

# GPTSee: Enhancing Moment Retrieval and Highlight Detection via Description-Based Similarity Features

Yunzhuo Sun, Yifang Xu, Zien Xie, Yukun Shu, and Sidan Du, *Member, IEEE*

**Abstract**—Moment retrieval (MR) and highlight detection (HD) aim to identify relevant moments and highlights in video from corresponding natural language query. Large language models (LLMs) have demonstrated proficiency in various computer vision tasks. However, existing methods for MR&HD have not yet been integrated with LLMs. In this letter, we propose a novel two-stage model that takes the output of LLMs as the input to the second-stage transformer encoder-decoder. First, MiniGPT-4 is employed to generate the detailed description of the video frame and rewrite the query statement, fed into the encoder as new features. Then, semantic similarity is computed between the generated description and the rewritten queries. Finally, continuous high-similarity video frames are converted into span anchors, serving as prior position information for the decoder. Experiments demonstrate that our approach achieves a state-of-the-art result, and by using only span anchors and similarity scores as outputs, positioning accuracy outperforms traditional methods, like Moment-DETR.

**Index Terms**—Image description, semantic similarity, video moment retrieval, video highlight detection.

## I. INTRODUCTION

As the internet and video production technology evolve rapidly, users upload hundreds of millions of videos to various platforms daily. How to effectively search and browse through such a vast amount of content has attracted widespread attention. Given a video and a natural language query, video moment retrieval (MR) [1], [2], [3] strives to retrieve the most relevant spans, each comprising a start and an end moment. On the other hand, video highlight detection (HD) [4], [5], [?] aims to predict moment-wise highlight scores across the whole video. In this letter, we focus on MR&HD simultaneously due to their shared characteristics, notably the need to learn the similarity between the textual query and video moments.

With the recent surge in large language models (LLMs) like LLaMA [6] and GPT-4 [7], an emerging trend is to adapt these expansive models to computer vision tasks [8], [9], [10], [11]. This transition has unveiled impressive capabilities; for instance, MiniGPT-4 [11] can create websites from handwritten drafts and generate detailed image captions. Moreover,

Manuscript received August 10, 2023; revised November 23, 2023; accepted November 23, 2023. Date of publication December 16, 2023; date of current version December 23, 2023. The associate editor coordinating the review of this manuscript and approving it for publication was Yue Gao (Corresponding author: Sidan Du.)

Yunzhuo Sun and Yukun Shu are with the School of Physics and Electronics, Hubei Normal University, Huangshi 435002, China (e-mail: sunyunzhuo98@outlook.com; Shuyk@hbnu.edu.cn)

Yifang Xu, Zien Xie and Sidan Du are with the School of Electronic Science and Engineering, Nanjing University, Nanjing 210000, China (e-mail: xyf@smail.nju.edu.cn; xze@smail.nju.edu.cn; coff128@nju.edu.cn).

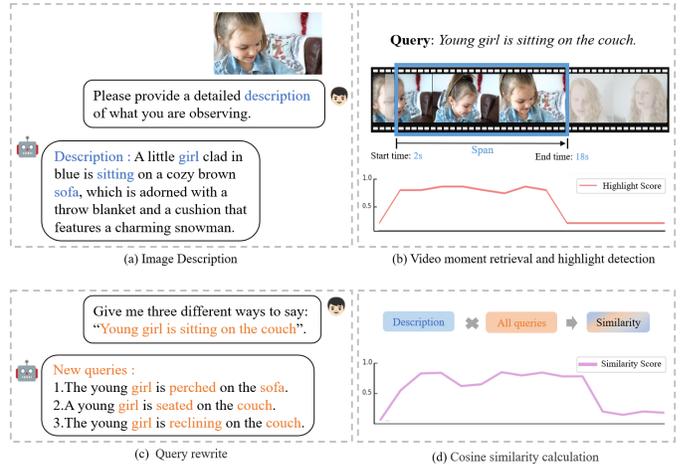


Figure 1. (a) Describe video frame content with GPT. (b) Examples of video moment retrieval and highlight detection (MR&HD) tasks. (c) Rewrite queries using GPT. (d) Calculate the cosine similarity between the image description and the rewritten query.

