

Attract me to Buy: Advertisement Copywriting Generation with Multimodal Multistructured Information

Anonymous Author(s)
Submission Id: 1606

ABSTRACT

Recently, online shopping has gradually become a common way of shopping for people all over the world. Attractive commodity advertisements can attract more people to buy, which often integrate multimodal multistructured information of commodities, such as visual spatial information and fine-grained structure information. However, traditional multimodal text generation focused on the routine description of what existed and happened, which did not match the real-world advertising copywriting. Advertisement copywriting has a vivid language style and higher requirements of faithfulness. Unfortunately, there is a lack of reusable evaluation frameworks and a scarcity of datasets. Therefore, we present a dataset, E-MMAD (e-commercial multimodal multistructured advertisement copywriting), which requires, and supports much more detailed information in text generation. Noticeably, it is one of the largest video captioning datasets in this field. Accordingly, we propose a baseline method and metric on the strength of structured information reasoning to solve the demand in reality on this dataset. We achieve SOTA performance on all metrics. We will release the dataset and method to promote further investigations on both multimodal text generation and e-commerce advertisement.

CCS CONCEPTS

• **Computing methodologies** → *Natural language generation*;
• **Video summarization**; • **Applied computing** → **Electronic commerce**.

KEYWORDS

datasets, multimodal, structure information

1 INTRODUCTION

Nowadays, online shopping has become one of the main ways for people to shop, such as Taobao, Amazon. The product advertisement is often an important factor in people’s shopping. Wonderful product advertising can attract people’s attention and promote sales. Commodity advertisement copywriting[?] presents commodities in a more concise and intuitive way, which is convenient for people to search and shop. Different from conventional text generation [6, 24], commodity advertisement copywriting often has vivid language style and flexible grammar. Meanwhile, it also needs to comprehensively consider multimodal information[?] and fine-grained structural information[32] parameters of commodities, which results in sellers often need to spend a lot of manpower, time and money on elaborate design to produce high-quality advertising copywriting. We named this more challenging problem as multi-modal e-commerce advertisement copywriting generation.

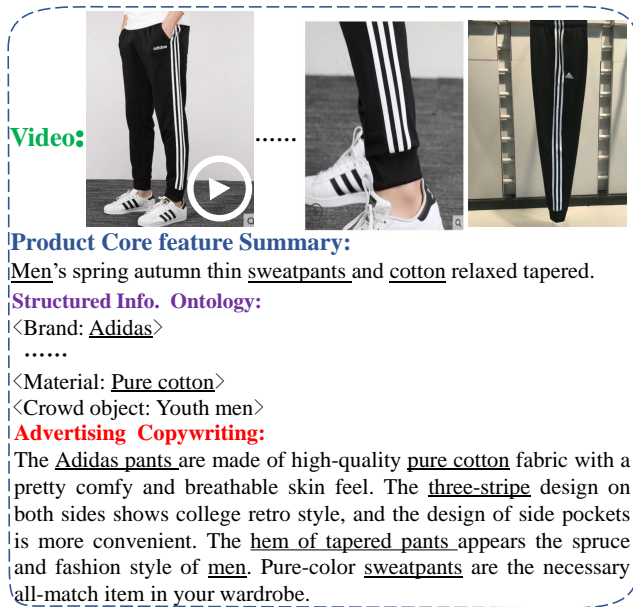


Figure 1: An illustration of our dataset. The four different parts of our dataset, from top to bottom are product information (commodity displaying video, the product core features summary, structured information) and commodity advertising description. We use multimodal multistructured information to assist in generating a semantically richer copywriting. The underlined words are closely related to the advertising copywriting, and are also important in terms of fidelity.

Multimodal product advertising text generation is different from video Caption [40, 41], natural language generation [6, 7], etc. The information sources needed for the commodity copywriting are diversified, which are closely related to the video display of commodity, commodity attribute table and commodity core feature summary. For example, as shown in Figure 1, the product video shows the spatial and visual color impact information of the product, while the structured information of the attribute word lists shows the accurate fine-grained information of the product, such as brand and the product core features summary (unstructured information), which summarizes the core features of the product in a general way. Through multimodal fusion, structured and unstructured information is combined to directly generate high-quality fine-grained text. This is also consistent with the fact that in real life, text generation of commodity advertisements is closely related to multi-modal information sources, and details of different modes

are considered comprehensively. Therefore, how to extract the required information from multimodal and fully integrate is one of the challenges in this task.

The other of the challenges is faithfulness. For commodity advertising copywriting, its fine-grained key information should greatly reduce factual errors of principle. For example, for an Adidas shoe, a model-dependent output branding error cannot occur because a similar Nike shoe appears in the training set. Faithfulness is an extremely critical core issue. This information from product structural attributes can make the description rigorous and reliable. However, the current metrics[5, 46], such as KOBE[5] are also based on N-gram for lexical diversity and BERTScore[46] are based on the knowledge pre-train, which cannot calculate the accuracy of attribute words in real-world advertisement application. To this end, we proposed a corresponding metric, hard homologous Metric, for auditing the fidelity of measured structured information.

