

# Multi-modal Representation Learning for Video Advertisement Content Structuring

Daya Guo\*

Sun Yat-sen University  
Guangzhou, Guangdong, China  
guody5@mail2.sysu.edu.cn

Zhaoyang Zeng\*

Sun Yat-sen University  
Guangzhou, Guangdong, China  
zengzhy5@mail2.sysu.edu.cn

## ABSTRACT

Video advertisement content structuring aims to segment a given video advertisement and label each segment on various dimensions, such as presentation form, scene, and style. Different from real-life videos, video advertisements contain sufficient and useful multi-modal content like caption and speech, which provides crucial video semantics and would enhance the structuring process. In this paper, we propose a multi-modal encoder to learn multi-modal representation from video advertisements by interacting between video-audio and text. Based on multi-modal representation, we then apply Boundary-Matching Network to generate temporal proposals. To make the proposals more accurate, we refine generated proposals by scene-guided alignment and re-ranking. Finally, we incorporate proposal located embeddings into the introduced multi-modal encoder to capture temporal relationships between local features of each proposal and global features of the whole video for classification. Experimental results show that our method achieves significantly improvement compared with several baselines and Rank 1 on the task of Multi-modal Ads Video Understanding in ACM Multimedia 2021 Grand Challenge. Ablation study further shows that leveraging multi-modal content like caption and speech in video advertisements significantly improve the performance.

## CCS CONCEPTS

• **Computing methodologies** → **Scene understanding; Video summarization.**

## KEYWORDS

Multi-Modal Representation; Transformer; Inception; Boundary-Matching Network; Proposal Located Classifier.

### ACM Reference Format:

Daya Guo and Zhaoyang Zeng\*. 2021. Multi-modal Representation Learning for Video Advertisement Content Structuring. In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21)*, October 20–24, 2021, Virtual Event, China. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3474085.3479218>

\*Equal contribution. Order determined by alphabet order.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '21, October 20–24, 2021, Virtual Event, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3479218>

## 1 INTRODUCTION

As the number of video advertisements in Internet grows rapidly, video ads content analysis methods have become more crucial and attracted more attention from both academia and industry. Video ads content structuring is an important task in video ads content analysis area, which aims to segment a given video ads in time and label each segment on various dimensions, such as presentation form, scene, and style. As shown in Figure 1, video advertisements are different from real-life videos in temporal action detection datasets like ActivityNet [2] and THUMOS [6]. They provide sufficient multi-modal content like caption and speech for the purpose of promoting and popularizing their products. Therefore, multi-modal contents like caption and speech provide crucial video semantic and are important for the structuring process [4].

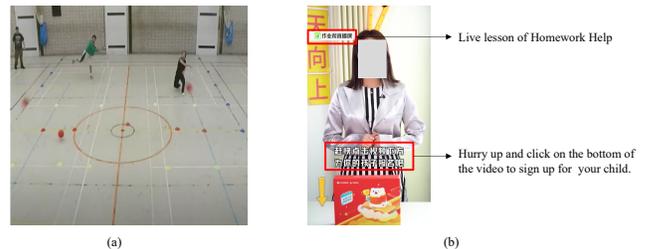


Figure 1: (a) The left part is a frame of video from ActivityNet dataset. (b) The right part is a frame of video advertisement from ACM Multimedia 2021 Grand Challenge.

In this paper, we propose a multi-modal encoder to learn multi-modal representation of video ads. The encoder consists of three components, including a video-audio encoder, a text encoder and a cross-modality encoder. The video-audio encoder contains several Inception modules [13] implemented by 1D convolutional layers, which takes video and audio features as the input and outputs context representation of video-audio. The text encoder is a powerful pre-trained model BERT [3]. A Transformer based cross-modality encoder is used for cross-modality interaction between the text and the video-audio to obtain multi-modal representation, which will be used for video advertisement content structuring.

