of clean Flickr8k, MAOP’s performance is not as good as MAG’s. Oneață and Cucu (2022) takes a sequence of image feature vectors from the layer preceding the global average pooling layer in the vision encoder, for leveraging more fine-grained characteristics of the image. However, they did not consider that some image vectors in the sequence have low correlation with the audio. Introducing these vectors fully into the backbone will still impact the model’s recognition ability. Han et al. (2023) uses ViT as a vision encoder and utilizes the image tokens for visual representation, which aligns with our approach. However, they do not reduce the weight of visual tokens with low importance, as we do. This resulted in the intro-

Dataset	Mask Ratio	Baseline		VHASR			
		WER (\downarrow)	RR (\uparrow)	WER _{ASR} (\downarrow)	RR _{ASR} (\uparrow)	WER _{M₂} (\downarrow)	RR _{M₂} (\uparrow)
Flickr8k	30%	30.68	79.78	29.20	81.76	23.89	81.84
	50%	49.05	68.27	47.43	71.19	40.33	71.59
	70%	63.93	58.00	65.33	58.52	56.54	58.99
ADE20k	30%	24.79	92.02	24.40	92.51	19.96	92.60
	50%	34.16	89.18	32.95	89.86	26.91	90.06
	70%	42.30	86.33	40.70	87.45	33.39	87.46
COCO	30%	25.60	92.02	24.23	92.85	20.13	92.87
	50%	35.59	89.42	33.22	91.05	27.06	91.05
	70%	44.00	87.76	41.35	89.26	33.84	89.32

Table 4: Experimental results of audio corruption with AWGN.

duction of visual information not improving the recognition performance of the model. Compared to these works that use ordinary multimodal fusion approach, our proposed method, which injects visual modality information by vision hotwords, have made improvements in refining image representation and eliminating irrelevant image information. Therefore, our proposed model can enhance performance using visual feature even when the dataset is of high quality and the baseline is strong.

4.4 Ablation Result

Dataset	WER _{M₁} (\downarrow)		WER _{M₂} (\downarrow)	
	w/o refine	w refine	w/o refine	w refine
Flickr8k	3.93	3.79	3.68	3.56
ADE20k	10.67	10.38	10.17	9.80
COCO	10.46	10.28	9.64	9.63
OpenImages	8.73	8.58	7.81	7.73

Table 5: Experimental results of ablation studies.

To demonstrate that the refined image representation extracted by the method proposed in Section 3.2.1 is more effective than the full image representation, we conduct the ablation experiments. The experimental results are presented in Table 5. On four datasets, whether it is M_1 or M_2 , the model using refined image representation has better performance. This not only shows the effectiveness of the method described in Section 3.2.1 but also offers one of reasons why our model is stronger than other benchmarks.

4.5 Audio Corruption

To further demonstrate that introducing image information related to audio can reduce transcription errors in proposed model, we conduct an audio corruption experiment proposed by Srinivasan et al. (2020a). We first use the timestamp prediction model proposed by Shi et al. (2023) to align audio

and transcribed text. Then, we mask the words in the audio to a certain proportion by replacing the audio segments corresponding to the masked words with Additive White Gaussian Noise (AWGN). We use the recovery rate (RR) defined in Srinivasan et al. (2020a) to calculate the proportion of masked words recovered in the model transcription results. Unlike Srinivasan et al. (2020a), our approach only masks the test data, while the training data remains unchanged.

We conduct this experiment on Flickr8k, ADE20k, and COCO, and the experimental results are shown in Table 4. In terms of baseline and ASR stream, regardless of the mask ratio, the ASR stream has lower WER and higher RR on all three datasets. This suggests that the jointly trained ASR stream exhibits stronger noise resistance and audio content prediction abilities compared to unimodal ASR. In terms of ASR stream and M_2 , by incorporating image information, M_2 significantly reduces WER and enhances RR, as evidenced by the mask ratio across the three datasets. This indicates that image information can assist the model in capturing image-related words in audio, enabling the model to accurately transcribe these words even if their corresponding audio is masked. Furthermore, we can argue that on normal unmasked data, image information can assist the model in correcting words related to image but with transcription errors.

