

E-MMAD: Multimodal Advertising Caption Generation Based on Structured Information

Anonymous ACL submission

Abstract

With multimodal tasks increasingly getting popular in recent years, datasets with large scale and reliable authenticity are in urgent demand. Therefore, we present an e-commercial multimodal advertising dataset, E-MMAD, which contains 120 thousand valid data elaborately picked out from 1.3 million real product examples in both Chinese and English. Noticeably, it is one of the largest video captioning datasets in this field, in which each example has its product video (around 30 seconds), title, caption and structured information table that is observed to play a vital role in practice. We also introduce a fresh task for vision-language research based on E-MMAD: e-commercial multimodal advertising generation, which requires to use aforementioned product multimodal information to generate textual advertisement. Accordingly, we propose a baseline method on the strength of structured information reasoning to solve the demand in reality on this dataset.

1 Introduction

Vision-and-Language has been drawing increasing attention from both computer vision and natural language processing communities, for there exists various multimodal information in real human life. As one of the most important tasks of vision-and-language (Uppal et al., 2021), multimodal text generation (Lin et al., 2021) is aimed to generate high-level text by fusing different modal effective information, such as video captioning (Lei et al., 2020a; Yang et al., 2019; Krishna et al., 2017).

However, there are few studies of multimodal text generation making full use of realistic multimodal inputs. One of the reasons is the lack of corresponding publicly available datasets, which can provide real-life multimodal information to help generate. Existing video-text generation datasets are mostly single modal input and are collected by manual batch-written templated descriptions such as MSR-VTT (Xu et al., 2016), Vatec (Wang

et al., 2019). While in practice, information can also be divided into structured information and unstructured information. Humans tend to use richer structured information to generate appropriate text. This information can make the description rigorous and reliable. In this case, a large-scale and reliable dataset with structured information are in urgent demand.

In this paper, we elaborately collect a large-scale e-commercial multimodal advertising dataset for multimodal text generation research, E-MMAD. To support in-depth research, we collect a rich set of product annotations. The E-MMAD dataset consists of 120,984 product instances in both Chinese and English, in which each example has a product video, a title, structured information and a caption. Figure 1 illustrates a sample of our E-MMAD dataset. As is shown in Figure 1, E-commercial multimodal advertising generation task is typically more challenging than existing multimodal text generation, as the advertising description is vivid and information sources are abundant. More importantly, the caption needs to cover the information mentioned in the structured information table but missed in the video.

In response to the realistic demand for advertising generation, we propose the e-commercial multimodal advertising generation task and approach, which is qualified for better performance in generating appropriate text by making full use of the rich information. In addition, considering that various types of information are often encountered in the process of model training and generalization, it will be difficult for the model to train. And in the generalization process, since a considerable part of the nouns do not appear in the training, the caption quality generated by the model is not good enough. For example, when faced with unknown information including new brand names appearing in structured information, the model is not able to effectively identify and judge. So we propose

