# VISUAL INFORMATION MATTERS FOR ASR ERROR CORRECTION

*Vanya Bannihatti Kumar, Shanbo Cheng, Ningxin Peng, Yuchen Zhang*

ByteDance

{vanya.bk, chengshanbo, nxpeng, zhangyuchen.zyc}@bytedance.com

## ABSTRACT

Aiming to improve the Automatic Speech Recognition (ASR) outputs with a post-processing step, ASR error correction (EC) techniques have been widely developed due to their efficiency in using parallel text data. Previous works mainly focus on using text or/ and speech data, which hinders the performance gain when not only text and speech information, but other modalities, such as visual information are critical for EC. The challenges are mainly two folds: one is that previous work fails to emphasize visual information, thus rare exploration has been studied. The other is that the community lacks a high-quality benchmark where visual information matters for the EC models. Therefore, this paper provides 1) simple yet effective methods, namely gated fusion and image captions as prompts to incorporate visual information to help EC; 2) large-scale benchmark dataset, namely `Visual-ASR-EC`[1], where each item in the training data consists of visual, speech, and text information, and the test data are carefully selected by human annotators to ensure that even humans could make mistakes when visual information is missing. Experimental results show that using captions as prompts could effectively use the visual information and surpass state-of-the-art methods by upto 1.2% in Word Error Rate(WER), which also indicates that visual information is critical in our proposed `Visual-ASR-EC` dataset.

***Index Terms***— Automatic Speech Recognition, Text correction, Multimodal

## 1. INTRODUCTION

Over the past years, automatic speech recognition (ASR) models have achieved great success [1, 2]. However, there are still many errors in the ASR outputs caused by inherent difficulties, such as grammatical incoherence, homophone errors, etc.

To alleviate grammatical errors, some previous work propose to use an error correction module, typically sequence-to-sequence models, to correct the ASR outputs, with the help of large-scale text data[3]. But since the error patterns vary irregularly based on context, pronunciation and language understanding, it is challenging to construct good quality of
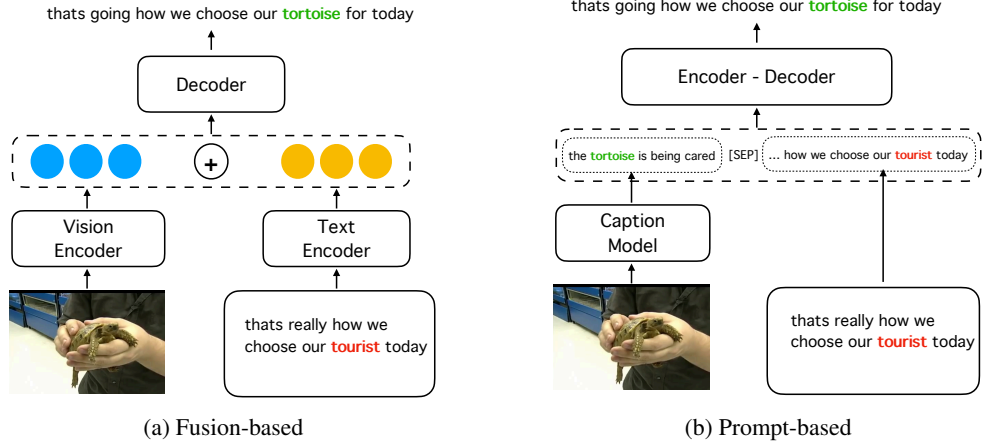
ASR hypotheses to reference text pairs to be used for supervised training. As the homophone errors are more difficult and might not be corrected with only text information, other works propose to utilize audio information in terms of phonetics for correction [4, 5]. Although great works have been studied, there are still certain errors that cannot be properly corrected with audio and text information. For example, the sentence "`if you plan a little bit`" and "`if you plant a little bit`" sounds quite similar to each other, and both are grammatically correct. In such cases, equipping the correction models with only audio and text information is not sufficient, and other modalities, for example, visual information, is critical.

As visual information is crucial for ASR in certain cases, there are some previous works that use visual information for ASR directly, as in [6] or for lip reading task [7]. However, even though visual information could help ASR, the performance gain is limited by the amount of training data, where each training instance should be the triplet of vision, audio, and text information. We believe ASR error correction models, which could leverage visual data and large-scale text data, will improve the performance further.