Inspired by state-of-the-art methods [7–9] in temporal action detection area, we decouple the task of video advertisement content structuring into two subtasks, i.e Temporal Segmentation and Proposal Tagging. In the temporal segmentation phrase, we apply Boundary-Matching Network (BMN) [8] to generate temporal proposals using multi-modal representation. We then propose to refine them by scene-guided alignment and re-ranking to make

the generated proposals more accurate. In the proposal tagging phrase, we take a whole video as the input and incorporate proposal located embeddings into the multi-modal encoder to capture temporal relationships between local features of each proposal and global features of the whole video. Finally, the mean pooling is adopted on the top of multi-modal representation for tagging.

We evaluate the proposed method on the dataset of Multi-modal Ads Video Understanding in ACM Multimedia 2021 Grand Challenge. Experiments show that our method achieves state-of-the-art performance and Rank 1 on the leaderboard. Further analysis shows that multi-modal information and newly introduced proposal located embeddings are helpful for video ads content structuring.

## 2 APPROACH

### 2.1 Overview

Figure 2 gives an overview of our approach. As shown in the Figure, we decouple the task of video ads content structuring into two subtasks, i.e. segment and tagging. We first propose a multi-modal encoder as the backbone of two subtasks that takes text with video-audio feature as the input to obtain multi-modal representation of video advertisement. In the segment phase, we adopt Boundary-Matching Network (BMN) [8] to generate proposals based on the multi-modal representation. These generated proposals are further refined by scene-guided alignment [1] and re-ranking, which will make proposals more accurate. In the tagging phase, we propose a proposal located embeddings (PLE) to capture temporal relationships between local features of each proposal and global features of the whole video for proposal tagging. In the next, we will introduce how to encode multi-modal content including video-audio and text in Section 2.2. The details about our solution on temporal segmentation and proposal tagging subtasks will be introduced in Sections 2.3 and Section 2.4, respectively.

### 2.2 Multi-modal Encoder

**Text Encoder.** Captions in video advertisements are of significant for video ads content structuring. Taking Figure 1 as an example, the caption ‘‘Hurry up and click on the bottom of the video to sign up for your child’’ could help the model infer this frame is a promotion page. To leverage text information in video ads, we first utilize Optical Character Recognition (OCR) technique to extract caption as the text input, denoted as  $X = \{x_0, x_1, \dots, x_{n-1}\}$ . Since pre-trained models [3, 12] have led to strong improvement on numerous natural language processing (NLP) tasks, we use a powerful pre-trained model BERT [3] as our text encoder to encode the text input and obtain hidden states of the text  $H_X = \{h_{x_0}, h_{x_1}, \dots, h_{x_{n-1}}\}$ .

**Video-Audio Encoder.** Inspired by [10], we combine video and audio features to obtain more discriminate representations. Given an input video, we first split it into video clips with length 0.5 second. For each video clip, we follow [15] to use S3D model pre-trained on HowTo100M dataset[11] to extract its visual feature. For the given audio input, we extract its feature by VGGish[5]. We then re-sample the video and audio feature sequence into the same temporal size  $m$  using bi-linear interpolation. We denote the re-sampled video and audio features as  $V = \{v_0, v_1, \dots, v_{m-1}\}$  and  $A = \{a_0, a_1, \dots, a_{m-1}\}$ , respectively. Since different video segment are variance in length,

we follow [13] to adopt Inception module to capture the information from different temporal sizes. Specifically, We concatenate  $V$  and  $A$  into  $Y = \{[v_0; a_0], \dots, [v_{m-1}; a_{m-1}]\}$  along the channel dimension to form the input video-audio features, then feed it into two Inception modules and produce  $H_Y = \{h_{y_0}, h_{y_1}, \dots, h_{y_{m-1}}\}$ . We follow [13] to design the Inception module, while only replacing the 2D convolutional layers into 1D ones.

**Cross-Modality Encoder.** The text encoder and video-audio encoder mainly focus on a part of modality. To fully leverage the text and video-audio, we adopt a 6-layer Transformer [14] as our cross-modality encoder for cross-modality interaction. The input of cross-modality encoder is constructed by summing type embeddings and the concatenation  $[H_X; H_Y]$  of the text and video-audio representation. Finally, we obtain multi-modal representation  $H' = \{h'_{x_0}, \dots, h'_{x_{n-1}}, h'_{v_0}, \dots, h'_{v_{m-1}}\}$ . We denote the multi-modal representation of the whole video ads as  $H = \{h'_{v_0}, h'_{v_1}, \dots, h'_{v_{m-1}}\}$ , which will be used for segment and tagging phases.