VideoChat [10] and Video-ChatGPT [11] have demonstrated adaptability to certain video understanding subtasks such as video summarization and video question answering. However, existing GPT-based video models encounter difficulties with more fine-grained subtasks like MR&HD. This is due to two main reasons. Firstly, MR&HD necessitates modeling of moment-level features. However, the upper limit on the context length in large models poses a significant constraint, reducing their performance. Secondly, these large models lack dedicated modules explicitly designed for MR&HD [12].

In this paper, we propose a two-stage stepwise optimization model, utilizing the output of the LLMs as input to the transformer encoder-decoder [13]. First, we extract a frame every two seconds from the video, converting them into textual descriptions using MiniGPT-4 [8]. Query rewrites with identical semantics are generated with the same model to explore semantic information. We then calculate the semantic similarity between the content description and queries, identifying the range for continuous video frames with high similarity, termed span anchors. Finally, the features and span anchors from MiniGPT-4 are input into the second-stage transformer's encoder and decoder, respectively. In fact, due to the instability of bipartite graph matching, DETR-like models tend to underperform. However, the high-quality features and span anchors facilitate the positioning in the second-stage model, thereby enhancing the final results.

Overall, our main contributions are as follows:

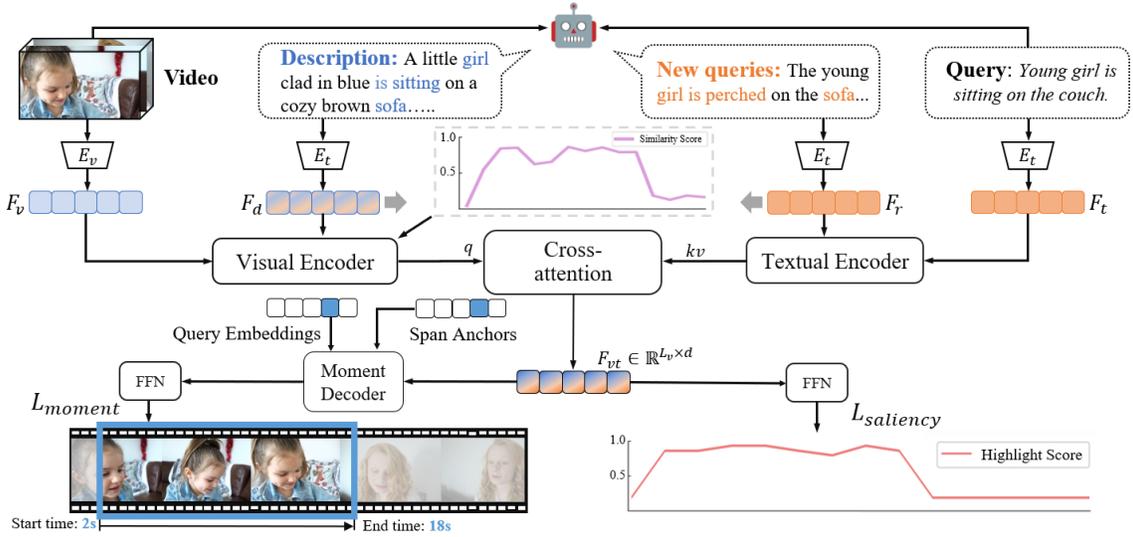


Figure 2. An overview of our proposed model GPTSee. Video frames and query text are initially fed into MiniGPT4, generating corresponding image content descriptions and semantically rewritten queries. Subsequently, the visual extractor  $E_v$  and text extractor  $E_t$  obtain features from these descriptions and rewritten queries, which are input into their respective visual and text encoders. In parallel, similarity scores are calculated based on the semantic similarity between the image content descriptions of key video frames and rewritten queries. The visual encoder jointly receives these scores, concatenated with the image description. The encoded visual and text features interact through a cross-attention mechanism, resulting in the cross-modal features  $F_{vt}$ . This feature is then directly processed by an FFN to derive the highlight scores for the HD task. Frames bearing consecutive high similarity scores form a range, referred to as span anchors, serving as prior position information for the moment decoder. Subsequently, for the MR task, this decoder establishes the start and end positions of video moments.