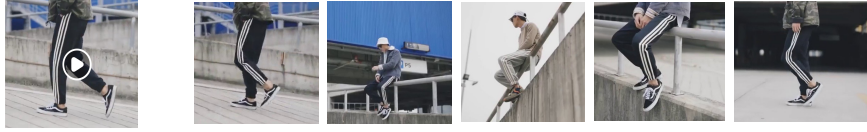

Product Information																																													
Video																																													
Title	春秋薄款 三条纹 运动裤男潮牌小脚收口休闲裤 纯棉 束脚宽松卫裤男 Men's spring autumn thin sweatpants with three stripes fashion brand tapered casual pants men's cotton relaxed tapered sweatpants																																												
Structured Information	<table border="0"> <tr> <td>1. <性别: [男]></td> <td><Gender: [Male]></td> <td>12. <服装版型: [直筒]></td> <td><Garment version: [Straight]></td> </tr> <tr> <td>2. <品牌: [叮目]></td> <td><Brand: [Dingmu]></td> <td>13. <基础风格: [青春流行]></td> <td><Basic style: [Youth Pop]></td> </tr> <tr> <td>3. <风格: [潮牌]></td> <td><Style: [Fashion brand]></td> <td>14. <裤长: [长裤]></td> <td><Pants Length: [Full length]></td> </tr> <tr> <td>4. <图案: [纯色]></td> <td><Pattern: [Pure Color]></td> <td>15. <适用对象: [青年]></td> <td><Target user: [Youth]></td> </tr> <tr> <td>5. <上市时间: [2017年]></td> <td><Time to market: [2017]></td> <td>16. <棉含量: [95%以上]></td> <td><Cotton composition: [Above 95%]></td> </tr> <tr> <td>6. <材质: [纯棉]></td> <td><Material: [Pure cotton]></td> <td>17. <款式元素: [宽松]></td> <td><Style elements: [Relaxed fit]></td> </tr> <tr> <td>7. <设计特点: [三条纹]></td> <td><Design feature: [Three stripes]></td> <td>18. <季节时间: [春秋]></td> <td><Season: [Spring and Autumn]></td> </tr> <tr> <td>8. <厚薄: [加绒加厚]></td> <td><Thickness: [Brushed and thicken]></td> <td></td> <td></td> </tr> <tr> <td>9. <款式细节: [嵌线/贴条]></td> <td><Style details: [Line inlaid/sticker]></td> <td></td> <td></td> </tr> <tr> <td>10. <品类: [运动裤]></td> <td><Category: [Sweatpants]></td> <td></td> <td></td> </tr> <tr> <td>11. <修饰: [收口, 束脚, 薄款]></td> <td><Workmanship: [Close up, Tapered, Thin]></td> <td></td> <td></td> </tr> </table>	1. <性别: [男]>	<Gender: [Male]>	12. <服装版型: [直筒]>	<Garment version: [Straight]>	2. <品牌: [叮目]>	<Brand: [Dingmu]>	13. <基础风格: [青春流行]>	<Basic style: [Youth Pop]>	3. <风格: [潮牌]>	<Style: [Fashion brand]>	14. <裤长: [长裤]>	<Pants Length: [Full length]>	4. <图案: [纯色]>	<Pattern: [Pure Color]>	15. <适用对象: [青年]>	<Target user: [Youth]>	5. <上市时间: [2017年]>	<Time to market: [2017]>	16. <棉含量: [95%以上]>	<Cotton composition: [Above 95%]>	6. <材质: [纯棉]>	<Material: [Pure cotton]>	17. <款式元素: [宽松]>	<Style elements: [Relaxed fit]>	7. <设计特点: [三条纹]>	<Design feature: [Three stripes]>	18. <季节时间: [春秋]>	<Season: [Spring and Autumn]>	8. <厚薄: [加绒加厚]>	<Thickness: [Brushed and thicken]>			9. <款式细节: [嵌线/贴条]>	<Style details: [Line inlaid/sticker]>			10. <品类: [运动裤]>	<Category: [Sweatpants]>			11. <修饰: [收口, 束脚, 薄款]>	<Workmanship: [Close up, Tapered, Thin]>		
1. <性别: [男]>	<Gender: [Male]>	12. <服装版型: [直筒]>	<Garment version: [Straight]>																																										
2. <品牌: [叮目]>	<Brand: [Dingmu]>	13. <基础风格: [青春流行]>	<Basic style: [Youth Pop]>																																										
3. <风格: [潮牌]>	<Style: [Fashion brand]>	14. <裤长: [长裤]>	<Pants Length: [Full length]>																																										
4. <图案: [纯色]>	<Pattern: [Pure Color]>	15. <适用对象: [青年]>	<Target user: [Youth]>																																										
5. <上市时间: [2017年]>	<Time to market: [2017]>	16. <棉含量: [95%以上]>	<Cotton composition: [Above 95%]>																																										
6. <材质: [纯棉]>	<Material: [Pure cotton]>	17. <款式元素: [宽松]>	<Style elements: [Relaxed fit]>																																										
7. <设计特点: [三条纹]>	<Design feature: [Three stripes]>	18. <季节时间: [春秋]>	<Season: [Spring and Autumn]>																																										
8. <厚薄: [加绒加厚]>	<Thickness: [Brushed and thicken]>																																												
9. <款式细节: [嵌线/贴条]>	<Style details: [Line inlaid/sticker]>																																												
10. <品类: [运动裤]>	<Category: [Sweatpants]>																																												
11. <修饰: [收口, 束脚, 薄款]>	<Workmanship: [Close up, Tapered, Thin]>																																												
 Advertising																																													
<p>采用优质纯棉面料, 亲肤舒适透气性强。裤线两边的三道纹设计更显学院复古风, 斜插口袋设计更显方便。收裤脚的裤边显出男神干净利落的时尚气质, 中腰的款式让穿着的舒适感更高。纯色系的运动裤是你衣橱中必不可少的百搭单品。</p> <p>The pants are made of high-quality pure cotton fabric with a pretty comfy and breathable skin feel. The three-stripe design on both sides shows college retro style, and the design of side pockets is more convenient. The hem of tapered pants appears the spruce and fashion style of men, while the mid waisted fit provides a more comfortable wear. Pure-color sweatpants are necessary all-match item in your wardrobe.</p>																																													

Figure 1: One example of our E-MMAD dataset. The four different colors represent the four different parts of our dataset, from top to bottom are product information (commodity displaying video, title, structured information) and commodity advertising description. The task of our model is to use the product information to generate corresponding advertising description. We add structured information to the original Video Caption to assist in generating a semantically richer caption. The colored words will be presented in the final generated caption.

Conceptual Operations to deal with complex and diverse information in real life.

In summary, our contributions concentrate on the following three aspects:

- (1) We collect a large-scale high-quality and reliable e-commercial multimodal advertising dataset. It is one of the largest video captioning datasets in this field. E-MMAD is completely collected from human real life and carefully selected so that it is qualified to meet the needs of real life.
- (2) We introduce a fresh task for vision-language research based on E-MMAD: e-commercial multimodal advertising generation, which requires to use the product multimodal information to generate textual advertisement.
- (3) We propose a simple but effective baseline method on the strength of structured information reasoning to solve the demand in reality on E-MMAD dataset.

2 Related Work

2.1 Multimodal video-text generation datasets

There are various datasets for multimodal video-text generation that cover a wide range of domains,

such as movies (Rohrbach et al., 2015), cooking (Das et al., 2013; Zhou et al., 2018a), and Activities (Xu et al., 2016). MSR-VTT (Xu et al., 2016) is a widely-used dataset for video captioning, which has 10,000 videos from 257 activities and was collected in 2016. MSVD (Chen and Dolan, 2011) was collected in 2011, containing 1970 videos. ActivityNet (Caba Heilbron et al., 2015) has 20,000 videos but is used for Dense Video Captioning (Krishna et al., 2017), which means to describe multiple events in a video. TVR (Lei et al., 2020b) is collected from movie clips whose text is mainly character dialogue. VateX (Wang et al., 2019) is a famous dataset released in 2019, whose caption is written by batch manpower. Compared with some mainstream datasets in Table 1, our dataset also provide an additional structured information table. And the generated caption needs to include the information mentioned in the structured information.

