# VISUAL ANSWER LOCALIZATION WITH CROSS-MODAL MUTUAL KNOWLEDGE TRANSFER

*Yixuan Weng\*†, Bin Li\*\*‡*

† National Laboratory of Pattern Recognition Institute of Automation, Chinese Academy Sciences
‡ College of Electrical and Information Engineering, Hunan University

## ABSTRACT

The goal of visual answering localization (VAL) in the video is to obtain a relevant and concise time clip from a video as the answer to the given natural language question. Early methods are based on the interaction modelling between video and text to predict the visual answer by the visual predictor. Later, using the textual predictor with subtitles for the VAL proves to be more precise. However, these existing methods still have cross-modal knowledge deviations from visual frames or textual subtitles. In this paper, we propose a cross-modal mutual knowledge transfer span localization (MutualSL) method to reduce the knowledge deviation. MutualSL has both visual predictor and textual predictor, where we expect the prediction results of these both to be consistent, so as to promote semantic knowledge understanding between cross-modalities. On this basis, we design a one-way dynamic loss function to dynamically adjust the proportion of knowledge transfer. We have conducted extensive experiments on three public datasets for evaluation. The experimental results show that our method outperforms other competitive state-of-the-art (SOTA) methods, demonstrating its effectiveness[1].

***Index Terms*—** Cross-modal, Mutual Knowledge Transfer, Visual Answer Localization

## 1. INTRODUCTION

The explosion of online videos has changed the way that people obtain information, and knowledge [1, 2]. Various video platforms make it more convenient for people to perform video queries [3, 4]. However, people who want to get direct instructions or tutorials from the video often need to browse the video content several times to locate relevant parts, which usually takes time and effort [5]. Visual answer localization (VAL) is an emerging technology to solve the above problem [6], and has received wide attention because of its practical value [7, 8]. As shown in Fig. 1(a), the task of VAL is to find a time clip that can answer the given question.

---

[1] All the experimental datasets and codes are open-sourced on the website https://github.com/WENGSYX/MutualSL.
\*: These authors contributed equally to this work.
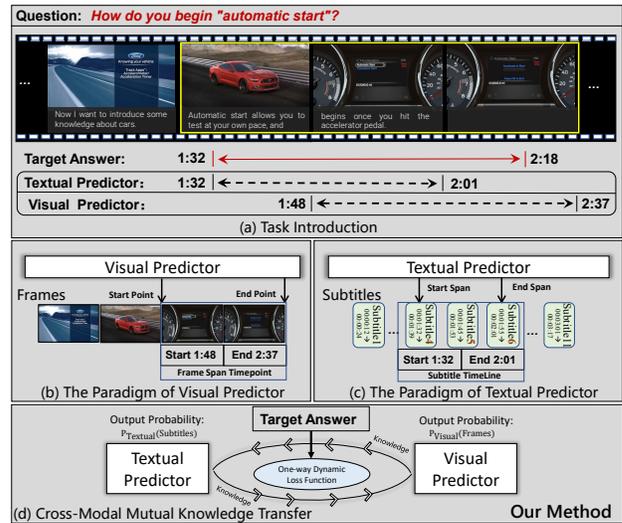⋆: Corresponding author.



**Fig. 1**. Task description of the visual answer localization, where the below is the paradigms of the previous methods and our method.

For example, when inputting "How do you begin 'automatic start'", you may need to find a clip according to voice content (or transcribed text subtitles) and visual frames. The VAL technology can not only recognize the relevant video clips to the text questions but also return the target visual answer (1:32 ˜ 2:18).

The existing VAL method can be mainly divided into visual predictor and textual predictor according to the prediction contents. The paradigm of visual predictor is shown in Fig. 1(b). The video information is first extracted according to the frame features, and then these frame features queried by the question are used to predict the relevant time points [9, 10]. The paradigm of textual predictor is shown in Fig. 1(c). The textual predictor adopts a span-based method to model the cross-modal information, where the predicted span intervals with subtitle timeline are used as the final results [8, 11].