- 1) We use LLMs to generate detailed descriptions of images and rewrite queries, then compute the semantic similarity scores between them. The mentioned operation introduces three novel features for the MR&HD task.
- 2) We optimized the decoder module by leveraging high-quality prior positional information from the first stage, enhancing model performance.
- 3) We have conducted extensive experiments on the QVHighlights dataset, demonstrating that our method performs better than the current state-of-the-art approaches

## II. METHOD

### A. Overview

Given an untrimmed video  $V \in \mathbb{R}^{N_v \times H \times W \times 3}$  containing  $N_v$  moments and a natural language query  $T \in \mathbb{R}^{N_t}$  with  $N_t$  words, the task of MR&HD aims to localize all boundaries  $B \in \mathbb{R}^{N_b \times 2}$ , each comprising a start moment and an end moment, that are highly relevant to  $T$ . Simultaneously, it predicts moment-wise highlight scores  $H \in \mathbb{R}^{N_v}$  for the entire video. The overall structure of our approach, designed based on the foundational principles of Moment-DETR [14], is depicted in Figure 2.

Our process begins with generating detailed image descriptions and query rewrites. Utilizing MiniGPT-4 [8], we produce natural language descriptions for each video frame and create query rewrites that retain semantic similarity while introducing syntactic variations. With the aid of CLIP [15], visual features  $F_v \in \mathbb{R}^{L_v \times d_v}$  and textual features  $F_t \in \mathbb{R}^{L_t \times d_t}$  are extracted from raw videos and queries, respectively. Similarly, the CLIP text encoder  $E_t$  extracts features from the descriptions  $F_d$  and query rewrites  $F_r$ . Corresponding encoders process these

features to produce visual tokens  $\tilde{F}_v \in \mathbb{R}^{N_v \times d}$  and textual tokens  $\tilde{F}_t \in \mathbb{R}^{N_t \times d}$ , which are then fused by a cross-modal interaction module. The computation of similarity scores  $S \in \mathbb{R}^{N_s \times 1}$  between the visual and textual features identifies span anchors  $A \in \mathbb{R}^{L_v \times 2}$ , used as prior positional information in the decoder. The final stage employs a linear layer and sigmoid activation to estimate moments and highlight scores in the prediction head.

### B. Image Detail Description

Videos contain rich semantic information. Previous work has predominantly focused on extracting features from specific aspects of images, such as optical flow [16], [17], depth maps [18], [19], achieving noteworthy results. However, these prior approaches have overlooked the potential of translating the visual content of images into natural language descriptions. Inspired by the way humans perceive content within videos, our approach translates visual content into comprehensive text descriptions and feeds them into the encoder as an innovative feature form. This method encodes video content from a textual standpoint, introducing a new dimension to the field.

We employ MiniGPT-4 [8] to generate natural language descriptions for each frame of the input videos, as shown in Figure. 1 (a). MiniGPT-4, which integrates advanced language models with visual perception components, can generate content-rich and semantically coherent text descriptions. For instance, given a video frame depicting a little girl sitting on a sofa, MiniGPT-4 might produce a description such as "A little girl clad in blue is sitting on a cozy brown sofa." These descriptions serve as additional contextual information, feeding into the model to enrich its comprehension of the video content.

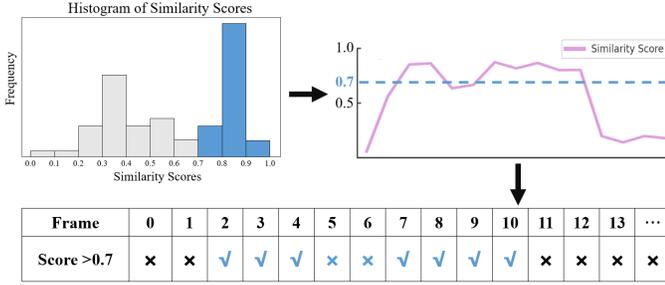


Figure 3. Determine the threshold and then collect the indices of all frames whose similarity scores exceed this threshold.

### C. Query Rewriting

To fully exploit the semantic information of queries, we have designed a query rewriting module, primarily relying on two strategies: semantic equivalent transformation and synonym substitution. For example, given an original query, "Young girl is sitting on the couch" can be transformed into "The young girl is perched on the sofa" through semantic equivalent transformation or into "A young girl is seated on the couch" by applying synonym substitution.

We guide MiniGPT-4 [8] to generate syntactically distinct sentences but semantically close sentences. In this way, the model can interpret and understand queries from various perspectives, thus enhancing its ability to handle ambiguous or unclear queries. The quality of these rewritten queries is validated by computing their semantic similarity to the original queries.

$$M(q, q_0) = \min \left( 1, \frac{P(q_0|q)}{P(q|q)} \right) \quad (1)$$

where  $M(q, q_0)$  is the measure of semantic similarity between the original query  $q$  and its paraphrase  $q_0$ ,  $P(q_0|q)$  is the probability of a paraphrase  $q_0$  given the original question  $q$ , and  $P(q|q)$  is used to normalize different distributions.