Error correction can be viewed as similar to machine translation (MT) task where a sequence-to-sequence model like transformers [8] with large amounts of high-quality data can lead to excellent results. Inspired by multi-modal MT and other related tasks like visual text correction [9, 10], in this paper, we propose a multi-modal ASR error correction method which utilizes visual information.

The contributions of this paper are: (1)We present the multi-modal dataset of ASR transcripts along with images to help its correction. The datasets were obtained from two sources, how2 dataset [11] and from the publicly available YouTube videos. Both the sources have reference transcripts annotated by humans. And the ASR transcripts were obtained from the Huggingface wav2vec model and the Google ASR API to show that the method is independent of the ASR model; (2)We propose two ways to utilize the visual information for ASR text correction. Firstly, a gated fusion method where the image features are concatenated with the textual embeddings, similar to previous works [9]. Secondly, we propose a prompt-based method to better utilize large-scale text data, where the captions from the images are used as prompts

---

[1] `https://github.com/VanyaBK/visual_ASR_EC`

thats going how we choose our **tortoise** for today

thats going how we choose our **tortoise** for today

**(a) Fusion-based**

**(b) Prompt-based**

**Fig. 1**: Illustration of the proposed methods for Visual ASR Correction. For the gated fusion method(left), the visual features and text embeddings are concatenated and sent to a decoder to get the corrected output. For the prompt-based method(right), given a source sentence to be corrected, we first generate the corresponding image caption as additional context to the source input, and then use a sequence-to-sequence model to generate the final result.

for ASR correction to provide more context.

## 2. METHOD

### 2.1. Prompt Based Method

Prompt-based methods have gained wide recognition following the success in several NLP tasks [12]. In this method, we use the caption data of the images as prompt for the ASR text correction task in order to provide more context. The caption data is obtained from the Flamingo model [13] trained for image captioning on the Google's Conceptual Captions datasets[2]. For example, consider Fig.1b), here the caption data, "the tortoise is being cared for", is sent as a prompt by modifying the source sentence as, "the tortoise is being cared for [SEP] thats really how we choose our tourist it for today". The modified source sentence and the target sentence are used as the parallel data for training and the modified source sentence is also used as the input while testing. Because of the presence of the caption data as prompt, more context is available to the model to correct `tourist` to `tortoise`, which otherwise would be difficult for humans too.

### 2.2. Gated Fusion Method

Gated fusion techniques are widely used in combining the representations from different modalities as is done in some of the previous works [9]. In this method, for any input sample which consists of image $I$, source text $S$ and target text $T$, the image features are obtained using the OpenAI CLIP's Vision Transformer(ViT) model [14] as ViT($I$) and the textual embeddings are obtained from the standard transformer encoder as $H^S$. These two representations are fused by a vector concatenation noted as :

$$H^{\text{fused}} = g([H^S; \mathbf{H}^I]) \tag{1}$$

where $H^I \in R^{L \times D}$ is the projected form of ViT($I$) to that of the length of the text representation $H^S \in R^{L \times D}$ using a linear projection layer. [L - sequence length; D - hidden dimension]. The fused representation $H^{\text{fused}} \in R^{L \times D}$ is then passed through a `tanh` gate to control the amount of visual information used as :

$$\Lambda = \tanh(f([H^S; H^{\text{fused}}])) \tag{2}$$

The gated fused information is then added to the original textual embeddings to get the multimodal fusion representation as :

$$H^{\text{out}} = H^S + \Lambda H^{\text{fused}} \tag{3}$$

Here tanh gate is used instead of sigmoid since it is centered at zero. This means that when $H^{\text{fused}}$, is close to zero, the output of $\lambda$ will be close to zero, which aligns naturally with the situation when the image is absent as in the synthetic data. Although we found that the results were very similar for both tanh and sigmoid gates.