The performance of the textual predictor is better than the visual one [7], because it uses the additional subtitle information, and embeds visual information into the text feature space with the visual information as an auxiliary feature. However, as shown in Fig. 1(a), results from both two predictors suffer cross-modal knowledge deviations. For the textual pre-
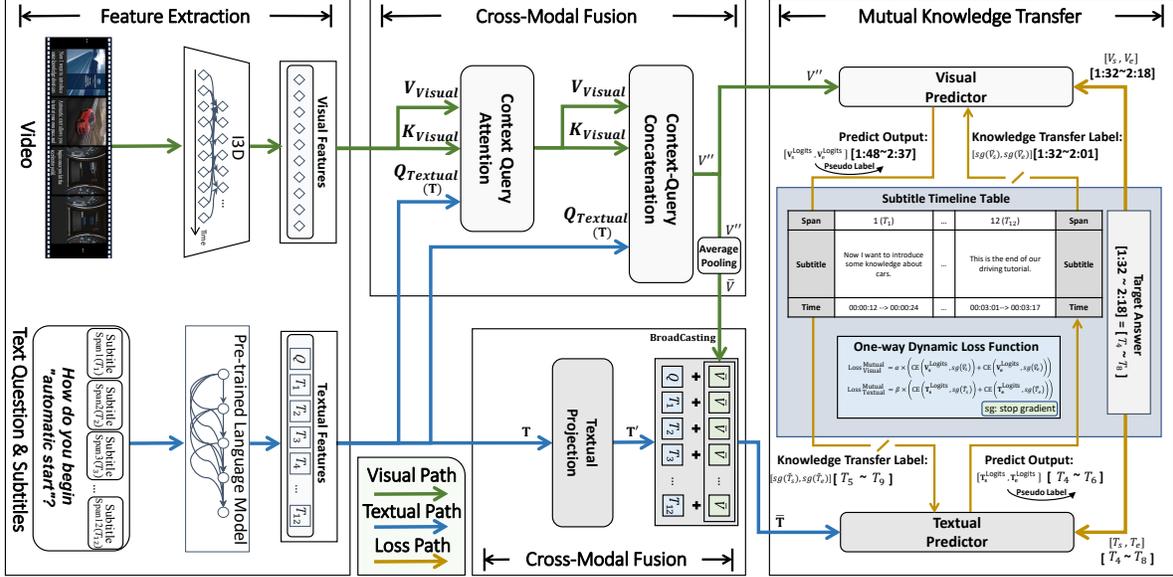
**Fig. 2**. Overview of the proposed cross-modal mutual knowledge transfer span localization (MutualSL).

dictor, if the video lacks subtitle information for a long clip, this clip cannot be located; For the visual predictor, it is difficult to have continuous clip prediction because of the frequent changing of video scenes and semantics to the question.

In this paper, we propose a novel cross-modal mutual knowledge transfer span localization (MutualSL) method to reduce the cross-modal knowledge deviation shown in Fig. 1(d). Specifically, the MutualSL uses both visual predictor and textual predictor, where these two predictors have different prediction targets so that they have different strength perceptions of different-modal information. We expect that these two predictors can enhance the information perception of their own modal. Each predictor needs to predict the output value of another predictor on the basis of the target answer in the training stage. Then we design a one-way dynamic loss function (ODL) to dynamically adjust the knowledge transfer, which can alleviate the difference of cross-modal knowledge transferring in the training process.

Our contributions are as follows: (1) we propose the MutualSL method, which for the first time uses two different predictors in VAL tasks meanwhile, and uses a Look-up Table to achieve cross-modal knowledge transfer; (2) We design ODL to dynamically adjust the knowledge transfer, which can alleviate the differences in knowledge transfer between different predictors; (3) We have conducted extensive experiments to prove the effectiveness of the MutualSL, where results show that the proposed method outperforms all other competitive SOTA methods in VAL tasks.