### D. Similarity Calculation and Span Anchors

When humans perform MR&HD tasks, it is natural to compare semantic similarity between video content and query text to determine relevance [20], [21]. Inspired by this, we analyze the similarity  $S \in \mathbb{R}^{N_s \times 1}$  and set thresholds to identify the relevant range. As shown in Figure 3, we calculate the cosine similarity between the image description and the rewritten query to quantify their relevance. Due to the inconsistent distribution of similarity scores among video groups, employing a fixed threshold for determining relevance is not advisable. Therefore, we analyze the distribution in each group, setting the threshold as the third most common value. Scores above this threshold are marked with a ✓, while those below are marked with a ×. If the number of × marks is less than six, the range is considered non-relevant and is defined as span anchors  $A \in \mathbb{R}^{L_v \times 2}$ . Experimental results demonstrate that, without additional training and solely utilizing the span anchors and similarity score, our model can surpass Moment-DETR [14] on the MR&HD task. According to [22], we feed  $S$  into the visual encoder and process  $A$  through a feed-forward network (FFN), incorporating it into query embeddings and then providing it into the moment decoder.

### E. Moment Decoder and Prediction Heads

We integrate visual tokens and textual tokens via cross-attention [23], [24], [25] to form  $F_{vt} \in \mathbb{R}^{L_v \times d}$ . An FFN with ReLU [26] predicts normalized moment center and width. Class label prediction utilizes a softmax linear layer. Predicted moments are assigned foreground or background based on alignment with ground truth. Another linear layer predicts highlight scores, as  $H \in \mathbb{R}^{N_v}$ .

*Moment Retrieval Loss.*  $L_m$ , measuring the between predicted  $\hat{m}$  and ground-truth moments  $m$ , is defined as:

$$L_m = \lambda_{L1} \|m - \hat{m}\|_1 + \lambda_{iou} L_{iou}(m, \hat{m}) \quad (2)$$

where  $\lambda_{L1}$  and  $\lambda_{iou}$  are real-valued hyperparameters. It combines the L1 loss and the generalized Intersection-over-Union (IoU) loss  $L_{iou}$  [27], which computes the temporal overlap between  $\hat{m}$  and  $m$ .

*Cross-entropy Loss.* We utilize the weighted binary cross-entropy loss to categorize predicted spans into foreground or background. This can be mathematically represented as:

$$L_{cls} = - \sum_{i=1}^{L_s} [w_p z_i \log(p_i) + (1 - z_i) \log(1 - p_i)] \quad (3)$$

In this equation,  $p_i$  and  $z_i$  denote the forecasted probability of the foreground and its respective label. The foreground label is attributed with a higher weight  $w_p$  to alleviate label imbalance.

*Highlight Detection Loss* The Highlight Detection Loss is designed to optimize the highlight score for each moment. This loss is computed using hinge loss across two distinct sets of segments:

$$L_h = \max(0, \delta + H(t_{low}) - H(t_{high})) + \max(0, \delta + H(t_{out}) - H(t_{in})) \quad (4)$$

Here, the first set comprises a high-score segment ( $t_{high}$ ) and a low-score segment ( $t_{low}$ ) within the actual temporal moments. The second set includes one segment ( $t_{in}$ ) inside and another ( $t_{out}$ ) outside the actual temporal moments.

*Total Loss.* The total loss is computed as a linear combination of the above losses:

$$L_{total} = L_m + \lambda_{cls} L_{cls} + \lambda_h L_h \quad (5)$$

## III. EXPERIMENTS

### A. Evaluation Dataset and Metrics Selection

Based on the QVHighlights [14], [28] dataset, we evaluated our model. QVHighlights is the only publicly accessible dataset with ground-truth labels for MR&HD. The dataset comprises 10,148 YouTube videos. Each video in the dataset is annotated with an unstructured textual query, associated time spans, and scores for significant moments. We followed a widely accepted data partitioning scheme (training, validation, testing) utilized in recent studies to ensure a fair comparison.