5 Conclusion

We propose VHASR, a multimodal speech recognition system that utilizes vision hotwords to strengthen the model’s speech recognition ability. To improve the effectiveness of cross-modal fusion, it calculates the similarity between different modalities. Through extensive experiments, we demonstrate that VHASR has powerful ASR performance.

528 Limitations

529 The Limitations of VHASR include: (1) currently,
530 VHASR can only introduce image information to
531 enhance the model’s speech recognition ability,
532 which does not have sufficient versatility. In the
533 future, we will enable VHASR to support input of
534 audio-related text information (such as hotwords, ti-
535 tles) and video information, enabling the model to
536 extract feature beneficial for speech recognition
537 from multiple modal information, and building
538 a more versatile multimodal speech recognition
539 model. (2) we have only demonstrated that vision
540 hotwords is a effective way to utilize image infor-
541 mation, and there may be other applicable methods.
542 We will design more in-depth experiments in the
543 following work to explore more feasible ideas.

544 References

545 Ozan Caglayan, Ramon Sanabria, Shruti Palaskar, Loic
546 Barraul, and Florian Metze. 2019. Multimodal
547 grounding for sequence-to-sequence speech recog-
548 nition. In *ICASSP 2019-2019 IEEE International
549 Conference on Acoustics, Speech and Signal Process-
550 ing (ICASSP)*, pages 8648–8652. IEEE.

551 William Chan, Navdeep Jaitly, Quoc V Le, and Oriol
552 Vinyals. 2015. Listen, attend and spell. *arXiv
553 preprint arXiv:1508.01211*.

554 Allen Chang, Xiaoyuan Zhu, Aarav Monga, Seoho Ahn,
555 Tejas Srinivasan, and Jesse Thomason. 2023. Multi-
556 modal speech recognition for language-guided em-
557 bodied agents. *arXiv preprint arXiv:2302.14030*.

558 Alexey Dosovitskiy, Lucas Beyer, Alexander
559 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
560 Thomas Unterthiner, Mostafa Dehghani, Matthias
561 Minderer, Georg Heigold, Sylvain Gelly, et al. 2020.
562 An image is worth 16x16 words: Transformers
563 for image recognition at scale. *arXiv preprint
564 arXiv:2010.11929*.

565 Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian
566 Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao
567 Du, Zhangyu Xiao, et al. 2023. Funasr: A funda-
568 mental end-to-end speech recognition toolkit. *arXiv
569 preprint arXiv:2305.11013*.

570 Zhifu Gao, Shiliang Zhang, Ming Lei, and Ian
571 McLoughlin. 2020. San-m: Memory equipped self-
572 attention for end-to-end speech recognition. *arXiv
573 preprint arXiv:2006.01713*.

574 Zhifu Gao, Shiliang Zhang, Ian McLoughlin, and Zhijie
575 Yan. 2022. Paraformer: Fast and accurate parallel
576 transformer for non-autoregressive end-to-end speech
577 recognition. *arXiv preprint arXiv:2206.08317*.

Minglun Han, Feilong Chen, Ziyi Ni, Linghui Meng,
Jing Shi, Shuang Xu, and Bo Xu. 2023. Vilas: In-
tegrating vision and language into automatic speech
recognition. *arXiv preprint arXiv:2305.19972*.

Minglun Han, Linhao Dong, Shiyu Zhou, and Bo Xu.
2021. Cif-based collaborative decoding for end-to-
end contextual speech recognition. In *ICASSP 2021-
2021 IEEE International Conference on Acoustics,
Speech and Signal Processing (ICASSP)*, pages 6528–
6532. IEEE.

David Harwath and James Glass. 2015. Deep multi-
modal semantic embeddings for speech and images.
In *2015 IEEE Workshop on Automatic Speech Recog-
nition and Understanding (ASRU)*, pages 237–244.
IEEE.

David Harwath, Antonio Torralba, and James Glass.
2016. Unsupervised learning of spoken language
with visual context. *Advances in Neural Information
Processing Systems*, 29.

Wei-Ning Hsu, David Harwath, Tyler Miller, Christo-
pher Song, and James Glass. 2021. Text-free image-
to-speech synthesis using learned segmental units.
In *Proceedings of the 59th Annual Meeting of the
Association for Computational Linguistics and the
11th International Joint Conference on Natural Lan-
guage Processing (Volume 1: Long Papers)*, pages
5284–5300.