Since the high-quality annotated set with images is of a smaller size(185k samples) to correct grammatical errors, this method is applied sequentially after first obtaining the result from the baseline transformers model(Transformers+Gated fusion method in Table 1). Finally, the changes made by the gated fusion method to the output from the baseline method (section 3.1), is filtered by including only those changes where the similarity probability(calculated by CLIP's ViT) between the image and the changed text is higher than that with the original text(i.e output from the baseline method)[referred to as Transformers+Gated fusion(Filter) in Table 1]. This is done to retain only image based corrections, since the grammar corrections are already done by the baseline method. Similar filtering and sequential correction is also applied after the prompt-based method [referred to as

---
[2]https://github.com/dhansmair/flamingo-mini

| ID | Models and Variants | Google ASR API | | Huggingface wav2vec | |
|---|---|---|---|---|---|
| | | WER | SER | WER | SER |
| 1 | Original | 36.80 | 100 | 31.94 | 97.79 |
| 2 | Transformer | 34.14 | 99.42 | 22.36 | 91.45 |
| 3 | Prompt-based | **33.5** | **98.94** | **21.13** | **90.59** |
| 4 | Gated fusion | 36.58 | 99.90 | 25.44 | 95.39 |
| 5 | Transformer + Gated fusion | 34.8 | 99.42 | 22.92 | 92.8 |
| 6 | Transformer + Gated fusion (Filter) | 34.54 | 99.42 | 22.66 | 92.51 |
| 7 | Prompt-based + Gated Fusion | 34.21 | 98.94 | 21.80 | 92.12 |
| 8 | Prompt-based + Gated Fusion (Filter) | 33.9 | 99.04 | 21.49 | 91.74 |
| 9 | Random Image captions(Caption Prompt) | 34.2 | 99.14 | 22.11 | 91.16 |
| 10 | Original(Random test set) | 34.59 | 99.56 | 30.70 | 97.6 |
| 11 | Transformer(Random test set) | 32.49 | 98.68 | 21.48 | 90.4 |
| 12 | Prompt-based (Random test set) | **30.79** | **97.8** | **20.6** | **89.10** |

**Table 1**: Measurement of error correction performance

as Prompt-based + Gated Fusion and Prompt-based + Gated Fusion(Filter) respectively in Table 1].

## 2.3. Dataset Retrieval

The datasets were obtained from mainly two sources, the how2 dataset and the youtube videos. The how2 dataset consists of 300h of videos with annotated transcripts in English which resulted in 220,000 samples. The youtube videos were collected using the Youtube-DL toolkit[3], where each of these videos had annotated transcripts and audio in English. The youtube videos accounted for 2.5 million samples of data with annotated transcripts. The images for these samples were obtained by capturing the frame of the video at exactly the mid of the start and end timestamps of that sample. To obtain high-quality annotated dataset from how2 and youtube videos for visual ASR correction, we filtered the dataset to only include the examples where the similarity between the caption data and the reference transcript was greater than 0.2 as measured by Huggingface's sentence-transformers. Similarity score of 0.2 was chosen after testing performance for other thresholds2. After this filtering, 58k samples were obtained from how2 dataset(accounting for 26% of the total dataset) and 127k samples from the youtube videos(accounting for 5% of the total dataset), indicating that a significant portion of the dataset could be corrected using the visual information. Totally, this high-quality annotated dataset consists of 185k samples with audio, images and ASR correction parallel data.

| Similarity Score | Size of train set | WER | SER |
|---|---|---|---|
| 0.1 | 380441 | 33.94 | 98.46 |
| 0.2 | 183306 | **33.5** | **98.94** |
| 0.3 | 85181 | 33.89 | 99.23 |

**Table 2**: Prompt-based method across similarity thresholds

The training set consists of this high-quality annotated dataset along with synthetic EC parallel data. This synthetic dataset is obtained from 3 sources : 1) **TED-LIUM3** - The TED-LIUM3 corpus [15] is built from the TED videos; 2) **DATA2** - The DATA2 corpus [16] is built for the named-entity recognition tasks; 3) **LibriSpeech** - The LibriSpeech corpus [17] is derived from audiobooks of the LibriVox project.

A weak checkpoint trained on the GigaSpeech small dataset [18] using Neural Speech Translation toolkit [19] was used to obtain the synthetic ASR transcripts by setting the beam size in the range of 8-16. Using this method, 5 million synthetic parallel data of candidate ASR transcripts and reference transcripts were obtained.

The test set was constructed to contain 1,041 ASR transcripts from two ASR models, Google ASR API and the Huggingface wav2vec model such that the corresponding image would help in the correction of these transcripts.

## 3. EXPERIMENTS

### 3.1. Baseline

The baseline used for this work is the standard BART-base model trained using Fairseq[4]. This model is first pre-trained with 5 million synthetic data obtained from *TED-LIUM3*, *DATA2* and *LibriSpeech*. Then it is fine-tuned on the high-quality annotated dataset of 185k samples which were filtered based on the similarity metric of sentence-transformers model as explained in section 2.3. GPT-2 based BPE was applied in the pre-processing stage.