2.2 Video Captioning Approaches

Video caption/description is one of the important tasks in multimodal text generation. Early video caption methods are all based on templates

(Mitchell et al., 2012; Krishnamoorthy et al., 2013). However, sentences made in this way tend to be rigid and stiff. The sequence-to-sequence model (Venugopalan et al., 2015) is a classic work, which includes an encoding phase and a decoding phase. After CNN extracts the image features of the video frames, an image feature is sent to the LSTM for encoding at each time step. What needs to be generated in the decoding stage is text. Some of the popular practices recently are based on data-driven (Zhang et al., 2021b) and transformer-based mechanisms (Yang et al., 2019; Zhou et al., 2018b; Lei et al., 2020a). MART (Lei et al., 2020a) can produce more coherent, non-repetitive, and relevant text to enhance the transformer architecture by using memory storage units. Vx2text (Lin et al., 2021) uses multimodal inputs for text generation. They use a backbone(Tran et al., 2018; Ghadiyaram et al., 2019) model to transform different modals information to nature language and then the problem turns to nature language generation. Although good progress has been made by them, the original information of the modal is not fully utilized and integrated.

3 Datasets

In this section, we will introduce our dataset in detail, including the statistic analysis, collecting process, and comparison.

3.1 Data Collection

1) Dataset sources. Our dataset sources are the Chinese largest e-commerce website shopping platform (www.taobao.com), from which we have collected nearly 1.3 million commodity examples with structured information. It comprised more than **4,000** merchandise categories to guarantee the diversity of the dataset, such as clothes, furniture, office supplies, etc. The information of each commodity data sample includes structured information, commodity displaying video, title of product and commodity advertising description. Different from previous works (Wang et al., 2019; Xu et al., 2016; Chen and Dolan, 2011), The sources of datasets are derived from what merchants themselves numerously design and select, which comply with the standard rules of the authenticity of product advertisements and are supervised by false product advertising rules of *Taobao*. Specifically, videos visually display the commodity performance and application.

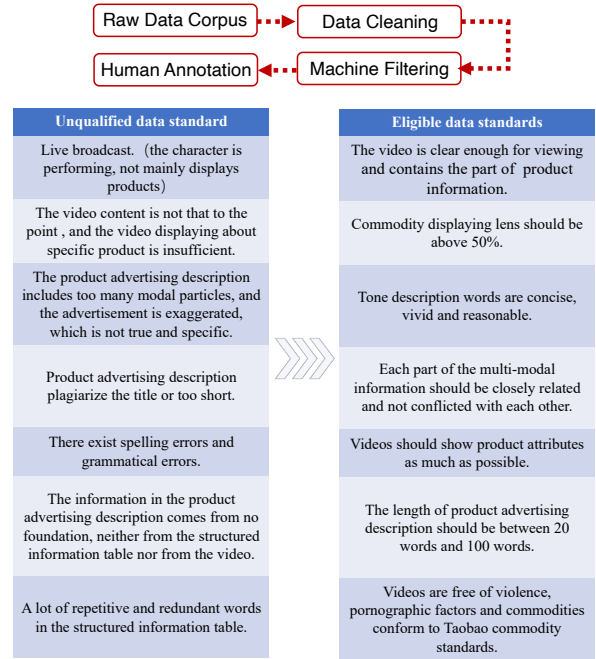


Figure 2: The process of creating a dataset, including cleaning, machine filtering, manual post-filtering, etc. and data specification of the dataset.

In addition, we fully consider ethical privacy issues to ensure that the dataset has no potential negative effects and legal issues(Gebru et al., 2018). All data is collected in *Taobao* shopping platform, which is a public platform for the general public. All information, even the characters in the video, is ensured to comply with *Taobao laws* including personal privacy, legal prohibitions, false information, protection of minors and women, and so on.

In consideration of data and ethics, we perform programmatic screening and manual cleaning again in accordance with the established data cleaning rules. Figure 2 shows our data collection process.

2) Data filtering. The intention for data filtering is to determine whether the product advertising description is closely related to the product displaying video, and whether the structured information of the product is in accordance with the composition of the product advertising description and ethical considerations. The product attributes structured information and product displaying video will be valid only if human being can write similar product advertising descriptions with them. We use programs to screen and judge at first, traversing the values of structured information. Our screening basis is the proportion of structured information words in the product advertising description. When the proportion is up to n words or more, the

data will be reserved as valid data. After copywriters' continuous attempt to generate advertising descriptions with structured information words that account for different proportions, we finally determine the structured information with more than five words in the product advertising description as valid data and form **207,852** machine-screened data.

By virtue of this, we respectively test different groups of random data to formulate screening and judgment rules. Several times our different copywriters have tested and discussed to make the manual evaluation criterion. Consequently, testers sample **100** examples randomly according to the judgment rules of **Figure 2**, and the pass rate is mostly about 60%. In this case, we validate the manual screening rules and draw the conclusion that random subjective factors hardly have any influence. So far, the manual data screening and judging rules have been formed, as is shown in **Figure 2**.