To address these challenges, we elaborately collect a large-scale e-commercial multimodal multistructured advertising dataset for multimodal text generation research. To support in-depth research, we collect a rich set of product annotations. Our dataset consists of 120,984 product instances in both Chinese and English, in which each instance has a product video, a product core feature summary, structured information and a caption. In response to the realistic demand for advertising generation, we propose the multimodal information fusion module and generation decoder module which make full use of the rich information. In faithfulness, we propose Conceptualization Operations 4.1 to conceptualize complex and diverse information in real life as ontology. An ontology models generalized data, that is, we take into consideration general objects that have common properties and not specified individuals. Dataset and code will be available at our Website. The proposed network leads to a significant improvement over existing practical application methods, on our constructed dataset.

In summary, our contributions concentrate on the following aspects:

- (1) We introduce a fresh task: e-commercial multimodal advertising generation. A new large-scale high-quality and reliable e-commercial multimodal advertising dataset is introduced, which requires requires and supports multimodal fusion and faithfulness accuracy. It is also one of the largest video-text datasets in this field. E-MMAD is collected from human real life scenes and carefully selected so that it is qualified to meet the needs of real life.
- (2) We propose a simple yet effective strong baseline method to solve the challenges in reality. Our approach achieves the Top-1 accuracy in faithfulness and other metrics, outperforming existing baseline methods.
- (3) As for the fidelity of advertising copy, we propose the hard homologous metric. This metric allows advertising to be sourced.

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 [30, 31], cooking [8, 49], and Activities [2, 41]. MSR-VTT [41] is a widely-used dataset for video captioning, which has 10,000 videos from 257 activities and was collected in 2016. MSVD [4] was collected in 2011, containing

1970 videos. ActivityNet [1] has 20,000 videos but is used for Dense Video Captioning [15, 18], which means to describe multiple events in a video. TVR [21] is collected from movie clips whose text is mainly character dialogue. Vatec [40] is a famous dataset released in 2019, whose caption is written by batch manpower. Poet[45] is an e-commerce dataset containing two small raw datasets BFVD and FFVD. The data was downloaded directly from the Internet and not carefully filtered by human multimodal alignment. As shown in Table 1, we generated a larger video dataset after a lot of time and manpower screening. In addition, the struct info that we emphasize is carefully selected and generated by us to solve the real fidelity problem, not from the rough data of the network. This part will be used to solve the fidelity in the advertisement. Compared with some mainstream datasets in Table 1, our dataset also provide an additional product structured information. We find that the advertising caption includes a lot of structured information in fact.

2.2 Video Captioning Approaches

Video caption/description is one of the important tasks in multimodal text generation[15]. Early video caption methods are all based on templates [19, 25, 33]. However, sentences made in this way tend to be rigid. The sequence-to-sequence model [38] 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 and text will be generated in the decoding stage[17]. Some of the popular practices recently are based on data-driven [48] and transformer-based mechanisms [20, 43, 50]. MART [20] can produce more coherent, non-repetitive, and relevant text to enhance the transformer architecture by using memory storage units[27, 42]. Vx2text [23] uses multimodal inputs for text generation. They use a backbone [12, 34] model to transform different modalities information to natural language and then the problem turns to natural language generation. Recently, there are works[14, 39, 44] extracting object-level features in representing the videos for video caption. 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, product core feature summary and commodity advertising description. Different from previous works [4, 40, 41], 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. So the E-MMAD are limited to e-commerce aspect. In addition, we fully consider ethical privacy issues to ensure that the dataset has no potential negative effects and legal issues [11]. 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