## 2. METHOD

### 2.1. Task Definition

Given an untrimmed video $V$ with a duration of $k$ seconds, the corresponding subtitle $S = \{T_i\}_{i=1}^r$ and the text question is $Q$, the VAL task requires us to predict the most relevant visual

clips within the video $[V_s^*, V_e^*] \subseteq V$ that answer the question $Q$, where $T_i$ is the subtitle of each span, $r$ is the subtitle span length, and $[V_s, V_e]$ is defined as the target time clip answer, $s, e \in [1, k]$. Moreover, it provides a subtitle timeline table, which is translated a span into corresponding timeline span from each subtitle set $S$. We can use the subtitle timeline table as the Look-up Table, such as providing accurate target answers for textual predictor transferred from the frame span timepoints, and vice versa.

$$[V_s^*, V_e^*] = \underset{V_s, V_e}{\mathrm{Argmin}}(P([V_s, V_e]|V, S, Q)) \quad (1)$$

### 2.2. Main Structure

As shown in Fig. 2, the MutualSL is divided into three parts, which are Feature Extraction, Cross-modal Fusion, and Mutual Knowledge Transfer. The Mutual Knowledge Transfer includes visual predictor and textual predictor.

**Feature Extraction.** Following the previous method [9, 11], we use the pre-trained visual model I3D [12] and the pre-trained language model (PLM) [13] to extract feature vectors from video $V$ and concatenated texts $T = [Q, T_1, \ldots, T_r]$ respectively. These pre-trained models can provide us with high-quality information representation.

$$\mathbf{V} = \mathrm{I3D}(V), \mathbf{T} = \mathrm{PLM}(T), \quad (2)$$

where $\mathbf{V} \in \mathbb{R}^{k \times d}$ and $\mathbf{T} \in \mathbb{R}^{n \times d}$, the $d$ is the dimension and $n$ is the length of the concatenated text tokens of T.

**Cross-modal Fusion.** We use context query attention (CQA) [9] to capture the cross-modal interaction between visual and textual to enhance the semantics in the visual path. CQA adopts two attention mechanism context to query ($\mathcal{D}$) and query to context ($\mathcal{F}$) processes for cross-modal modeling, where $\mathcal{G}_r \in \mathbb{R}^{k \times n}$ and $\mathcal{G}_c \in \mathbb{R}^{k \times n}$ represent the row- and column-wise normalization of $\mathcal{G}$ by SoftMax.

$$\mathcal{D} = \mathcal{G}_r \cdot \mathbf{T} \in \mathbb{R}^{k \times d}, \mathcal{F} = \mathcal{G}_c \cdot \mathcal{G}_r^T \cdot \mathbf{V} \in \mathbb{R}^{k \times d}$$

We use one layer of feedforward neural network (FFN$_\text{C}$) and convolution layer (in_channels = $2d$, out_channels = $d$) as Context-Query Concatenation to capture deeper semantic information, where $\{\mathbf{V'}, \mathbf{V''}\} \in \mathbb{R}^{k \times d}$.

$$\mathbf{V'} = \text{FFN}_\text{C}\big([\mathbf{V}; \mathcal{D}; \mathbf{V} \odot \mathcal{D}; \mathbf{V} \odot \mathcal{F}]\big) \quad (3)$$

$$\mathbf{V''} = \text{Conv1d}\big(\text{Concat}[\text{Attention}(\mathbf{V'}, \mathbf{T}); \mathbf{T}]\big) \quad (4)$$

We use a textual projection layer (FFN$_\text{P}$) to extract text features $\mathbf{T'} \in \mathbb{R}^{n \times d}$ in the text path. Then we embed the average pooled visual feature $\overline{\mathbf{V}}$ into each token $T_j$ in $\mathbf{T'}$ through the broadcast mechanism, where $\overline{\mathbf{V}} \in \mathbb{R}^d$ and $T_j \in \mathbb{R}^{1 \times d}$,

$$\mathbf{T'} = \text{FFN}_\text{P}(\mathbf{T}), \overline{\mathbf{V}} = \text{AvgPool}(\mathbf{V''})) \quad (5)$$

$$\overline{\mathbf{T}} = \{\overline{\mathbf{V}} + T_j\}_{j=1}^n, \overline{\mathbf{T}} \in \mathbb{R}^{n \times d} \quad (6)$$