To measure the effectiveness of moment retrieval (MR), we employed metrics like Recall@1 with thresholds of 0.5 and 0.7, mean average precision (mAP) with intersection over union (IoU) thresholds of 0.5 and 0.75, and consolidation of mAP at various IoU thresholds in the range of 0.5 to 0.95 incremented by 0.05 were used. For highlight detection (HD),

mAP along with HIT@1 was employed, wherein HIT@1 accounts for instances where the moment with the highest score is correctly identified.

### B. Details of Implementation

Our model integrates a visual encoder, a textual encoder, and a cross-modal interaction module, all equipped with a singular attention layer, allowing seamless communication between the visual and textual information. Additionally, our moment decoder includes four self-attention layers. We apply a dropout rate 0.1, enhanced by post-normalization style layer normalization [29] and ReLU [26] activation functions. The loss function’s hyperparameters are set as follows:  $\lambda_{L1} = 10$ ,  $\lambda_{iou} = 1$ ,  $\lambda_{cls} = 4$ ,  $\lambda_h = 1$  and  $w_p = 10$ . We use the AdamW optimizer [30], with a learning rate set at  $2e-4$  and a weight decay parameter of  $1e-4$ . We conducted training over 200 epochs with a batch size of 32, utilizing 8 RTX 3090.

### C. Experimental Results

Initially, we provide an extensive comparison of our proposed GPTSee model with preceding state-of-the-art models on the QVHighlights dataset, as presented in Table I. Across all metrics, our model consistently excels over the UniVTG [31]. Notably, there are increments of 5.86% and 1.03% in MR-mAP Avg. and HD-mAP, respectively, underscoring our model’s robustness and efficacy in the MR&HD task.

Following that, we perform a targeted comparison between our GPTSee model and Moment-DETR [14] in Table II, demonstrating the benefits of employing span anchors for evaluation. GPTSee surpasses Moment-DETR, even when using only span anchors (*A*) and similarity (*S*) as output, particularly for MR-mAP, HD-mAP, and HD-HIT@1 metrics.

### D. Ablation Studies

**Ablations on Transfer Capability:** As shown in Table II, the integration of *A* and *S* into Moment-DETR and UMT was tested to validate the generality of our two-stage approach. The results reveal that this method improves accuracy. Notably, using *A+S* alone for final localization, the HD-HIT@1 scores excel, showing strong performance in local highlight detection but only average overall. By incorporating the similarity into the encoder, our method achieved a significant enhancement in overall HD tasks, along with a slight improvement in local.

**Ablations on LLMs:** As depicted in Table III, when selecting models for generating image detail descriptions and rewriting queries, a comparison was made among VideoChat [10], Video-ChatGPT [11], and MiniGPT-4, utilizing the same prompt for generating descriptions and rewriting queries. The conclusion ascertains that MiniGPT-4 slightly outperforms the other two. This is likely due to MiniGPT-4 generating fewer irrelevant words during sentence creation, thus having an advantage in semantic similarity computation.

**Ablations on Module Configuration:** As shown in Table IV, to validate the effectiveness of each model component, several baseline models were constructed with varying components. The analysis reveals that the description and span anchors contribute most significantly to the overall model performance, while query rewrite provides only a marginal contribution.

TABLE I  
PERFORMANCE COMPARISON ON QVHIGHLIGHTS TEST SPLIT.

Methods	R1		MR			HD	
	@0.5	@0.7	@0.5	@0.75	Avg.	$\geq$ Very Good	HIT@1
CAL [32]	25.49	11.54	23.40	7.65	9.89	-	-
XML [33]	41.83	30.35	44.63	31.73	32.14	34.49	55.25
XML+ [33]	46.69	33.46	47.89	34.67	34.90	35.38	55.06
Moment-DETR [14]	52.89	33.02	54.82	29.40	30.73	35.69	55.60
UMT [34]	56.23	41.18	53.83	37.01	36.12	38.18	59.99
UniVTG [31]	58.86	40.86	57.60	35.59	35.47	38.20	60.96
GPTSee (Ours)	<b>62.84</b>	<b>48.01</b>	<b>61.92</b>	<b>42.55</b>	<b>41.33</b>	<b>39.23</b>	<b>62.80</b>

TABLE II  
PERFORMANCE COMPARISON OF SPAN ANCHORS AND SIMILARITY SCORES AS DIRECT OUTPUTS ON QVHIGHLIGHTS TEST SPLIT.