### 2.3 Temporal Segmentation

**Boundary-Matching Network.** In the task of temporal action detection, Lin et al. [8] propose the Boundary-Matching Network (BMN) to generate high-quality temporal proposals. The network consists of two components, including Temporal Evaluation Module (TEM) and Proposal Evaluation Module (PEM). TEM aims to predict precise boundaries for all temporal locations in untrimmed video and PEM aims to provide confidence for each proposal.

Different from datasets like ActivityNet [2] and THUMOS [6], there are no backgrounds in video ads content structuring. Therefore, TEM only predicts boundary probability of two segments for all temporal locations, shown in the top left of segment part in Figure 2. We use 3-layer CNN following by a 6-layer Transformer as the TEM and take multi-modal representation  $H$  as the input to calculate boundary probability of each video clip  $p^{TEM} = \{p_0^{TEM}, \dots, p_{m-1}^{TEM}\}$ . We use the cross entropy (CE) loss to train TEM, where  $y_i \in \{0, 1\}$  is the boundary label for  $i$ -th video clip.

$$loss_{TEM} = -\frac{1}{m} \sum_{i=0}^{m-1} [y_i \log p_i^{TEM} + (1 - y_i) \log(1 - p_i^{TEM})] \quad (1)$$

For proposal evaluation module shown in the top right of segment part in Figure 2, we use the same network architecture and loss function  $loss_{PEM}$  as Lin et al. [8] and provide a confidence  $p_{ij}^{conf}$  for each proposal  $prop_{ij}$  from  $i$ -th to  $j$ -th video clip. The final loss function of Boundary-Matching Network is empirically as:

$$loss_{BMN} = 5 \cdot loss_{TEM} + loss_{PEM} \quad (2)$$

In the inference phase of BMN, for each proposal  $prop_{ij}$ , we fuse its boundary probabilities and confidence scores by multiplication to generate the final confidence score  $s_{ij}$ :

$$p_{ij}^{prop} = p_i \cdot p_j \cdot p_{ij}^{conf} \cdot \min\{\overline{p_{i+1}}, \overline{p_{i+2}}, \dots, \overline{p_{j-1}}\} \quad (3)$$

where  $\overline{p_k} = 1 - p_k$ . Different from the final confidence score of Lin et al. [8], we add the last term to indicate the approximate probability of no existing boundary in the middle of the proposal  $prop_{ij}$ . The main reason for using approximate probability is that true probability will over-punish proposals that too long.