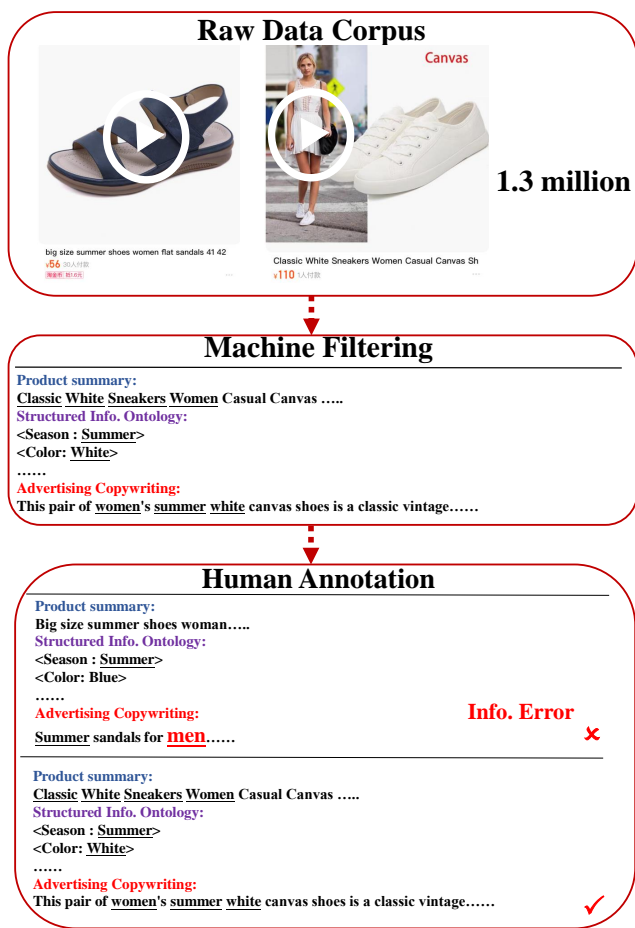


Figure 2: The process of creating a dataset, including machine filtering, manual post-filtering, etc. and data specification of the dataset. Each set of data needs to be carefully filtered and annotated manually in order to produce high-quality multimodal dataset.

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. Multiple copywriters tested and discussed to make the manual evaluation criterion several times. Finally, different testers sample 100 examples randomly according to the judgment rules of Figure 6, 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 6.

3) **Data annotation.** We invited 25 professional advertising copywriters as data screening and annotation staff to conduct manual screening under the rules of Figure 6 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 all people are consistently passable, can this data be qualified. Finally, 120,984 valid data has been generated. Simultaneously, we also translate the filtered 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 WMT 2019 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.

After 25 people's diligent work of manual data labeling and cleaning, there are 120,984 valid data selected finally.

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 Figure3. 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 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

Table 1: Comparison with other datasets. *Videos*, *Average Time*, *Caption Length*, *Classes* respectively represent the total number of videos in the dataset, the average video time in the dataset, the average length of the captions in the dataset and the number of instance types in the dataset. *Input Modality* indicates the input of the dataset, e.g. from Video to Text, Multimodal to Text. *Structure info.* means whether the dataset contains structured information. There are 3,876 keys of the structure information in E-MMAD dataset. en means English version dataset and zh means Chinese version dataset.

Datasets	#Videos	Average Time	Caption Length	#Classes	Input Modality
MSR-VTT [41]	10,000	14.8s	9	257	Video
MSVD [4]	1,970	9.0s	8	-	Video
TVR [21]	21,800	9.0s	13	-	Video-query
VaTEX (en/zh) [40]	41,269	10.0s	15/13	600	Video
FFVD (zh) [45]	32,763	27.7s	62	-	Video - Attribute
BFVD (zh) [45]	43,166	11.7s	93	-	Video - Attribute
E-MMAD (en/zh)	120,984	30.4s	97/67	4,863	Video - Summary - Structure info.

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 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. 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[40] 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 [10, 21], Product Search [3] 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 core feature summary, structured information(structured words) and the displaying video’s frames and fuse them to generate captions. Faced with various information words, our model uses ontology, a method of conceptualizing information. That is to pre-process the various data, conceptualize and extract information from the complex information words to Key as highly conceptual 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:

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

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

$$token_{gr} \rightarrow k_{gr},$$

$$\forall token_{gr} \in Values_{gr}. \quad (3)$$

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

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

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

$$rc_token \rightarrow v_{gr},$$

$$\forall rc_token \in Values_{rc}. \quad (6)$$

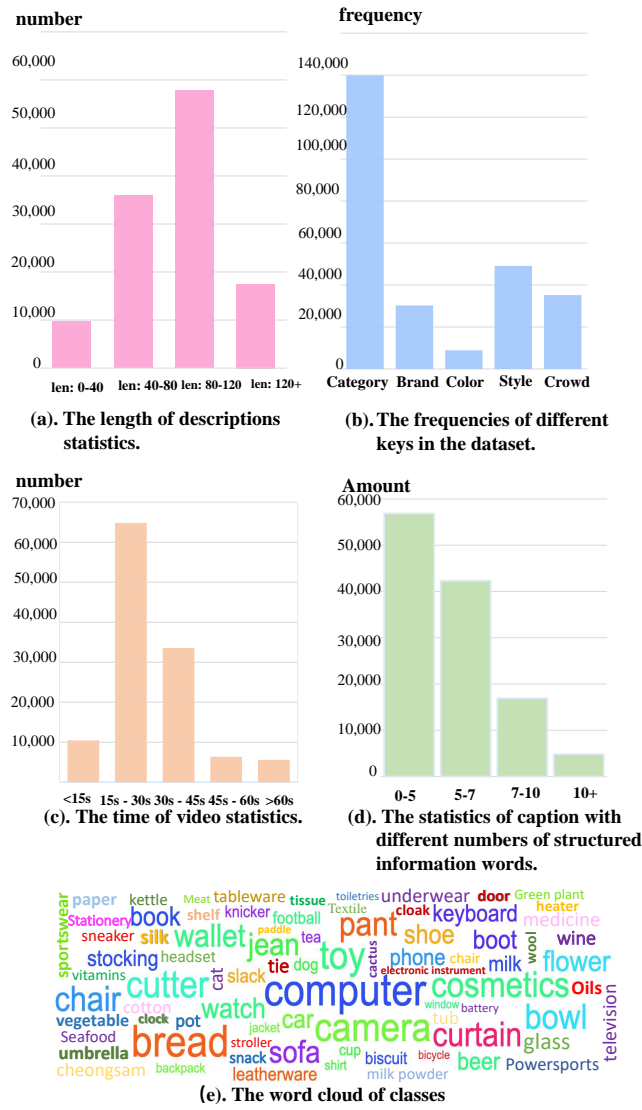


Figure 3: Statistics about the five different forms of data in our dataset. The data statistics are presented in terms of video, structured information, caption, and the main classes of the dataset contained, respectively.

Among them, Equation 3/6, $A \rightarrow B$ 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, the ontology of the structured information part is conceptual value words. An ontology models generalized data, that is, we take into consideration general objects that have common properties and not specified individuals. By this, the 3,876 types of Keys represent the various information words as the highly conceptual network feature input. We also

reference the summary as the basis to determine the priority position of each key according to the order in which the structured information appears in the summary.

4.2 Representation

Textual Information. Given a product core feature summary 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 [9, 47]. In addition, since the keys of structured information are prioritized, we use *position embedding* to represent the priority position of the keys.

Visual Information. Given a sequence of video frames/clips of length S , we feed it into pre-trained 3D CNNs[16] 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:

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

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

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

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

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 multimodal information consisting of a list of all $K + N + S$ modalities from $\{v_S^{\text{frames}}\}$, $\{v_K^{\text{words}}\}$ and $\{v_N^{\text{keys}}\}$. 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 [6, 28]. The decoder accesses multimodal fusion outputs at each layer with a multi-head attention [36]. Specifically, the decoder applies a multi-headed self-attention over the caption textual feature. After that, the position-wise feed forward layer was used to produce a distribution probability of each generation tokens for the final generated caption. There is a

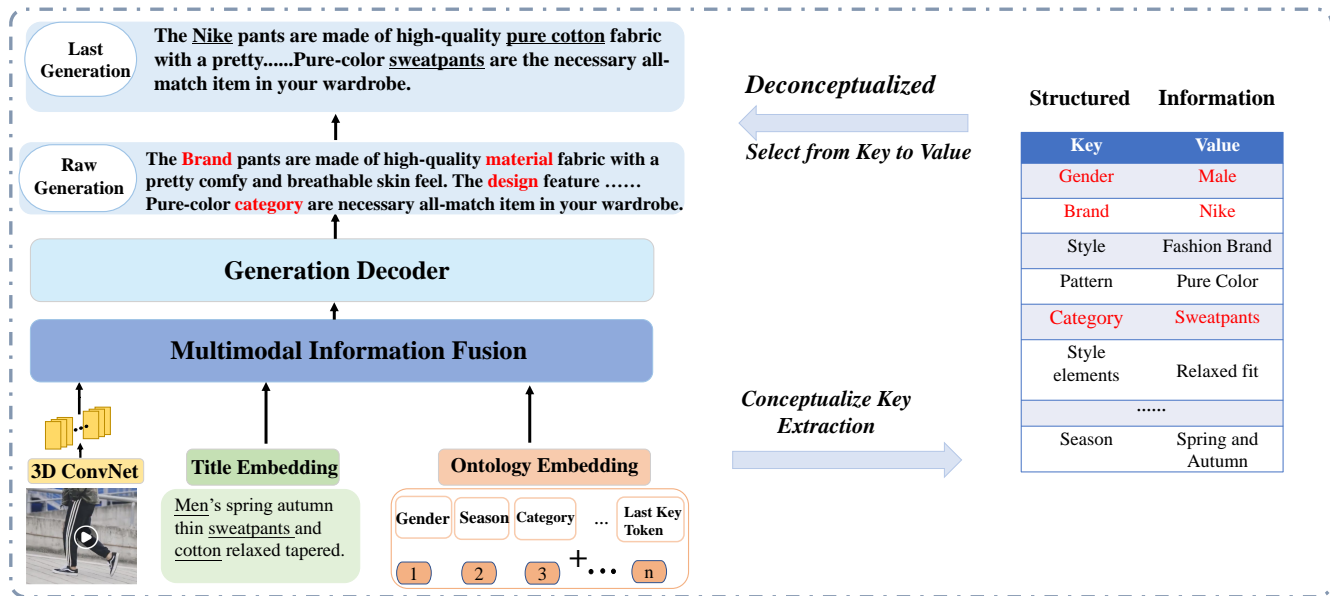


Figure 4: The overall architecture of our model, which contains three main parts: the representation for multimodal information, the multimodal fusion module based on self-attention and the generation decoder module on the basis of [29]. According to the Key-Value, the used Structure information words are conceptualized as ontology to face the various words such as assorted brands in real life.