**Visual Predictor.** We use two unidirectional LSTMs, including LSTM$_\text{Start}$ and LSTM$_\text{End}$, where the in_channels = $k \times d$, out_channels = $k \times d$. Two feedforward layers FFN$_\text{Start}^\text{Visual}$ and FFN$_\text{End}^\text{Visual}$( in_channels = $k \times d$, out_channels = $k$) are adopted to construct a visual span predictor. We input the features $\mathbf{V''}$ into LSTMs, then use the feedforward layer to calculate the predicted time point logits, including the start time point and end time point.

$$\mathbf{V_s^{Logits}} = \text{FFN}_\text{Start}^\text{Visual}(\text{LSTM}_\text{Start}(\mathbf{V''})) \quad (7)$$

$$\mathbf{V_e^{Logits}} = \text{FFN}_\text{End}^\text{Visual}(\text{LSTM}_\text{End}(\mathbf{V''})) \quad (8)$$

**Textual Predictor.** We follow to the structure of QANet [14] and calculate the probability of outputting the start and the end subtitle point through two different feedforward layers FFN$_\text{Start}^\text{Textual}$ and FFN$_\text{End}^\text{Textual}$, where the in_channels = $n \times d$, out_channels = $n$.

$$\mathbf{T_s^{Logits}} = \text{FFN}_\text{Start}^\text{Textual}(\overline{\mathbf{T}}), \mathbf{T_e^{Logits}} = \text{FFN}_\text{End}^\text{Textual}(\overline{\mathbf{T}}) \quad (9)$$

**Loss Function.** We adopt Cross-Entropy (CE) function to maximize the visual predictor logits of the target span-point $[V_s, V_e]$. Also, we convert $[V_s, V_e]$ to $[T_s, T_e]$ by subtitle timeline Look-up Table for the loss calculation.

$$\text{Loss}_\text{Visual} = \text{CE}(\mathbf{V_s^{Logits}}, V_s) + \text{CE}(\mathbf{V_e^{Logits}}, V_e) \quad (10)$$

$$\text{Loss}_\text{Textual} = \text{CE}(\mathbf{T_s^{Logits}}, T_s) + \text{CE}(\mathbf{T_e^{Logits}}, T_e) \quad (11)$$

### 2.3. Look-up Table

The outputs of different predictors are inconsistent shown in Fig. 1. In order to solve the problem of semantic information deviation between cross-modalities, we design a Look-up Table $\mathbb{Q}$ to convert the output probabilities of one predictor as the target answer of another (such as converting the prediction subtitle timelines $T_{(s/e)}$ of the textual predictor to the corresponding frame span timepoints $V_{(s/e)}$, which realized information alignment of the cross-modal target answer.

$$\breve{T}_s = \text{Argmin}\,(V_s - \mathbb{Q}(T_i)), \breve{T}_e = \text{Argmin}\,(V_e - \mathbb{Q}(T_i)) \quad (12)$$

$$\breve{V}_s = \text{Argmin}\,(T_s - \mathbb{Q}(V_i)), \breve{V}_e = \text{Argmin}\,(T_e - \mathbb{Q}(V_i)) \quad (13)$$

### 2.4. Mutual Knowledge Transfer

In order to perform cross-modal mutual knowledge transfer, we introduce auxiliary objectives, and expect that the predictor can effectively learn the cross-modal information by predicting the output produced by the another-modal predictor. The whole procession is shown as follows, where the notation $\mathbb{E}$ means the Expectation. This represents that we want to transfer the knowledge as the pseudo label from one predictor to another.

$$\mathbb{E}([\mathbf{T_s^{Logits}}, \mathbf{T_e^{Logits}}]) = [\breve{T}_s, \breve{T}_e],\ \mathbb{E}([\mathbf{V_s^{Logits}}, \mathbf{V_e^{Logits}}]) = [\breve{V}_s, \breve{V}_e] \quad (14)$$

### 2.5. One-way Dynamic Loss Function

In the training stage, the cross-modal knowledge reserves for each predictor are different. We design a one-way dynamic loss function (ODL) that can adjust knowledge transferring. On the right of Fig. 2, ODL can dynamically adjust the proportion of knowledge transferring by comparing the matching between the prediction result and the target answer via the IoU function shown in equation (15).