### 3.2. Prompt Based Method

For the prompt-based method, a similar setup to that of the baseline was followed. The only change from the baseline was the addition of caption data as prompts in the source sentence in all three train, valid and test sets.

| | | |
|---|---|---|
| **Ex.1** | **Source sentence** | to the best strings for whatever **gucci** youre going to restring what i have on here |
| | **Baseline** | to the best strings for whatever **gucci** youre going to restring what i have on here |
| | **Caption prompt** | to the best strings for whatever **guitar** player youre going to restring what i have on here |
| | **Target sentence** | to the best strings for whatever the **guitar** youre going to restring what i have on here |
| **Ex.2** | **Source sentence** | boxing and yeah it does because boxing kickboxing rg condola training methods |
| | **Baseline** | boxing and yeah it does because boxing kickboxing are condola training methods |
| | **Caption prompt** | boxing and yeah it does because boxing kickboxing are condola training methods |
| | **Target sentence** | boxing and ya it does because boxing kick boxing are **jeet kune do** training methods |

**Table 3**: Sample Analysis

### 3.3. Gated Fusion Method

This method was implemented based on Multimodal Machine Translation where the sigmoid gate function was replaced with tanh. All the parameters were kept constant as in [9], except for the learning rate which was changed to 0.001 and the max updates to 800,000. For evaluation, the average of last 10 checkpoints was used for more reliable results.

### 4. RESULTS AND ANALYSIS

All the experiments are conducted on two different ASR models, Google ASR API and Huggingface wav2vec to verify the generality of the methods. As can be seen from Table 1, the best results are obtained from the **prompt-based method**, thus verifying that visual information helps in ASR EC.

### 4.1. Relevance of Images

To study the relevance of image information in ASR correction, we further conduct more experiments by assigning a random image to the parallel data and seeing if it can help in the ASR correction. As we can see from Table 1, including random image's caption as prompt to the transformer model, leads to a decrease in performance of WER from the prompt-based method, and is on par with the baseline transformer model. This shows that the right image of the video captured during the speech is essential to improve performance.

### 4.2. Why is Caption Prompt Based Method Better?

Because it's easier to use large amounts of synthetic data for text-only methods, the prompt-based method works better than gated fusion. In the gated fusion method it is hard to correct the grammatical errors because of the lack of a large amount of parallel data with the images. In order to emphasize why caption is better, we perform experiments where, with an increase in synthetic data(Table 4), the prompt-based method performs better for Google ASR text correction. This shows that a better representation of the caption is learnt with increasing synthetic data and hence it can be better used as context for ASR text correction.

### 4.3. Sample Analysis

From Ex.1 of Table 3, it is clear that with additional information from the videos, it is easy to correct `gucci` to `guitar`

| ID | Method | 2M Synthetic | | 5M Synthetic | |
|---|---|---|---|---|---|
| | | WER | SER | WER | SER |
| 1 | Prompt-based | 33.94 | 99.33 | **33.5** | **98.94** |

**Table 4**: EC performance with different size of synthetic data

which otherwise would be hard to correct for the transformers model. We also note the limitations of the captioning model which leads to a decrease in the performance of ASR EC. For instance, Ex.2 of Table 3 shows that neither the transformer model nor the caption prompt method is able to predict "jeet kune do because the caption prompt for this particular example is "person a former professional boxer is a trainer and trainer", which do not help in predicting "jeet kune do". Instead, if we have a human annotator providing the caption as "person practicing jeet kune do", it would better help in the ASR EC. Thus, due to the limitations of the captioning model, we observe that the performance of the ASR EC can be further improved with better captioning data.

### 4.4. Random Test Set

Instead of constructing the test set to have cherry-picked examples where images are needed for ASR EC, we test our proposed methods on a test set built by sampling 1000 random sentences from the how2 dataset. From the results of Table 1, we can see that the gap between the baseline and caption prompt methods increases by 2x for the google ASR API and decreases only slightly for huggingface wav2vec(with still a considerable gap of 0.88%). This shows that the caption prompt model has the potential to improve the performance of ASR text correction for instructional videos in general.