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

Version	Input	Method	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
en	Text	NLG [6]	13.6	6.8	3.1	1.9	13.0	10.1
	Video	NACF [43]	18.9	7.9	3.9	2.2	15.3	14.8
	Multimodal	NACF + multi-input	20.0	8.5	4.3	2.4	17.8	18.5
		TVC [21]	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
zh-CN	Text	CPM (zh) [47]	7.9	4.6	1.1	0.5	7.2	8.3
	Multimodal	ours (zh)	11.6	6.5	4.4	2.2	12.5	15.3

description of part of the formula for the decoder module:

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

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

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

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

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

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 [36]. We refer to [29] as the model decoder 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 [43] for fairness and opt for the same

type, first extract 3D features with 2048 dimensions, 2048-D image features from ResNet-101 [13] pre-trained on ImageNet dataset. For generation decoder, we use `<sep>` to separate the input from the ground truth of caption. We adopt diverse automatic evaluation metrics to compare with other model: BLEU [26], Rouge-L [22], and CIDEr [37]. 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. For the sake of fairness of comparison, we simply modify the input part 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 substantially improve text generation tasks. It indicates that multimodal information indeed helps captioning by modal information between the mutual enhancement. 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 importance of our proposed dataset containing structured information. As shown in Table 3, we analyze the gap between the generated caption of the model and the real commodity advertising description in the absence of partial information. As we can see, the absence of any of the three input components significantly degrades the final generated caption result. From our analysis of the generated caption, we can conclude that: 1) the lack of structured information will make the generated caption less informative, rigorous and reliable.

2) The lack of a commodity core feature summary or displaying video will impair the foundation of generated text. In addition, the structured information is like a knowledge base, which can promote inference and judgment to generate appropriate caption.

Conceptual Operation. Considering that writing product descriptions in real life often involves a great number of unfamiliar words, which makes it hard for the model to identify and remember its feature when facing a new word, such as new brand name. The predecessor's approach tend to use as much corpus and large model

Table 3: Performance comparison with our proposed model by masking different parts of input and only using the remainder as input. Here "Summary", "SI" and "Video" indicates commodity core feature summary, structured information and commodity displaying video respectively.

Input	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
SI & Video	22.8	14.8	6.9	5.5	22.2	25.3
Summary & 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
Summary & SI	22.0	13.8	5.8	4.9	20.6	23.7

parameters as possible, which brings huge difficulties to natural language generation. In this case, we proposed the Conceptualization operation. As shown in Table 4, we conduct ablation experiments about Conceptualization on the Chinese and English datasets. As for models without conceptual operations, we use unconceptualized captions as the ground truth to train. We directly input unordered structured words for the input of the model. Experiments have proved that the Conceptualization operation can indeed bring a significant effect improvement, because this method can conceptualize and extract information from complex information in the dataset, and thus highly conceptualize network features. We expect this discovery to inspire the community.

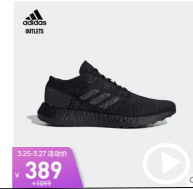
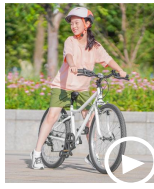
Table 4: Performance comparison of whether our proposed model has conceptual operations (CO).

Operation	Bleu1	Bleu2	Bleu3	Bleu4	Rouge_L	CIDEr
ours w/o CO (en)	23.8	15.4	8.1	6.4	24.2	27.3
ours w/o CO (zh)	9.9	5.5	2.8	1.5	10.1	12.4
ours w/ CO (en)	25.0	16.6	9.6	7.2	25.3	29.1
ours w/ CO (zh)	11.6	6.5	4.4	2.2	12.5	15.3

5.4 Hard Homologous Metric

As shown in Table 5.2, we adopt three common NLG metrics. BLEU [26] for sanity check, ROUGE_L [22] based on the longest common subsequence co-occurrence, CIDEr [37] based on human-like consensus. However, in reality application, as shown in Figure 5, we found that, the advertisement copywriting has the characteristics of flexible language style and rich vocabulary, and fidelity is particularly important. In terms of faithfulness, product advertising especially focuses on key information such as product brand and color, which is not fully reflected in the above indicators. For example, for Adidas shoes, the model can be lazy and output Nike brand because there is a similar training corpus in the training dataset, which is a common and serious mistake. To address this problem, we proposed a hard homologous metric. In terms of accuracy, we traverse the attribute words in each groundtruth (dataset has marked the key ontology of each value word), and compare with generation to calculate the proportion of correct words in generated words. In terms of error rate, according to figure 3 statistical data and realistic requirements, we successively select brand, color, material, people, time and season as the five key labels as top-5 core attribute words. If they are inconsistent, they will be regarded as errors. In the meantime, we'll call the rest unknown. Under such