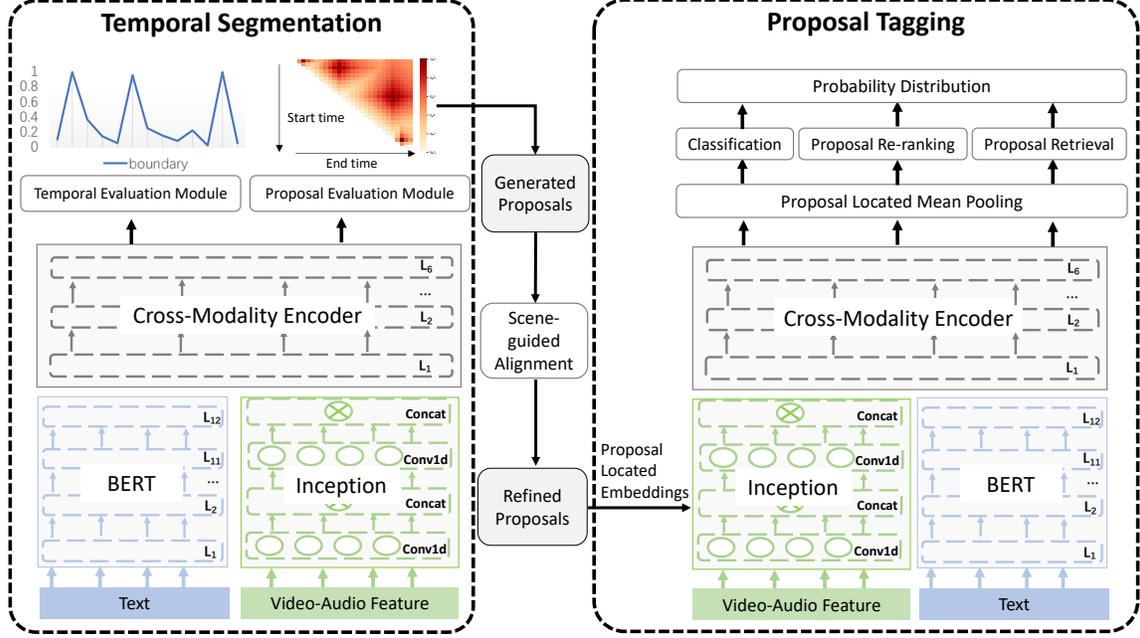


Figure 2: Overview of our proposed framework. We decouple the task of video ads content structuring into two subtasks. The first subtask is segment that generates proposals, which is shown in the left part. The second subtask is tagging that classifies each proposal, which is shown in the right part.

Finally, we use Non-Maximum Suppression (NMS) algorithm to generate non-overlap proposal by the confidence score  $s_{ij}$ .

**Scene-Guided Alignment.** The duration of each video clip we use for temporal segmentation is only 0.5 second, thus the temporal boundaries predicted by BMN are rough. To make the predicted boundaries more precise, we propose to use the scene-changed frames to fine-adjust the predicted boundaries. We call this method scene-guided alignment (SGA). We follow [1] to extract all frames with scene probability greater than 0.1 for each given video. For each predicted boundary, if the temporal error between it and its nearest scene frame is less than 0.5 second, we will move the boundary to the position of its nearest scene frame.

## 2.4 Proposal Tagging

**Proposal Located Embedding.** To incorporate proposal temporal information for proposal tagging, we introduce proposal located embedding before the Inception modules. Specially, given a proposal  $prop_{ij}$ , the input of Inception module changes from  $Y$  to  $Y^t = \{[v_0; a_0] + t_0, \dots, [v_{m-1}; a_{m-1}] + t_{m-1}\}$ , where  $t_k$  is a trainable randomly initialized embedding to indicate proposal location if  $i \leq k \leq j$  otherwise  $t_k$  is another trainable randomly initialized embedding to indicate non-proposal location. After obtaining multi-modal representation of video-audio from multi-modal encoder, we get the final vector  $v_{ij}$  by mean pooling over the location of the proposal  $prop_{ij}$ . Finally, we leverage the final vector for proposal tagging in a multi-task manner.

**Classification.** The proposal tagging task can be formulate as a multi-label classification problem. We use a fully connected layer

with sigmoid activation for classification, and adopt binary cross entropy as loss function. The probabilities of classification for the proposal  $prop_{ij}$  are denoted as  $p_{ij}^{cls}$ .

**Proposal Re-ranking.** To make proposal more accurate, we introduce a re-ranking task in the classifier. The task aims to re-rank generated proposals, which will improve the precision of generated proposals. Specially, we use a fully connected layers with sigmoid activation to predict Intersection over Union (IoU) score for each generated proposal  $prop_{ij}$ , denote as  $p_{ij}^{iou}$ .

**Proposal Retrieval.** For video ads content structuring task, it would have more application value if it can be salable to some new scenes or categories. Inspired by Yang et al. [16], we propose to utilize multi-modal representation  $v_{ij}$  of the proposal  $prop_{ij}$  to retrieve top 10 most similar segments using cosine similarity. We denote labels and similarity score of these retrieved segments as  $\{g_{ij}^0, g_{ij}^1, \dots, g_{ij}^9\}$  and  $\{c_{ij}^0, c_{ij}^1, \dots, c_{ij}^9\}$ , respectively. The retrieved result is calculated by weighted summing retrieved labels:

$$p_{ij}^{ret} = \frac{\sum_{k=0}^9 c_{ij}^k g_{ij}^k}{\sum_{k=0}^9 c_{ij}^k} \quad (4)$$