3) Data annotation. We invited 25 professional advertising copywriters as data screening and annotation staff to conduct manual screening under the rules of **Figure 2** and the Toronto Declaration. Manual screening of all data also ensures that each piece of data complies with the Toronto Declaration and *Taobao* laws to protect gender equality, racial equality, etc. In order to ensure the reliability of the data, we use the following two methods to sample and verify: (1). Add verification steps. We will send back samples that have been annotated right answers to annotators from time to time to check their work quality. (2). Multiple people Choices. The data is sent to different people randomly. Only if the answers of multiple people are consistent, can they be passed. Finally, **120,984** valid data has been generated. Simultaneously, we also translate the filtrated valid data into English so that both Chinese and English versions can be provided in the dataset. To ensure the quality of the English version, we use the WMT2019 Chinese-English translation champion, *Baidu machine translation*. We also monitor the translation quality in the manual screening section, such as random checking in batch translation, using text error correction to monitor retranslation, and back translation comparison. Of course, we mainly encourage and urge people to study natural language research in the Chinese e-commerce market.

After the consumption of 25-people toil for manual data labeling and cleaning, **120,984** carefully

chosen valid data have finally been generated.

3.2 Dataset Analysis

Among the **207,852** data we send for annotation, there are **120,984** eligible samples passing the screening. We make an elaborate analysis on these valid data and the result is shown in **Figure 3**. In addition to this, **Figure 3** reveals the distribution of the product videos' duration and advertising descriptions. By **Table 1** comparison, we can find that

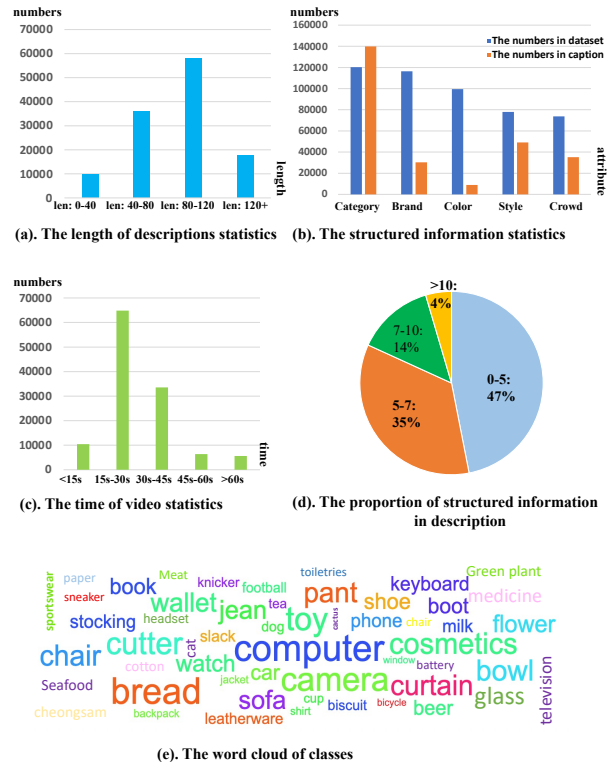


Figure 3: Statistics about the five different forms of data in our dataset. The data statistics are presented in terms of video, structured information, and the main classes of the dataset contained, respectively.

our product advertising descriptions are not only at least twice longer than others, but also root in more vivid and realistic ones used in practice. The whole statistics about the structured information in our dataset is displayed in **Figure 3(d)**. What's more, there exist average 21 structured information words in each sample and 6.2 words of them are finally displayed in its product advertising description. The **(e)** shows the abundance of our datasets source classes.

3.3 Dataset Comparison

In **Table 1**, we make a comparison between our dataset and others from the following several

Name.	Data Size	Video Time	Length	Classes	Modality	SI
MSR-VTT(Xu et al., 2016)	10,000	10s - 20s	9	257	Video to Text	×
MSVD(Chen and Dolan, 2011)	1,970	9s	8	-	Video to Text	×
TVR(Lei et al., 2020b)	21800	9s	13	-	Multimodal to Text	×
VaTEX(en/zh)(Wang et al., 2019)	41,269	10s	15/13	600	Video to Text	×
Ours(en/zh)	120,984	30s	97/67	4,863	Multimodal to Text	✓

Table 1: Comparison with other datasets. Each column represents a different piece of information about the dataset. *Data Size*, *Video Time*, *Length*, *Classes* respectively represent the total number of videos in the dataset, the range of video lengths in the dataset, the average length of the captions in the dataset and the number of video types in the dataset. *Modality* indicates the use of the dataset, e.g. from Video to Text, Multimodal to Text. *SI* means whether the dataset contains structured information.

perspectives: dataset scale, dataset diversity and dataset reliability.

1)Dataset scale: As shown in **Table 1**, the size of our **E-MMAD** is the largest multimodal dataset among those we have already known so far, with the longest video duration and text length, and the richest structured information in the dataset.

2)Dataset Diversity: In terms of types, our dataset consists of **4,863** categories, which is currently the most various data sources in the datasets. Our dataset is also available in Chinese and English two versions, to support multi-language research, which cannot be satisfied by a single language dataset. At the same time, our Chinese and English corpus is richer in vocabulary, which can generate more natural and diversified video descriptions.

3)Dataset Reliability: Compared with other manual batch-written descriptions(Wang et al., 2019) and mechanically generated data, our data annotation is derived from the real society. Each of them is an exclusive description genuinely written by corresponding store. Besides, the videos in our dataset are from the real product shooting scene, other than clips from Youtube or movies. We firmly believe that only resorting to reliable dataset, can we train models better. Therefore, we invest considerable amount of manpower and time in order to promote our dataset quality.