$$\text{IOU}(A, B) = \frac{A \cap B}{A \cup B} \quad (15)$$

Meanwhile, the knowledge difference between the textual predictor and the visual predictor will lead to inconsistent learning progress. Therefore, we use stop gradient (sg) for the two predictors to learn independently (this means one-way).

$$\text{Loss}_\text{Visual}^\text{Mutual} = \alpha \times (\text{CE}(\mathbf{V_s^{Logits}}, sg(\breve{V}_s)) + \text{CE}(\mathbf{V_e^{Logits}}, sg(\breve{V}_e))) \quad (16)$$

$$\text{Loss}_\text{Textual}^\text{Mutual} = \beta \times (\text{CE}(\mathbf{T_s^{Logits}}, sg(\breve{T}_s)) + \text{CE}(\mathbf{T_e^{Logits}}, sg(\breve{T}_e))) \quad (17)$$

where the $\alpha$ and $\beta$ can be calculated dynamically, which are presented as follows.

$$\alpha = IOU([\breve{V}_s, \breve{V}_e], [V_s, V_e]), \beta = IOU([\breve{T}_s, \breve{T}_e], [T_s, T_e]) \quad (18)$$

Finally, our loss function is:

$$\text{Loss} = \text{Loss}_\text{Visual} + \text{Loss}_\text{Textual} + \text{Loss}_\text{Visual}^\text{Mutual} + \text{Loss}_\text{Textual}^\text{Mutual} \quad (19)$$

## 3. EXPERIMENT

### 3.1. Experimental Setting

We evaluate MutualSL in three different public VAL datasets, where these datasets are formed with the text questions and corresponding visual answer clips as the target answers. The MedVidQA [6] is a medical instructional dataset that contains 3,010 question-and-answer (QA) pairs and 899 videos; TutorialVQA [5] contains 76 tutorial videos about software editing tutorials with 6,195 QA pairs; The VehicleVQA [17] dataset has a series of *How-To* videos that introduce practical instructions on vehicles, including 9,482 QA pairs within 107 videos. Following previous works [6, 18, 19], we use

**Table 1**. Performance on three public datasets compared with several SOTA methods.

| Method | MedVidQA | | | | TutorialVQA | | | | VehicleVQA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU |
| VSLBase [9] | 27.66 | 14.19 | 6.99 | 21.01 | 10.84 | 9.58 | 0.37 | 8.71 | 18.95 | 8.64 | 4.28 | 20.11 |
| TMLGA [15] | 23.87 | 14.84 | 6.21 | 20.49 | - | - | - | - | - | - | - | - |
| VSLNet [9] | 30.32 | 16.61 | 8.39 | 22.41 | 9.96 | 9.21 | 0.00 | 8.58 | 16.53 | 8.47 | 4.03 | 20.07 |
| ACRM [10] | 24.83 | 16.55 | 10.96 | 22.89 | 12.61 | 5.17 | 1.26 | 11.18 | 20.77 | 12.10 | 8.27 | 22.28 |
| RaNet [16] | 32.90 | 20.64 | 15.48 | 27.48 | - | - | - | - | - | - | - | - |
| MoR [8] | 47.10 | 22.74 | 10.97 | 30.67 | - | - | - | - | - | - | - | - |
| VPTSL [11] | 77.42 | **61.94** | **44.52** | 57.81 | 50.07 | 40.01 | 25.79 | 40.20 | 74.15 | 67.15 | 54.59 | 64.51 |
| MutualSL | **80.65** | **61.94** | 39.99 | **58.32** | **60.14** | **43.59** | **28.28** | **43.48** | **78.74** | **69.81** | **53.14** | **65.74** |

**Table 2**. We report the effect of whether to conduct cross-modal mutual knowledge transfer for Visual Predictor (VP) and Textual Predictor (TP) respectively. Both predictors are output from MutualSL, but their VAL performance is different. Therefore, we report the results on whether using Mutual Knowledge Transfer (MKT).