In the inference phase, the probability distribution  $p_{ij}^{cat}$  of categories of the proposal  $prop_{ij}$  is calculated as:

$$p_{ij}^{cat} = p_{ij}^* \cdot p_{ij}^{iou} \cdot p_{ij}^{prop}, \quad (5)$$

where  $p_{ij}^*$  can be  $p_{ij}^{cls}$ ,  $p_{ij}^{ret}$ , or the combination of them.

### 3 EXPERIMENTS

#### 3.1 Experiments Setup

**Dataset.** We evaluate our proposed approach on the dataset of Multi-modal Ads Video Understanding in ACM Multimedia 2021 Grand Challenge. The dataset consists of 5,000 videos. Each video is split into one or several clips by annotators. Each video clip is annotated by at least one categories. The total number of category is 82. When performing evaluation, the task constraints that predicted proposals can not have overlap with each other. To avoid redundancy prediction, the task also limits that a proposal only produce 20 category labels.

**Training Detail.** All our experiments are conducted on one NVIDIA Tesla V100-32G GPU. We use AdamW optimizer with 1e-4 learning rate to train all models for 10 epochs and evaluate the model using 5-fold cross-validation.

**Evaluation.** The goal of temporal segmentation is to generate high quality proposals to cover ground truth segments with high recall and high temporal overlap. To evaluate proposal quality, we follow Lin et al. [8] to use AUC under IoU thresholds [0.5 : 0.05 : 0.95] as a metric. Beyond that, we also use F1-score between predicted and ground truth boundaries as another metric to evaluate the precision of proposals. Given a prediction of a video, if a predicted boundary can match any ground truth boundaries within 0.5s error, it will be consider as a true positive prediction, and otherwise will be consider as a false positive prediction. Note that one ground truth boundary will be only matched once. Finally, the overall performance of generated proposals is the product of AUC and F1-score. To evaluate the performance of proposal tagging, we follow Caba Heilbron et al. [2] to use  $mAP@[0.5:0.05:0.95]$  as the metric.

#### 3.2 Evaluation on Temporal Segmentation

Model	Video	Audio	Text	AUC	F1	Overall
BMN [8]	✓			72.1	78.7	56.7
Ours	✓			74.8	79.0	59.1
	✓	✓		<b>75.1</b>	78.6	59.0
	✓	✓	✓	74.4	<b>80.9</b>	<b>60.2</b>

Table 1: Results on video temporal segmentation task.

We take the state-of-the-art model in temporal action detection as our baseline, i.e. BNN [8]. The difference between our segment model and the baseline includes: (1) We add the last term to indicate the approximate probability of no existing boundary in the middle of the proposal in Equation 3. (2) we leverage multi-modal content.

Table 1 show experiment results on temporal segmentation task. We can see that our approach achieve a 3.5 gain of overall score, which significantly outperforms the baseline. From Table 1, we find that incorporating the audio feature improves the AUC score by 0.3 but will hurt the F1-score, which shows that audio feature may not have much effect on temporal segmentation task. After leveraging text information, we can see that the interaction with text bring 1.1 gain of overall score compared with the model that only leverages video, which reveals the importance of multi-modal representation.

Model	Video	Audio	Text	mAP
Inception [13] w/o PLE	✓			27.3
Ours	✓			27.9
	✓	✓		28.1
	✓	✓	✓	<b>29.5</b>

Table 2: Results on proposal tagging task.

#### 3.3 Evaluation on Proposal Tagging

We report the experiment results on proposal tagging task in Table 2. In this experiment, we take Inception [13] without proposal located embedding (**Inception w/o PLE**) as a baseline, which only use video feature and remove PLE in classification module. For fair comparison, all settings in Table 2 use the same generated proposals generated by best model in Table 1.