Model.	MR			HD ( $\geq$ VG)	
	R1@0.5	R1@0.7	mAP Avg.	mAP	HIT@1
A + S	56.55	31.60	32.27	36.11	<b>61.65</b>
Moment-DETR [14]	54.82	29.40	30.73	35.69	55.60
Moment-DETR + A + S	<b>60.20</b>	<b>34.88</b>	<b>36.29</b>	<b>36.45</b>	59.04
UMT [34]	53.83	37.01	36.12	38.18	59.99
UMT [34] + A + S	<b>60.74</b>	<b>40.05</b>	<b>38.01</b>	<b>38.58</b>	61.07

TABLE III  
COMPARISON OF SPAN ANCHORS AND SIMILARITY SCORES FROM DIFFERENT LLMS IN QVHIGHLIGHTS TEST SPLIT.

Model.	MR			HD ( $\geq$ VG)	
	R1@0.5	R1@0.7	mAP Avg.	mAP	HIT@1
VideoChat[10]	60.23	47.34	39.10	37.64	60.14
Video-ChatGPT[11]	<b>63.20</b>	47.89	41.12	38.78	62.00
MiniGPT-4[8]	62.84	<b>48.01</b>	<b>41.33</b>	<b>39.23</b>	<b>62.80</b>

TABLE IV  
ABLATION STUDY OF DIFFERENT FEATURES ON QVHIGHLIGHTS TEST SPLIT.

Description	Features			MR	HD ( $\geq$ VG)
	Rewritten Queries	Similarity	Span Anchors	mAP Avg.	mAP
✓				36.54	36.71
	✓			34.40	36.40
		✓		36.15	37.28
			✓	37.52	37.18
✓	✓			37.56	36.99
✓		✓		38.36	37.74
✓			✓	39.30	37.66
	✓	✓		37.04	37.44
		✓	✓	37.00	37.53
		✓	✓	38.09	38.10
✓	✓	✓		39.65	38.01
✓	✓		✓	38.76	37.98
✓		✓	✓	40.87	38.61
	✓	✓	✓	38.32	38.50
✓	✓	✓	✓	<b>41.33</b>	<b>39.23</b>

## IV. CONCLUSION

In this letter, we introduce an innovative two-stage model, GPTSee, which integrates LLMs’ output to assist a transformer encoder-decoder architecture. The effectiveness of our model is significantly enhanced by incorporating image descriptions and rewritten queries as novel inputs and using span anchors as a priori positional information. This framework amplifies video frames and query text data utilization, eliminating the need for intricate feature extraction or elaborate training schemas characteristic of earlier approaches. Experiments on the QVHighlights dataset substantiate our model’s superiority and efficacy. Future work may focus on designing more powerful LLMs or improving highlight score calculation.