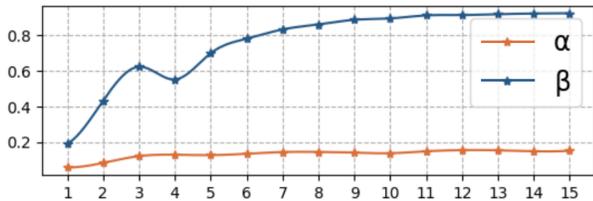| Method | MedVidQA | | | | TutorialVQA | | | | VehicleVQA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU | IoU=0.3 | IoU=0.5 | IoU=0.7 | mIoU |
| Ours (VP) W/O MKT | 18.63 | 11.53 | 8.12 | 16.42 | 12.02 | 6.31 | 4.57 | 12.66 | **24.40** | **9.48** | 2.62 | 18.53 |
| Ours (VP) | **28.24** | **14.68** | **9.51** | **21.45** | **12.36** | **6.47** | **4.92** | **13.48** | 16.53 | 8.87 | **4.53** | **20.19** |
| Ours (TP) W/O MKT | 78.06 | 61.29 | **43.87** | 57.78 | 56.00 | 38.62 | 23.45 | 40.44 | 71.74 | 65.46 | 49.52 | 62.22 |
| Ours (TP) | **80.65** | **61.94** | 39.99 | **58.32** | **60.14** | **43.59** | **28.28** | **43.48** | **78.74** | **69.81** | **53.14** | **65.74** |



**Fig. 3**. The changing trend of $\alpha$ and $\beta$ during training.

IoU-0.3/0.5/0.7 and mIoU as the evaluation metrics to compare several state-of-the-art (SOTA) methods on VAL tasks. We use the same visual extractor and text extractor in each baseline to ensure fairness, and follow the original author's parameter settings. We compare with several SOTA methods, VSLBase/VSLNet [9], TMLGA [15], ACRM [10] and RaNet [16] use visual predictor to predict frame span timepoint. MoR [8] and VPTSL [11] use a textual predictor to predict textual subtitle span, which are the competitive SOTA methods. In the parameter settings of MutualSL, we set $d = 1024$ and use the AdamW optimizer [20], where lr = 1e-5. We use Pytorch in three A100 GPUs for experiments, where the batch size = 4 and training epoch = 15. For all experiments, we repeat three-time experiments to reduce the random errors.

### 3.2. Results

As shown in Table 1, we compare the performance of different methods in three datasets of the VAL task. The MutualSL achieves SOTA performance in most metrics, which shows the effectiveness of our method, especially the mIoU increases by 0.51, 3.28, and 1.23 respectively. The reason may be that mutual knowledge transfer (MKT) can guide the predictor to understand different information, thus alleviating the deviations of knowledge in different modalities.

To further analyze the impact of MKT on the ODL under different predictors, we perform the ablation study of hyperparameters $\alpha$ and $\beta$ shown in Fig. 3. We can clearly see that both $\alpha$ and $\beta$ increase with the increase of the epoch. The $\beta$ is more than $\alpha$, because the textual predictor has better answering localization ability, which is also in line with the experimental results shown in Table 1. The Fig. 3 also shows that the ODL can dynamically adapt the ability of knowledge transfer from different predictors.

We've conducted extensive ablation experiments to analyze the MKT in Table 2. In the visual predictor, using MKT can improve the mIoU indicators of the three datasets by 0.69 on average; The average increase in the textual predictor is 2.33 mIoU. These prove that the use of MKT can enhance the model's perception of different modal information, thus improving the performance of VAL. We find that although the improvement of the visual predictor is low, the MKT can greatly improve the performance of the textual predictor. Meanwhile, the performance of the textual predictor outperforms the visual predictor. Therefore, we use the result of the textual predictor as the output of MusicalSL in the prediction phase.

## 4. CONCLUSION

In this paper, we proposed a cross-modal mutual knowledge transfer method (MutualSL) for VAL tasks. This method alleviates the problem of knowledge deviation, which uses visual predictor and textual predictor for cross-modal mutual knowledge transfer. We compare and ablate the proposed methods in three public datasets of the VAL task, where the proposed method outperforms all competitive SOTA methods. This proves the effectiveness of the MutualSL. In the future, we hope to explore more methods such as knowledge distillation for understanding cross-modal knowledge to promote the development of related fields.

# 5. REFERENCES

[1] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah, "Transformers in vision: A survey," *ACM Computing Surveys*, 2022.