3.4 Dataset Significance

To the extent of our knowledge, the dataset we propose is the largest multi-modal dataset so far, and the information involved is also the most diverse, which can better optimize and improve the performance of multi-modality models and promote their generalization ability to adapt to different scenarios in real world. For subsequent work, with the abundant and diverse information involved, our dataset

can be dedicated to several multi-modality domain tasks, such as Video Retrieval(Lei et al., 2020b), Product Search(Chang et al., 2021) and so on. In our future work, we will build more versatile e-commerce datasets which can cover most tasks in this field based on this dataset.

4 Method

In this work, we present a novel approach called the Multi-modal Fusion and Generation algorithm as shown in **Figure 4**, which extracts feature representations from three sources: the product title, structured information(structured words) and the displaying video’s frames and fuse them to generate captions. And to process various information, our model use a method of conceptualizing information. That is to pre-process the data, conceptualize and extract information from the complex information to highly conceptualize network features. For the restoration of complex information in the generation phase, we only need to perform the inverse conceptualization operation at the end.

4.1 Conceptualization

During the training process, we pre-conceptualize the true product descriptions. The formula is as follows:

$$Values_{gr} = SW.values \cap GR.tokens \quad (1)$$

$$k_{gr} \in SW.keys \quad (2)$$

$$\begin{aligned} token_{gr} &\rightarrow k_{gr} \\ \forall token_{gr} &\in Values_{gr} \end{aligned} \quad (3)$$

In the generation process, the raw caption with conceptualized information generated by the model is de-conceptualized to obtain the final caption. The

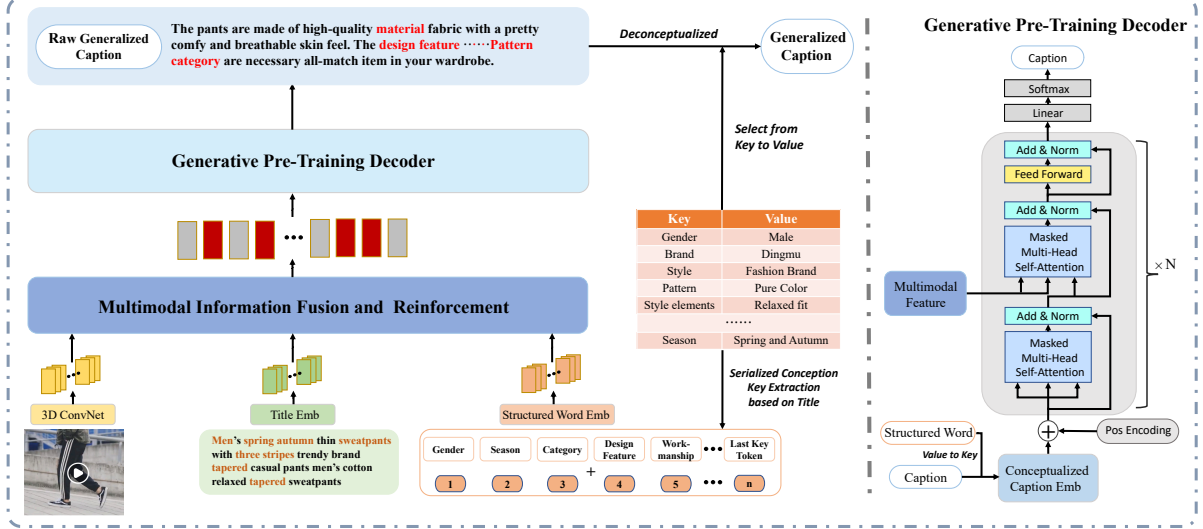


Figure 4: The overall architecture of our model, which contains three main parts: the representation for multimodal information, the multimodal fusion module based on self-attention and the generative pre-training decoder module on the basis of (Radford et al., 2019).

de-conceptualization is as follows:

$$Values_{rc} = SW.keys \cap RC.tokens \quad (4)$$

$$v_{gr} \in SW.values \quad (5)$$

$$rc_token \rightarrow v_{gr} \quad (6)$$

$$\forall rc_token \in Values_{rc}$$

Among them, $A \rightarrow B$ means replacing token A with token B . $A \in C$ means token A is an element of set C . $GR.tokens$ and $RC.tokens$ are the sets of corresponding n-gram phrases in ground truth and raw caption, respectively. $SW.values$ and $SW.keys$ respectively correspond to the sets of keys and values in the structured information. In terms of the model input, for the structured words part, we only extract the key components, and reference the title as the basis to determine the priority of each key according to the order in which the structured information appears in the title.

4.2 Representation

Textual Information. Given a product title as a list of K words, conceptualized product attributes as a list of N keys, we embed these words and keys into the corresponding sequence of d -dimensional feature vectors using trainable embeddings (Zhang et al., 2021a; Devlin et al., 2018). In addition, since the keys of structured words are prioritized, we use *position embedding* to represent the priority of the keys.