Video:



	CK chain bag for women.	Decathlon children's 20 inch bicycle.	Adidas black sports shoes.
Summary:	<Material: Leather> <Color: White> <Crowd: Women> <Season: Summer>	<Brand: Decathlon> <Color: White> <Crowd: Children> <Wheel size: 20 inches>	<Brand: Adidas> <Color: Dark> <Season: Summer> <Category: Sports shoes>
Structured Info. Ontology :			
Copywriter:	The white women's bag is the latest fashion trend published by Calvin Klein. It adopts one shoulder chain design to make shopping more convenient for you. Essential small bag in summer!	Decathlon kids bikes use 20-inch wheels, which is designed to be safer for kids. The white bike body highlights children's youth and vitality in the sun. It is the best choice for children's outdoor trips.	The Adidas sports shoe is designed with knitting technology, which is more comfortable to wear. Limited summer running time event price is only 389 RMB, now or never, come to snap up!
Ours :	CK white women's bag is designed in a one shoulder chain style. It is exquisite and portable. Summer shopping essential small bag, show goddess demeanor!	Decathlon children's bikes use 20-inch wheels. The child riding a white bicycle looks handsome outside the house. It is essential children's toys for families!	The Adidas sports shoe is a pure black simple design style. Knitting technology design conforms to human body mechanics. Essential sneakers for summer running!

Figure 5: Some example results generated by our methods.

hard homologous metric, significant statistical analysis results were shown in Table 5 for the faithfulness of product advertisements copywriting.

Table 5: The results of our propose hard homologous metric.

		Correct Rate	Erro Rate	Unknown
Ours	w/ conception	18.7%	20.2%	61.1%
	w/o conception	13.8%	27%	59.2%
NLG (CPM)	w/ conception	15.6%	23.1%	62.3%
	w/o conception	9.8%	30%	60.2%
Video caption (NACF)	w/ concaption	10.9%	29%	60.1%
	w/o conception	5%	38%	57%

5.5 Human Assessment

It is well-known that the human evaluation metrics[35] for video captioning are required due to the inaccurate evaluation by automatic metrics. We especially focus on advertising generation, which depend on human aesthetics. So we invite the people involved in the data annotation and new advertising slogan designers to conduct the human evaluation. We select 200 samples from the test dataset and each evaluator evaluate each of these 200 samples to reflect the performance of our model by rating whether the caption generated by our model can be used as a description of the product. As the result shows in Table 6, the caption generated by our model has a certain degree of pass rating, whose results can be approbated by people. Therefore, this is also acceptable that our experiments

on Table 2 did not achieve high scores for mechanical evaluation indicators. We also test the human performance. The human test results were generated by the merchant copywriters.

Table 6: The pass rate results of the human evaluation, reflecting the proportion of the 200 reality application test examples where the model generated caption could be used as a product description that describes the reasonableness of the generated caption. Annotators are from the dataset annotation and persons are from the frequent online shopping masses.

	Annotator 1	Annotator 2	Annotator 3	Person 1	Person 2	Person 3
Ours	42%	44%	43%	48%	56%	53%
CPM	30%	23%	27%	40%	47%	39%
Human Performance	74%	77%	69%	90%	81%	89%

6 CONCLUSION AND FUTURE WORK

This research sets out to provide an e-commercial multimodal multistructured advertisement copywriting dataset, E-MMAD, which is one of the largest video captioning datasets in this field. Based on E-MMAD, we also present a novel task: e-commercial multimodal multistructured advertising generation, and propose a baseline method on the strength of multistructured information reasoning to solve the realistic demand. We hope the release of our E-MMAD would facilitate the development of multimodal generation problems. However, there still exist limitations about our dataset and

method that should be acknowledged as shown in Figure 5. We cannot identify the price information of the video in (c), which may require video OCR or ASR technology. Moving forward, we are planning to extend E-MMAD to better performance and more diversified tasks by exploring new model structures, fine-grained and so on.