When only leveraging the same video feature in Table 2, we can see that incorporating PLE into the Inception bring 0.6 gain of mAP score, which demonstrates that the proposal located embedding could help proposal tagging. After leveraging multi-modal content like speech and caption, results show that our multi-modal encoder significantly outperforms the model that only uses single-modal, which shows the effectiveness of our multi-modal representation.

#### 3.4 Classification-Vs. Retrieval-based Classifier

Classifier	mAP
Classification-based method	29.8
Retrieval-based method	30.3
Ensemble method	<b>31.7</b>

Table 3: Results of various classifiers.

In real application scenario, the category number may be updated frequently. Retrieval-based method that retrieves similar examples from training dataset and infers probability distribution of categories of the proposal may bring more application value. In the Table 3, we show the performance of various classifiers. The probability distribution of categories of the proposal  $prop_{ij}$  in classification-based method is calculated as  $p_{ij}^{cls}$  described in Section 2.4, while the probability distribution  $p_{ij}^{ret}$  in retrieval-based method is calculated by retrieving similar proposal as Equation 4.

We can find that, the retrieval-based method can achieve comparable result with the classification-based method. To achieve higher evaluation score, we ensemble the results from both classification-based and retrieval-based methods, and find that the ensemble model brings further improvements with 1.4% absolute gain.

### 4 CONCLUSION

In this paper, we propose a multi-modal encoder to learn multi-modal representation from video advertisements by interacting between video-audio and text. Experiments show that multi-modal representation significantly improve temporal segmentation and proposal tagging tasks. Based on multi-modal representation, we present an efficient framework for the task of video ads content structuring. The framework achieves Rank 1 on the task of Multi-modal Ads Video Understanding in ACM Multimedia 2021 Grand Challenge. In future work, we would like to explore how to pre-train a powerful multi-modal encoder using video ads for empowering video ads content analysis.

## REFERENCES

- [1] 2017. PySceneDetect. <https://github.com/Breakthrough/PySceneDetect>.
- [2] 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*. 961–970.
- [3] 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).
- [4] Daya Guo, Jiangshui Hong, Binli Luo, Qirui Yan, and Zhangming Niu. 2019. Multi-modal representation learning for short video understanding and recommendation. In *Proceedings of the IEEE International Conference on Multimedia & Expo Workshops*. IEEE, 687–690.
- [5] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. 2017. CNN architectures for large-scale audio classification. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 131–135.
- [6] Haroon Idrees, Amir R Zamir, Yu-Gang Jiang, Alex Gorban, Ivan Laptev, Rahul Sukthankar, and Mubarak Shah. 2017. The THUMOS challenge on action recognition for videos “in the wild”. *Computer Vision and Image Understanding* 155 (2017), 1–23.
- [7] Chuming Lin, Jian Li, Yabiao Wang, Ying Tai, Donghao Luo, Zhipeng Cui, Chengjie Wang, Jilin Li, Feiyue Huang, and Rongrong Ji. 2020. Fast learning of temporal action proposal via dense boundary generator. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11499–11506.
- [8] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. 2019. Bmn: Boundary-matching network for temporal action proposal generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3889–3898.
- [9] Tianwei Lin, Xu Zhao, Haisheng Su, Chongjing Wang, and Ming Yang. 2018. Bsn: Boundary sensitive network for temporal action proposal generation. In *Proceedings of the European Conference on Computer Vision*. 3–19.
- [10] Shuang Ma, Zhaoyang Zeng, Daniel McDuff, and Yale Song. 2020. Active Contrastive Learning of Audio-Visual Video Representations. In *Proceedings of the International Conference on Learning Representations*.
- [11] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. 2019. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2630–2640.
- [12] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. *arXiv preprint arXiv:1802.05365* (2018).
- [13] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1–9.
- [14] 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 *Proceedings of the Advances in Neural Information Processing Systems*. 5998–6008.
- [15] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. 2018. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European Conference on Computer Vision*. 305–321.
- [16] Jianfei Yang, Zhaoyang Zeng, Kai Wang, Han Zou, and Lihua Xie. 2021. GarbageNet: A Unified Learning Framework for Robust Garbage Classification. *IEEE Transactions on Artificial Intelligence* (2021).