Visual Information. Given a sequence of video frames/clips of length S , we feed it into pre-trained 3D CNNs (Ji et al., 2012) to obtain visual features $V = \{v_1, v_2, \dots, v_K\} \in \mathbb{R}^{S \times d_v}$, which are further encoded to compact representations $R \in \mathbb{R}^{S \times d}$, which have the same dimension as the representation of textual information via a *Visual Embedding Layer*. The *Visual Embedding Layer* can be formalized as following:

$$f_{VEL}(v) = BN(g \circ \bar{v} + (1 - g) \circ \hat{v}) \quad (7)$$

$$\bar{v} = W_1 v^\top \quad (8)$$

$$\hat{v} = \tanh(W_2 \bar{v}) \quad (9)$$

$$g = \sigma(W_3 \bar{v}) \quad (10)$$

where BN denotes batch normalization, \circ is the element-wise product, σ means sigmoid function, $W_1 \in \mathbb{R}^{d \times d_v}$ and $\{W_2, W_3\} \in \mathbb{R}^{d \times d}$ are learnable weights.

4.3 Multimodal Fusion

After embedding all information from each modality as vectors in the d -dimensional joint embedding space, we use a stack of L transformer layers with a hidden dimension of d to fuse the multi-modal information consisting of a list of all $K + N + S$ modalities from $\{v_S^{frames}\}$, $\{v_K^{words}\}$ and $\{v_N^{keys}\}$. Through the self-attention mechanism in transformer, we can model inter- and intra- modality context. The output from our Multimodal Information Fusion and Reinforcement module is a list of

d -dimensional feature vectors for entities in each modality, which can be seen as their interrelated embedding in multimodal context. In this work, the parameters chosen for our the module are consistent with the parameters of *BERT-base* ($L=12$, $H=768$, $A=12$), where L , H , A represents the number of layers, the hidden size, and the number of self-attention heads respectively.

4.4 Generation Decoder

Our model’s decoder is a left-to-right Transformer decoder, which is similar to the model architecture of (Chen et al., 2019). The decoder access multimodal fusion outputs at each layer with a multi-head attention (Vaswani et al., 2017). Specifically, the decoder applies a multi-headed self-attention over the caption textual feature. After that, the position-wise feedforward layer was used to produce a distribution probability of each generation tokens for the final generated caption. There is a description of part of the formula for the decoder module:

$$h_0 = V^{\text{cap}} \cdot W_t + PE \cdot W_p \quad (11)$$

$$h_l = \text{Trans_Block} (h_{l-1}) \quad (12)$$

$$P(w) = \text{Softmax} (h_n W_e^T) \quad (13)$$

$$PE_{(pos,2i)} = \sin \left(pos/10000^{2i/d_{\text{model}}} \right) \quad (14)$$

$$PE_{(pos,2i+1)} = \cos \left(pos/10000^{2i/d_{\text{model}}} \right) \quad (15)$$

where $V^{\text{cap}} = \{v_1, v_2, \dots, v_x\}$ is the textual vector of caption, n is the number of layers, $\forall l \in [1, n]$, and W_t, W_p is the learnable weight for caption embedding feature and position encoding respectively. *Trans_Block* represents a block of the decoder in the Transformer (Vaswani et al., 2017). We refer to (Vaswani et al., 2017; Radford et al., 2018, 2019; Chen et al., 2019) for a more detailed explanation of the model architecture.

5 Experiments

In this section, we will show a series of experiments of our proposed model on E-MMAD, including ablation studies, comparison experiments and state-of-the-art video caption methods and human evaluation.

5.1 Implementation Details

All the experiments are conducted on Nvidia TitanX GPU. The proposed model is implemented with PyTorch. For the representations of videos,

we follow (Yang et al., 2019) for fairness and opt for the same type, first extract 3D features with 2048 dimensions, 2048-D image features from ResNet-101 (Hara et al., 2017) pre-trained on ImageNet dataset. For generation decoder, we use $\langle \text{sep} \rangle$ to separate the input from the ground truth of caption. And, we adopt diverse automatic evaluation metrics to compare with other model: BLEU (Papineni et al., 2002), Rouge-L (Lin, 2004), and CIDEr (Vedantam et al., 2015). It is worth noticing that the focus of the CIDEr evaluation metric is on whether the generated caption captures the major information or not. Since the major information captured by each model is different, the key information component of the generated caption will not be the same, but it is cognitive at the semantic level, so the CIDEr evaluation metric will have a relatively large fluctuation. Our model introduces structured information so that the generated caption can include most of the major information. Therefore, the caption generated by our model can achieve significant results in the evaluation index of CIDEr.

5.2 Comparison with Other Approaches

During the comparison experiments, we uniformly divided the Chinese and English versions of our dataset into training set, validation set and test set in the ratio of 6:2:2 for training and testing. Since the current mainstream models do not use multimodal data for captioning, we use unimodal data for captioning on some classic and available methods, such as video caption, NLG, etc. Also, for the sake of fairness of comparison, we simply modify the input part (NACF*) of the above experimental model to accommodate multimodal data. As we can see from **Table 2**, the comparison of the results before and after the model modification shows that multimodal data can be substantially improved for text generation tasks. It indicates that multimodal information indeed helps captioning by modal information between the mutual enhancement. And as shown in **Table 2** our algorithm achieves a better performance than other methods because our model makes better use of multimodal data in the means of fusing different modalities and structured information to reason.

5.3 Ablation studies

Multimodal Input. We perform ablation studies based on changing the input components of our proposed model as a way to validate the impor-

Table 2: Performance (%) comparison with our proposed model and others. The NACF* means that we concat the structured information with video feature directly. On the premise of fair comparison, the following methods are relatively classic and available, which are applicable on E-MMAD by our objective attempts.