REFERENCES

- [1] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Nibbles. 2015. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 961–970.
- [2] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. 2018. A short note about kinetics-600. *arXiv preprint arXiv:1808.01340* (2018).
- [3] Wei-Cheng Chang, Daniel Jiang, Hsiang-Fu Yu, Choon Hui Teo, Jiong Zhang, Kai Zhong, Kedarnath Kolluri, Qie Hu, Nikhil Shandilya, Vyacheslav Levgraf, 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*. 2643–2651.
- [4] 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*. 190–200.
- [5] Qibin Chen, Junyang Lin, Yichang Zhang, Hongxia Yang, and Jie Tang. 2019. Towards Knowledge-Based Personalized Product Description Generation in E-commerce. *ACM* (2019).
- [6] 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* (2019).
- [7] Zewen Chi, Li Dong, Furu Wei, Wenhui Wang, Xian-Ling Mao, and Heyan Huang. 2020. Cross-lingual natural language generation via pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 7570–7577.
- [8] 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*. 2634–2641.
- [9] 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* (2018).
- [10] Maksim Dzabrac, Maksim Kalashnikov, Stepan Komkov, and Aleksandr Petushko. 2021. Mdmmt: Multidomain multimodal transformer for video retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3354–3363.
- [11] 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* (2018).
- [12] 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*. 12046–12055.
- [13] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh. 2017. Learning Spatio-Temporal Features with 3D Residual Networks for Action Recognition. *arXiv preprint arXiv:1708.07632* (2017).
- [14] Yaosi Hu, Zhenzhong Chen, Zheng-Jun Zha, and Feng Wu. 2019. Hierarchical global-local temporal modeling for video captioning. In *Proceedings of the 27th ACM International Conference on Multimedia*. 774–783.
- [15] Vladimir Iashin and Esa Rahtu. 2020. Multi-modal dense video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 958–959.
- [16] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 2012. 3D convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence* 35, 1, 221–231.
- [17] Xu Jia, Efstratios Gavves, Basura Fernando, and Tinne Tuytelaars. 2015. Guiding the long-short term memory model for image caption generation. In *Proceedings of the IEEE international conference on computer vision*. 2407–2415.
- [18] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Nibbles. 2017. Dense-captioning events in videos. In *Proceedings of the IEEE international conference on computer vision*. 706–715.
- [19] Niveda Krishnamoorthy, Girish Malkamkar, 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*.
- [20] Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L Berg, and Mohit Bansal. 2020. Mart: Memory-augmented recurrent transformer for coherent video paragraph captioning. *arXiv preprint arXiv:2005.05402* (2020).
- [21] Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2020. 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*. Springer, 447–463.
- [22] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- [23] 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*. 7005–7015.
- [24] Sidi Lu, Yaoming Zhu, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Neural text generation: Past, present and beyond. *arXiv preprint arXiv:1803.07133* (2018).
- [25] 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*. 747–756.
- [26] 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*. 311–318.
- [27] Wenjie Pei, Jiyuan Zhang, Xiangrong Wang, Lei Ke, Xiaoyong Shen, and Yu-Wing Tai. 2019. Memory-attended recurrent network for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8347–8356.
- [28] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. (2018).
- [29] 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 (2019), 9.
- [30] 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*. 3202–3212.
- [31] Anna Rohrbach, Atousa Torabi, Marcus Rohrbach, Niket Tandon, Christopher Pal, Hugo Larochelle, Aaron Courville, and Bernt Schiele. 2017. Movie description. *International Journal of Computer Vision* 123, 1 (2017), 94–120.
- [32] Linfeng Song, Ante Wang, Jinsong Su, Yue Zhang, Kun Xu, Yubin Ge, and Dong Yu. 2021. Structural information preserving for graph-to-text generation. *arXiv preprint arXiv:2102.06749* (2021).
- [33] Xiaoou Tang, Xinbo Gao, Jianzhuang Liu, and Hongjiang Zhang. 2002. A spatial-temporal approach for video caption detection and recognition. *IEEE transactions on neural networks* 13, 4 (2002), 961–971.
- [34] 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*. 6450–6459.
- [35] Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Krahmer. 2020. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Computer Speech & Language* (2020), 101151.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [37] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4566–4575.
- [38] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to sequence-video to text. In *Proceedings of the IEEE international conference on computer vision*. 4534–4542.
- [39] Huiyun Wang, Youjiang Xu, and Yahong Han. 2018. Spotting and aggregating salient regions for video captioning. In *Proceedings of the 26th ACM international conference on Multimedia*. 1519–1526.
- [40] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. 2019. Vatec: A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 4581–4591.
- [41] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5288–5296.
- [42] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*. PMLR, 2048–2057.
- [43] Bang Yang, Yuexian Zou, Fenglin Liu, and Can Zhang. 2019. Non-Autoregressive Coarse-to-Fine Video Captioning. *arXiv preprint arXiv:1911.12018* (2019).
- [44] Ziwei Yang, Yahong Han, and Zheng Wang. 2017. Catching the temporal regions-of-interest for video captioning. In *Proceedings of the 25th ACM international conference on Multimedia*. 146–153.