Denis Ivanko, Dmitry Ryumin, and Alexey Karpov.
2023. A review of recent advances on deep learning
methods for audio-visual speech recognition. *Mathe-
matics*, 11(12):2665.

Vanya Bannihatti Kumar, Shanbo Cheng, Ningxin Peng,
and Yuchen Zhang. 2023. Visual information matters
for asr error correction. In *ICASSP 2023-2023 IEEE
International Conference on Acoustics, Speech and
Signal Processing (ICASSP)*, pages 1–5. IEEE.

Yasufumi Moriya and Gareth JF Jones. 2018. Lstm
language model adaptation with images and titles for
multimedia automatic speech recognition. In *2018
IEEE Spoken Language Technology Workshop (SLT)*,
pages 219–226. IEEE.

Dan Oneață and Horia Cucu. 2022. Improving mul-
timodal speech recognition by data augmentation
and speech representations. In *Proceedings of the
IEEE/CVF Conference on Computer Vision and Pat-
tern Recognition*, pages 4579–4588.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-
try, Amanda Askell, Pamela Mishkin, Jack Clark,
et al. 2021. Learning transferable visual models from
natural language supervision. In *International confer-
ence on machine learning*, pages 8748–8763. PMLR.

Alec Radford, Jong Wook Kim, Tao Xu, Greg Brock-
man, Christine McLeavey, and Ilya Sutskever. 2023.
Robust speech recognition via large-scale weak su-
pervision. In *International Conference on Machine
Learning*, pages 28492–28518. PMLR.

635 Xian Shi, Yanni Chen, Shiliang Zhang, and Zhijie Yan.
636 2023. Achieving timestamp prediction while recog-
637 nizing with non-autoregressive end-to-end asr model.
638 In *arXiv preprint arXiv:2301.12343*.

639 Xian Shi, Yexin Yang, Zerui Li, Yanni Chen, Zhifu
640 Gao, and Shiliang Zhang. 2024. Seaco-paraformer:
641 A non-autoregressive asr system with flexible and
642 effective hotword customization ability. In *ICASSP*
643 *2024-2024 IEEE International Conference on Acous-*
644 *tics, Speech and Signal Processing (ICASSP)*, pages
645 10346–10350. IEEE.

646 Tejas Srinivasan, Ramon Sanabria, and Florian Metze.
647 2020a. Looking enhances listening: Recovering
648 missing speech using images. In *ICASSP 2020-2020*
649 *IEEE International Conference on Acoustics, Speech*
650 *and Signal Processing (ICASSP)*, pages 6304–6308.
651 IEEE.

652 Tejas Srinivasan, Ramon Sanabria, Florian Metze, and
653 Desmond Elliott. 2020b. Fine-grained grounding for
654 multimodal speech recognition. In *Findings of the*
655 *Association for Computational Linguistics: EMNLP*
656 *2020*, pages 2667–2677.

657 Tejas Srinivasan, Ramon Sanabria, Florian Metze, and
658 Desmond Elliott. 2020c. Multimodal speech recog-
659 nition with unstructured audio masking. In *Proceed-*
660 *ings of the First International Workshop on Natural*
661 *Language Processing Beyond Text*, pages 11–18.

662 Felix Sun, David Harwath, and James Glass. 2016.
663 Look, listen, and decode: Multimodal speech recog-
664 nition with images. In *2016 IEEE Spoken Language*
665 *Technology Workshop (SLT)*, pages 573–578. IEEE.

666 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
667 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
668 Kaiser, and Illia Polosukhin. 2017. Attention is all
669 you need. *Advances in neural information processing*
670 *systems*, 30.

A Appendix

A.1 Case Study

In Section 4.5, we demonstrated that VHASR can use image information to correct words which is related to images and has transcription errors. In this section, we will use examples to explain how VHASR achieves this.