Method	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
NLG(Chen et al., 2019)	13.6	6.8	3.1	1.9	13	10.1
NACF(Yang et al., 2019)	18.9	7.9	3.9	2.2	15.3	14.8
NACF*	20	8.5	4.3	2.4	17.8	18.5
TVC(Lei et al., 2020b)	21.3	12.4	6.2	3.7	19.3	22.5
Ours(en)	25.0	16.6	9.6	7.2	25.3	29.1
CPM(zh)(Zhang et al., 2021a)	7.9	4.6	1.1	0.5	7.2	8.3
Ours(zh)	11.6	6.5	4.4	2.2	12.5	15.3

498 tance of our proposed dataset containing structured
 499 information. As shown in **Table 3**, we analyze the
 500 gap between the generated caption of the model
 501 and the real commodity advertising description in
 502 the absence of partial information. As we can see,
 503 the absence of any of the three input components
 504 significantly degrades the final generated caption
 505 result. From our analysis of the generated caption,
 506 we can conclude that: 1) the lack of structured
 507 information will make the generated caption less
 508 informative, rigorous and reliable. 2) The lack of a
 509 commodity title or displaying video will impair the
 510 foundation of generated text. In addition, the struc-
 511 tured information is like a knowledge base, which
 512 can promote inference and judgment to generate
 513 appropriate caption.

514 **Conceptual Operation.** Considering that writ-
 515 ing product descriptions in real life often involves
 516 a great number of unfamiliar words, which makes
 517 it hard for the model to identify and remember
 518 its feature when facing a new word, such as new
 519 brand name. The predecessor’s approach tend to
 520 use as much corpus and large model parameters as
 521 possible, which brings huge difficulties to natural
 522 language generation. In this case, we proposed the
 523 Conceptualization operation. As shown in **Table 4**
 524 , we conduct ablation experiments about Conceptu-
 525 alization on the Chinese and English datasets. As
 526 for models without conceptual operations, we use
 527 unconceptualized captions as the ground truth to
 528 train. And for the input of the model, we directly in-
 529 put unordered structured words. Experiments have
 530 proved that the Conceptualization operation can
 531 indeed bring a significant effect improvement, be-
 532 cause this method can conceptualize and extract in-
 533 formation from complex information in the dataset,
 534 and thus highly conceptualize network features.
 535 We expect this discovery to inspire the community.

5.4 Human Assessment

536 It is well-known that the human evaluation metrics
 537 for video captioning are required due to the inaccur-
 538 ate evaluation by automatic metrics. We especially
 539 focus on advertising generation, which depend on
 540 human aesthetics. So we invite the people involved
 541 in the data annotation and new advertising slogan
 542 designers to conduct the human evaluation. We se-
 543 lect 200 samples from the test dataset and each eval-
 544 uator evaluate each of these 200 samples to reflect
 545 the performance of our model by rating whether
 546 the caption generated by our model can be used as
 547 a description of the product. As the result shows
 548 in **Table 5**, the caption generated by our model has
 549 a certain degree of usability, whose results were
 550 generally recognized by people. Therefore, this is
 551 also acceptable that our experiments on **Table 2** did
 552 not achieve high scores for mechanical evaluation
 553 indicators.
 554

6 Conclusion and Future Work

555 This research sets out to provide an e-commercial
 556 multimodal advertising dataset, E-MMAD, which
 557 is one of the largest video captioning datasets in this
 558 field. Based on E-MMAD, we also present a fresh
 559 task: e-commercial multimodal advertising genera-
 560 tion, and propose a baseline method on the strength
 561 of structured information reasoning to solve the
 562 realistic demand. We hope the release of our E-
 563 MMAD would facilitate the development of multi-
 564 modal generation problems in the real world. How-
 565 ever, there still exist limitations about our dataset
 566 and method that should be acknowledged. Mov-
 567 ing forward, we are planning to extend E-MMAD
 568 to better performance and more diversified tasks
 569 by exploring new model structures, using different
 570 language data and so on.
 571

References

Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970.

Wei-Cheng Chang, Daniel Jiang, Hsiang-Fu Yu, Choon Hui Teo, Jiong Zhang, Kai Zhong, Kedarnath Kolluri, Qie Hu, Nikhil Shandilya, Vyacheslav Ievgrafov, et al. 2021. Extreme multi-label learning for semantic matching in product search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2643–2651.

David Chen and William B Dolan. 2011. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th annual meeting of the association for computational linguistics: human language technologies*, pages 190–200.

Zhiyu Chen, Harini Eavani, Wenhui Chen, Yinyin Liu, and William Yang Wang. 2019. Few-shot nlg with pre-trained language model. *arXiv preprint arXiv:1904.09521*.

Pradipto Das, Chenliang Xu, Richard F Doell, and Jason J Corso. 2013. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2634–2641.

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*.

Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint arXiv:1803.09010*.

Deepti Ghadiyaram, Du Tran, and Dhruv Mahajan. 2019. Large-scale weakly-supervised pre-training for video action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12046–12055.

Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. 2017. Learning spatio-temporal features with 3d residual networks for action recognition. *arXiv preprint, arXiv:1708.07632*.

Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 2012. 3d convolutional neural networks for human action recognition. volume 35, pages 221–231. IEEE.

Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*, pages 706–715.