1045	[45]	Shengyu Zhang, Ziqi Tan, Jin Yu, Zhou Zhao, Kun Kuang, Jie Liu, Jingren Zhou, Hongxia Yang, and Fei Wu. 2020. Poet: Product-oriented video captioner for E-commerce. In <i>Proceedings of the 28th ACM International Conference on Multimedia</i> . 1292–1301.	1103
1046			1104
1047			1105
1048	[46]	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> (2019).	1106
1049			1107
1050	[47]	Zhengyan Zhang, Xu Han, Hao Zhou, Pei Ke, Yuxian Gu, Deming Ye, Yujia Qin, Yusheng Su, Haozhe Ji, Jian Guan, et al. 2021. CPM: A large-scale generative Chinese pre-trained language model. <i>AI Open</i> 2 (2021), 93–99.	1108
1051			1109
1052	[48]	Ziqi Zhang, Zhongang Qi, Chunfeng Yuan, Ying Shan, Bing Li, Ying Deng, and Weiming Hu. 2021. Open-book Video Captioning with Retrieve-Copy-Generate Network. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> . 9837–9846.	1110
1053			1111
1054			1112
1055	[49]	Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018. Towards automatic learning of procedures from web instructional videos. In <i>Thirty-Second AAAI Conference on Artificial Intelligence</i> .	1113
1056			1114
1057	[50]	Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. 2018. End-to-end dense video captioning with masked transformer. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> . 8739–8748.	1115
1058			1116
1059			1117
1060			1118
1061			1119
1062			1120
1063			1121
1064			1122
1065			1123
1066			1124
1067			1125
1068			1126
1069			1127
1070			1128
1071			1129
1072			1130
1073			1131
1074			1132
1075			1133
1076			1134
1077			1135
1078			1136
1079			1137
1080			1138
1081			1139
1082			1140
1083			1141
1084			1142
1085			1143
1086			1144
1087			1145
1088			1146
1089			1147
1090			1148
1091			1149
1092			1150
1093			1151
1094			1152
1095			1153
1096			1154
1097			1155
1098			1156
1099			1157
1100			1158
1101			1159
1102			1160

A DATASET SUPPLEMENT.

There are the data standards examples as shown in Figure 6. More data will be released on the Source Link: <https://github.com/E-MMAD/E-MMAD>

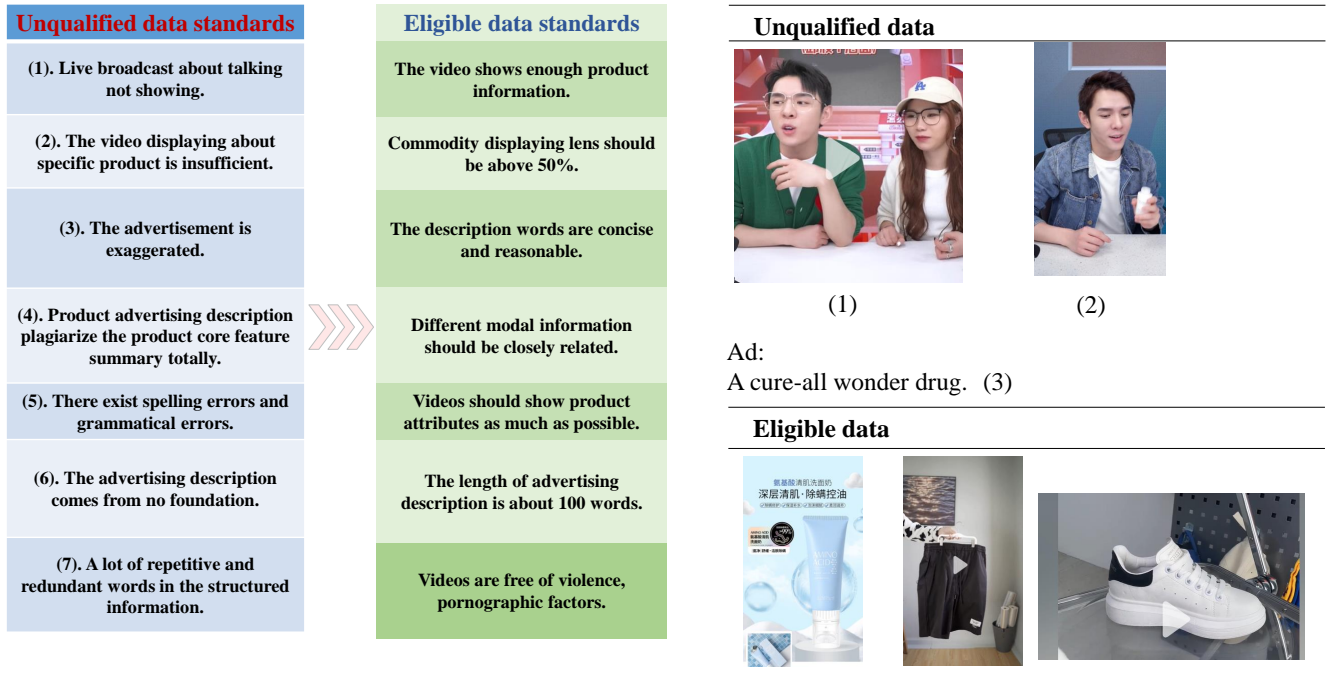


Figure 6: The overall architecture of our model, which contains three main parts: the representation for multimodal information, the multimodal fusion module based on self-attention and the generation decoder module on the basis of [29]. According to the Key-Value, the used Structure information words are conceptualized as ontology to face the various words such as assorted brands in real life.