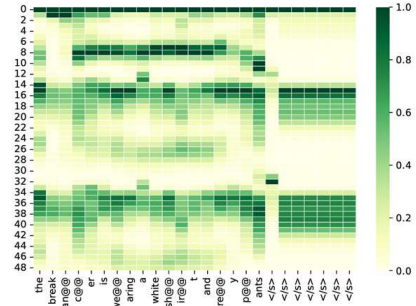
Figure 5 shows three examples from Flickr8k. "A" refers to the transcription of the ASR stream, "V" refers to the transcription of the VH stream, "M" refers to the transcription obtained by M_3 , and "T" refers to the real transcription. We extract the attention score matrix from the last layer of the VH decoder and create a heatmap. The horizontal axis of the heatmap represents the subtoken, while the vertical axis represents the number of vision hotwords. We identify the subtokens that are transcribed incorrectly by the ASR stream but corrected by the VH stream. Then, we extract the top 5 vision hotwords that have the highest attention scores with them. Chosen vision hotwords are marked on the original image.

In the first example, the ASR stream incorrectly transcribes "grey" as "gry", while the VH stream doesn't make this mistake. The combination of the two streams helps correct the error. specifically, the subtokens corresponding to "grey" focus on six vision hotwords, five of which are background, and one includes the grey pants of the dancer. Therefore, the vision encoder successfully extracts information about "grey" and helps the VH stream transcribe "grey" accurately. Furthermore, by merging the ASR stream and VH stream with M_3 , error in the ASR stream is rectified. In the second example, the ASR stream incorrectly transcribes "girls" as "girl", which was also corrected by the accurate VH stream. Among the vision hotwords corresponding to "girls", three are related to background, and two include the heads of the girls, so the VH stream successfully identified "girls". In the third example, the ASR stream incorrectly transcribes "river" as "room", but the VH stream correctly transcribes "river" by utilizing the information about "river" contained in the vision hotwords. By merging, the VH stream helps correct error in the ASR stream. These examples are not unique, and the same phenomenon occurs in many utterances. In Figure 6, we show another three examples from COCO for readers' reference.

Although the VH stream of VHASR has less speech recognition ability than the ASR stream, it

can extract features from key vision hotwords and capture keywords in transcription, thereby correctly identifying words that may be difficult for the ASR stream to recognize. After token-by-token merging based on visual-token similarity, the VH stream can correct some transcription errors in the ASR stream, leading to a more accurate transcription.

722
723
724
725
726
727
728

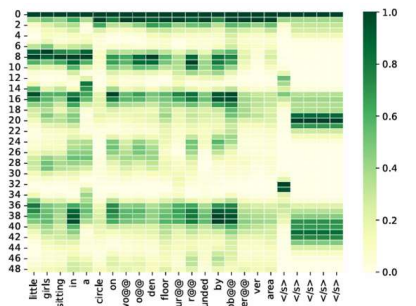
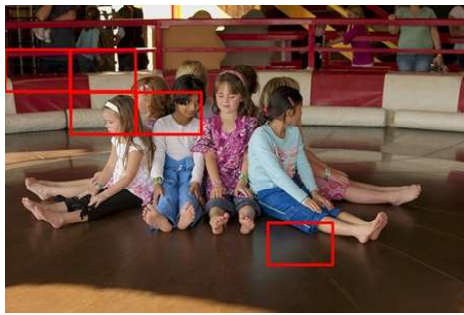


A: The break dancer is wearing a white shirt and **gry** pants.

M: The break dancer is wearing a white shirt and **grey** pants.

V: The break dancer is wearing a white shirt and **grey** pants.

T: The break dancer is wearing a white shirt and **grey** pants.

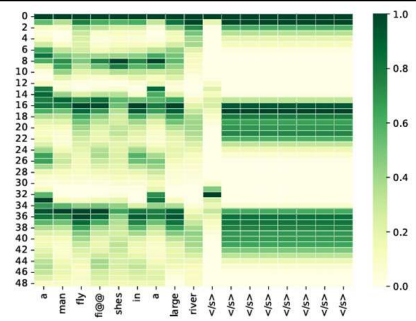
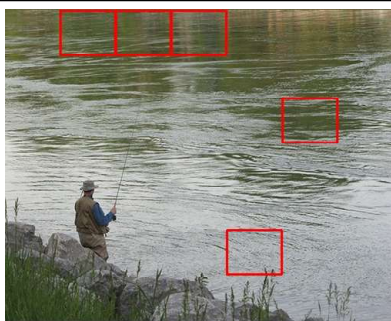


A: Little **girl** sitting in a circle on wooden floor surrounded by observer area.