[2] Yifan Xu, Huapeng Wei, Minxuan Lin, Yingying Deng, Kekai Sheng, Mengdan Zhang, Fan Tang, Weiming Dong, Feiyue Huang, and Changsheng Xu, "Transformers in computational visual media: A survey," *Computational Visual Media*, vol. 8, no. 1, pp. 33–62, 2022.

[3] Mahnaz Koupaee and William Yang Wang, "Wikihow: A large scale text summarization dataset," *arXiv preprint arXiv:1810.09305*, 2018.

[4] Maria Törhönen, Max Sjöblom, Lobna Hassan, and Juho Hamari, "Fame and fortune, or just fun? a study on why people create content on video platforms," *Internet Research*, vol. 30, no. 1, pp. 165–190, 2019.

[5] Anthony Colas, Seokhwan Kim, Franck Dernoncourt, Siddhesh Gupte, Daisy Zhe Wang, and Doo Soon Kim, "Tutorialvqa: Question answering dataset for tutorial videos," *language resources and evaluation*, 2019.

[6] Deepak Gupta, Kush Attal, and Dina Demner-Fushman, "A dataset for medical instructional video classification and question answering," *arXiv preprint arXiv:2201.12888*, 2022.

[7] Deepak Gupta and Dina Demner-Fushman, "Overview of the MedVidQA 2022 shared task on medical video question-answering," in *Proceedings of the 21st Workshop on Biomedical Language Processing*. Association for Computational Linguistics.

[8] Wojciech Kusa, Georgios Peikos, Óscar Espitia, Allan Hanbury, and Gabriella Pasi, "DoSSIER at MedVidQA 2022: Text-based approaches to medical video answer localization problem," in *Proceedings of the 21st Workshop on Biomedical Language Processing*, Dublin, Ireland, May 2022, Association for Computational Linguistics.

[9] Hao Zhang, Aixin Sun, Wei Jing, and Joey Tianyi Zhou, "Span-based localizing network for natural language video localization," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 6543–6554.

[10] Haoyu Tang, Jihua Zhu, Meng Liu, Zan Gao, and Zhiyong Cheng, "Frame-wise cross-modal matching for video moment retrieval," *IEEE Transactions on Multimedia*, vol. 24, pp. 1338–1349, 2021.

[11] Bin Li, Yixuan Weng, Bin Sun, and Shutao Li, "Towards visual-prompt temporal answering grounding in medical instructional video," *arXiv preprint arXiv:2203.06667*.

[12] Joao Carreira and Andrew Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.

[13] Pengcheng He, Jianfeng Gao, and Weizhu Chen, "Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing," *arXiv preprint arXiv:2111.09543*, 2021.

[14] Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V. Le, "Qanet: Combining local convolution with global self-attention for reading comprehension," *international conference on learning representations*, 2018.

[15] Cristian Rodriguez-Opazo, Edison Marrese-Taylor, Fatemeh Sadat Saleh, Hongdong Li, and Stephen Gould, "Proposal-free temporal moment localization of a natural-language query in video using guided attention," *workshop on applications of computer vision*, 2020.

[16] Jialin Gao, Xin Sun, Mengmeng Xu, Xi Zhou, and Bernard Ghanem, "Relation-aware video reading comprehension for temporal language grounding," *empirical methods in natural language processing*, 2021.

[17] Hongyin Luo, Mitra Mohtarami, James Glass, Karthik Krishnamurthy, and Brigitte Frances Mora Richardson, "Integrating video retrieval and moment detection in a unified corpus for video question answering.," *conference of the international speech communication association*, 2019.

[18] Meng Liu, Xiang Wang, Liqiang Nie, Xiangnan He, Baoquan Chen, and Tat-Seng Chua, "Attentive moment retrieval in videos," *international acm sigir conference on research and development in information retrieval*.

[19] Yitian Yuan, Tao Mei, and Wenwu Zhu, "To find where you talk: Temporal sentence localization in video with attention based location regression," *national conference on artificial intelligence*, 2019.

[20] Ilya Loshchilov and Frank Hutter, "Decoupled weight decay regularization," *Learning*, 2017.