### 5. CONCLUSION AND FUTURE WORK

In this work, we have used visual information to aid ASR EC, a method that has not been previously explored. We conducted several experiments to show that visual information can help in ASR EC, which would otherwise be hard to correct for strong baseline models like transformers or even humans without the context. We have introduced simple methods like using captions as prompts, which do not need any modification to the original architecture of transformers, and it improves the WER by upto 1.2% over the baseline methods.

Although we focused on text and visual methods in this work, we believe that incorporating audio information could further enhance our results, which we aim to explore in the future.

## 6. REFERENCES

[1] Alexei Baevski, Wei-Ning Hsu, Alexis CONNEAU, and Michael Auli, "Unsupervised speech recognition," in *Advances in Neural Information Processing Systems*. 2021, vol. 34, pp. 27826–27839, Curran Associates, Inc.

[2] Chung-Cheng Chiu, Tara Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J. Weiss, Kanishka Rao, Katya Gonina, Navdeep Jaitly, Bo Li, Jan Chorowski, and Michiel Bacchiani, "State-of-the-art speech recognition with sequence-to-sequence models," 2018.

[3] Jinxi Guo, Tara Sainath, and Ron Weiss, "A spelling correction model for end-to-end speech recognition," 05 2019.

[4] Shuai Zhang, Jiangyan Yi, Zhengkun Tian, Ye Bai, Jianhua Tao, Xuefei Liu, and Zhengqi Wen, "End-to-end spelling correction conditioned on acoustic feature for code-switching speech recognition," 08 2021, pp. 266–270.

[5] Fan Zhang, Mei Tu, Song Liu, and Jinyao Yan, "Asr error correction with dual-channel self-supervised learning," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 7282–7286.

[6] Dan Oneata and Horia Cucu, "Improving multimodal speech recognition by data augmentation and speech representations," 2022.

[7] Prajwal K R, Triantafyllos Afouras, and Andrew Zisserman, "Sub-word level lip reading with visual attention," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5152–5162, 2022.

[8] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *ArXiv*, vol. abs/1706.03762, 2017.

[9] Bei Li, Chuanhao Lv, Zefan Zhou, Tao Zhou, Tong Xiao, Anxiang Ma, and JingBo Zhu, "On vision features in multimodal machine translation," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Dublin, Ireland, May 2022, pp. 6327–6337, Association for Computational Linguistics.

[10] Amir Mazaheri and Mubarak Shah, "Visual text correction," in *ECCV*, 2018.

[11] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze, "How2: A large-scale dataset for multimodal language understanding," *ArXiv*, vol. abs/1811.00347, 2018.

[12] Yafu Li, Yongjing Yin, Jing Li, and Yue Zhang, "Prompt-driven neural machine translation," in *Findings of the Association for Computational Linguistics: ACL 2022, Dublin, Ireland, May 22-27, 2022*. 2022, pp. 2579–2590, Association for Computational Linguistics.

[13] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Andrew Zisserman, and Karen Simonyan, "Flamingo: a visual language model for few-shot learning," 2022.

[14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," *ArXiv*, vol. abs/2010.11929, 2021.

[15] François Hernandez, Vincent Nguyen, Sahar Ghannay, Natalia A. Tomashenko, and Y. Estève, "Ted-lium 3: twice as much data and corpus repartition for experiments on speaker adaptation," in *SPECOM*, 2018.

[16] Hemant Yadav, Sreyan Ghosh, Yi Yu, and Rajiv Ratn Shah, "End-to-end named entity recognition from english speech," in *INTERSPEECH*, 2020.

[17] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: An asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.

[18] Guoguo Chen, Shuzhou Chai, Guan-Bo Wang, Jiayu Du, Weiqiang Zhang, Chao Weng, Dan Su, Daniel Povey, Jan Trmal, Junbo Zhang, Mingjie Jin, Sanjeev Khudanpur, Shinji Watanabe, Shuaijiang Zhao, Wei Zou, Xiangang Li, Xuchen Yao, Yongqing Wang, Yujun Wang, Zhao You, and Zhiyong Yan, "Gigaspeech: An evolving, multi-domain asr corpus with 10, 000 hours of transcribed audio," in *Interspeech*, 2021.

[19] Chengqi Zhao, Mingxuan Wang, Qianqian Dong, Rong Ye, and Lei Li, "NeurST: Neural speech translation toolkit," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, Online, Aug. 2021, pp. 55–62, Association for Computational Linguistics.