M: Little **girls** sitting in a circle on wooden floor surrounded by observer area.

V: Little **girls** sitting in a circle on wooden floor surrounded by observer area.

T: Little girls sitting in a circle on wooden floor surrounded by observer area.



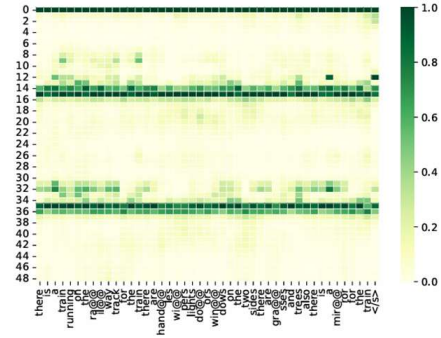
A: A man flv fishes in a large **room**.

M: A man fly fishes in a large **river**.

V: A man fly fishes in a large **river**.

T: A man fly fishes in a large **river**.

Figure 5: Three examples about how VH stream helps to rectify ASR stream's error.

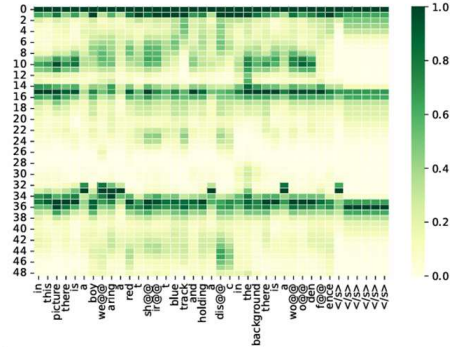
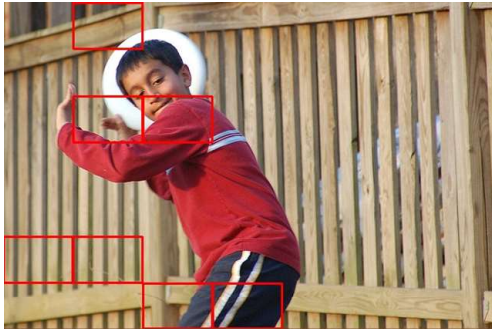


A: There is a train running on the railway track. For the train, there are handles, wipers, lights, doors, windows. On the two sides, there are grasses and trees, also there is a mirror for the **tree**.

M: There is a train running on the railway track. For the train, there are handles, wipers, lights, doors, windows. On the two sides, there are grasses and trees, also there is a mirror for the **train**.

V: There is a train running on the railway track. For the train, there are handles, wipers, lights, doors, windows. On the two sides, there are grasses and trees, also there is a mirror for the **train**.

T: There is a train running on the railway track. For the train, there are handles, wipers, lights, doors, windows. On the two sides, there are grasses and trees, also there is a mirror for the train.

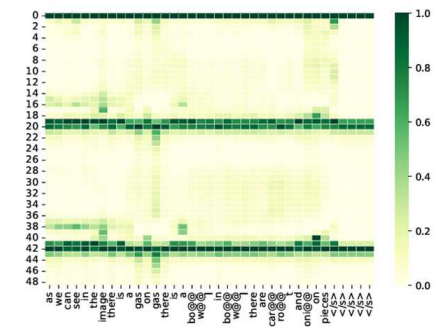


A: In this picture there is a boy wearing a red shirt, blue track and holding a **disk**. in the background there is a wooden fence.

M: In this picture there is a boy wearing a red shirt, blue track and holding a **disc**. in the background there is a wooden fence.

V: In this picture there is a boy wearing a red shirt, blue track and holding a **disc**. in the background there is a wooden fence.

T: In this picture there is a boy wearing a red shirt, blue track and holding a disc. in the background there is a wooden fence.



A: As we can see, in the image there is a **glass**. On gas, there is a bowl. In bowl, there are carrot and onion pieces.

M: As we can see, in the image there is a **gas**. On gas, there is a bowl. In bowl, there are carrot and onion pieces.

V: As we can see, in the image there is a **gas**. On gas, there is a bowl. In bowl, there are carrot and onion pieces.

T: As we can see, in the image there is a gas. On gas, there is a bowl. In bowl, there are carrot and onion pieces.

Figure 6: More examples about case study.