Niveda Krishnamoorthy, Girish Malkarnenkar, Raymond Mooney, Kate Saenko, and Sergio Guadarrama. 2013. Generating natural-language video descriptions using text-mined knowledge. In *Twenty-Seventh AAAI Conference on Artificial Intelligence*. 626
627
628
629
630

Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L Berg, and Mohit Bansal. 2020a. Mart: Memory-augmented recurrent transformer for coherent video paragraph captioning. *arXiv preprint arXiv:2005.05402*. 631
632
633
634
635

Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2020b. Tvr: A large-scale dataset for video-subtitle moment retrieval. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXI 16*, pages 447–463. Springer. 636
637
638
639
640
641

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81. 642
643
644

Xudong Lin, Gedas Bertasius, Jue Wang, Shih-Fu Chang, Devi Parikh, and Lorenzo Torresani. 2021. Vx2text: End-to-end learning of video-based text generation from multimodal inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7005–7015. 645
646
647
648
649
650

Margaret Mitchell, Jesse Dodge, Amit Goyal, Kota Yamaguchi, Karl Stratos, Xufeng Han, Alyssa Mensch, Alexander Berg, Tamara Berg, and Hal Daumé III. 2012. Midge: Generating image descriptions from computer vision detections. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 747–756. 651
652
653
654
655
656
657
658

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318. 659
660
661
662
663

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. 664
665
666

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9. 667
668
669
670

Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. 2015. A dataset for movie description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3202–3212. 671
672
673
674

Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. 2018. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459. 675
676
677
678
679
680

681 Shagun Uppal, Sarthak Bhagat, Devamanyu Hazarika,
682 Navonil Majumder, Soujanya Poria, Roger Zimmer-
683 mann, and Amir Zadeh. 2021. Multimodal research
684 in vision and language: A review of current and
685 emerging trends. *Information Fusion*.

686 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
687 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
688 Kaiser, and Illia Polosukhin. 2017. Attention is all
689 you need. In *Advances in neural information pro-
690 cessing systems*, pages 5998–6008.

691 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi
692 Parikh. 2015. Cider: Consensus-based image de-
693 scription evaluation. In *Proceedings of the IEEE
694 conference on computer vision and pattern recogni-
695 tion*, pages 4566–4575.

696 Subhashini Venugopalan, Marcus Rohrbach, Jeffrey
697 Donahue, Raymond Mooney, Trevor Darrell, and
698 Kate Saenko. 2015. Sequence to sequence-video
699 to text. In *Proceedings of the IEEE international
700 conference on computer vision*, pages 4534–4542.

701 Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-
702 Fang Wang, and William Yang Wang. 2019. Vatex:
703 A large-scale, high-quality multilingual dataset for
704 video-and-language research. In *Proceedings of the
705 IEEE/CVF International Conference on Computer
706 Vision*, pages 4581–4591.

707 Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-
708 vtt: A large video description dataset for bridging
709 video and language. In *Proceedings of the IEEE con-
710 ference on computer vision and pattern recognition*,
711 pages 5288–5296.

712 Bang Yang, Yuexian Zou, Fenglin Liu, and Can Zhang.
713 2019. Non-autoregressive coarse-to-fine video cap-
714 tioning. *arXiv preprint arXiv:1911.12018*.

715 Zhengyan Zhang, Xu Han, Hao Zhou, Pei Ke, Yuxian
716 Gu, Deming Ye, Yujia Qin, Yusheng Su, Haozhe
717 Ji, Jian Guan, et al. 2021a. Cpm: A large-scale
718 generative chinese pre-trained language model. *AI
719 Open*, 2:93–99.

720 Ziqi Zhang, Zhongang Qi, Chunfeng Yuan, Ying Shan,
721 Bing Li, Ying Deng, and Weiming Hu. 2021b. Open-
722 book video captioning with retrieve-copy-generate
723 network. In *Proceedings of the IEEE/CVF Confer-
724 ence on Computer Vision and Pattern Recognition*,
725 pages 9837–9846.

726 Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018a.
727 Towards automatic learning of procedures from web
728 instructional videos. In *Thirty-Second AAAI Confer-
729 ence on Artificial Intelligence*.

730 Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard
731 Socher, and Caiming Xiong. 2018b. End-to-end
732 dense video captioning with masked transformer. In
733 *Proceedings of the IEEE Conference on Computer
734 Vision and Pattern Recognition*, pages 8739–8748.

A Appendix Ablation Results Tables

Source Link: <https://github.com/E-MMAD/E-MMAD>

Table 3: Performance comparison with our proposed model by masking different parts of input and only using the remainder as input. Here "Title", "SW" and "Video" indicates commodity title, attribute structured word table and commodity displaying video respectively.

Input	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
SW & Video	22.8	14.8	6.9	5.5	22.2	25.3
Title & Video	19.5	9.4	4.5	3.1	16.4	15.7
Video	15.9	6.4	3.4	2.1	15	13.2
Title&SW	22.0	13.8	5.8	4.9	20.6	23.7

Table 4: Performance comparison of whether our proposed model has conceptual operations(CO).

Operation	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
No CO(en)	23.8	15.4	8.1	6.4	24.2	27.3
No CO(zh)	9.9	5.5	2.8	1.5	10.1	12.4
With CO(en)	25.0	16.6	9.6	7.2	25.3	29.1
With CO(zh)	11.6	6.5	4.4	2.2	12.5	15.3

Table 5: The results of the human evaluation, reflecting the proportion of the 200 examples where the model generated caption could be used as a product description that describes the reasonableness of the generated caption.

	Annotator1	Annotator2	Annotator3	Person1	Person2	Person3
Pass	42%	44%	43%	48%	56%	53%