## REFERENCES

- [1] S. Ghosh, A. Agarwal, Z. Parekh, and A. G. Hauptmann, "ExCL: Extractive Clip Localization Using Natural Language Descriptions," in *NAACL*, 2019.
- [2] S. Zhang, H. Peng, J. Fu, and J. Luo, "Learning 2d temporal adjacent networks for moment localization with natural language," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, 2020, pp. 12 870–12 877, issue: 07.
- [3] Y. Xu, Y. Sun, Z. Xie, B. Zhai, and S. Du, "Vtg-gpt: Tuning-free zero-shot video temporal grounding with gpt," *Applied Sciences*, vol. 14, no. 5, p. 1894, 2024.
- [4] M. Xu, H. Wang, B. Ni, R. Zhu, Z. Sun, and C. Wang, "Cross-category video highlight detection via set-based learning," in *CVPR*, 2021, pp. 7970–7979.
- [5] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, "Video summarization with long short-term memory," in *ECCV*. Springer, 2016, pp. 766–782.
- [6] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, and F. Azhar, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [7] "Introducing ChatGPT." [Online]. Available: <https://openai.com/blog/chatgpt>
- [8] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models," *arXiv preprint arXiv:2304.10592*, 2023.
- [9] T. Gong, C. Lyu, S. Zhang, Y. Wang, M. Zheng, Q. Zhao, K. Liu, W. Zhang, P. Luo, and K. Chen, "Multimodal-gpt: A vision and language model for dialogue with humans," *arXiv preprint arXiv:2305.04790*, 2023.
- [10] K. Li, Y. He, Y. Wang, Y. Li, W. Wang, P. Luo, Y. Wang, L. Wang, and Y. Qiao, "Videochat: Chat-centric video understanding," *arXiv preprint arXiv:2305.06355*, 2023.
- [11] M. Maaz, H. Rasheed, S. Khan, and F. S. Khan, "Video-ChatGPT: Towards Detailed Video Understanding via Large Vision and Language Models," Jun. 2023, arXiv:2306.05424 [cs]. [Online]. Available: <http://arxiv.org/abs/2306.05424>
- [12] H. Zhang, A. Sun, W. Jing, and J. T. Zhou, "Temporal Sentence Grounding in Videos: A Survey and Future Directions," Oct. 2022, arXiv:2201.08071 [cs] version: 2. [Online]. Available: <http://arxiv.org/abs/2201.08071>
- [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, \. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 5998–6008.
- [14] J. Lei, T. L. Berg, and M. Bansal, "Detecting Moments and Highlights in Videos via Natural Language Queries," *NeurIPS*, vol. 34, pp. 11 846–11 858, 2021.
- [15] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, and J. Clark, "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021, pp. 8748–8763.
- [16] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *ECCV*. Springer, 2020, pp. 402–419.
- [17] H. Xu, J. Zhang, J. Cai, H. Rezatofighi, and D. Tao, "Gmflow: Learning optical flow via global matching," in *CVPR*, 2022, pp. 8121–8130.
- [18] Y. Xu, M. Li, C. Peng, Y. Li, and S. Du, "Dual attention feature fusion network for monocular depth estimation," in *CAAI International Conference on Artificial Intelligence*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:245639385>
- [19] Y. Xu, C. Peng, M. Li, Y. Li, and S. Du, "Pyramid feature attention network for monocular depth prediction," in *2021 IEEE International Conference on Multimedia and Expo (ICME)*. Shenzhen, China: IEEE, Jul 2021, pp. 1–6. [Online]. Available: <https://ieeexplore.ieee.org/document/9428446/>
- [20] L. Zhen, P. Hu, X. Wang, and D. Peng, "Deep supervised cross-modal retrieval," in *CVPR*, 2019, pp. 10 386–10 395. [Online]. Available: <https://api.semanticscholar.org/CorpusID:198906771>
- [21] S. Chun, S. J. Oh, R. S. de Rezende, Y. Kalantidis, and D. Larlus, "Probabilistic embeddings for cross-modal retrieval," in *CVPR*, June 2021, pp. 8415–8424.
- [22] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "DAB-DETR: Dynamic Anchor Boxes are Better Queries for DETR," in *ICLR*, 2022.
- [23] L. Zhen, P. Hu, X. Peng, R. S. M. Goh, and J. T. Zhou, "Deep multimodal transfer learning for cross-modal retrieval," *IEEE Trans Neural Netw Learn Syst*, vol. 33, no. 2, pp. 798–810, Feb 2022.
- [24] M. Cheng, Y. Sun, L. Wang, X. Zhu, K. Yao, J. Chen, G. Song, J. Han, J. Liu, E. Ding, and J. Wang, "Vista: Vision and scene text aggregation for cross-modal retrieval," in *CVPR*, June 2022, pp. 5184–5193.
- [25] Y. Xu, Y. Sun, Z. Xie, B. Zhai, Y. Jia, and S. Du, "Query-guided refinement and dynamic spans network for video highlight detection and temporal grounding in online information systems," *Int. J. Semant. Web Inf. Syst.*, vol. 19, no. 1, pp. 1–20, Jun 2023.
- [26] G. E. Dahl, T. N. Sainath, and G. E. Hinton, "Improving deep neural networks for LVCSR using rectified linear units and dropout," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 8609–8613.
- [27] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *CVPR*, 2019, pp. 658–666.
- [28] Y. Xu, Y. Sun, Y. Li, Y. Shi, X. Zhu, and S. Du, "Mh-detr: Video moment and highlight detection with cross-modal transformer," *arXiv preprint arXiv:2305.00355*, 2023.
- [29] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [30] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [31] K. Q. Lin, P. Zhang, J. Chen, S. Pramanick, D. Gao, A. J. Wang, R. Yan, and M. Z. Shou, "Univtg: Towards unified video-language temporal grounding," in *arXiv preprint*, 2023.
- [32] V. Escorcia, M. Soldan, J. Sivic, B. Ghanem, and B. Russell, "Temporal localization of moments in video collections with natural language," *arXiv preprint arXiv:1907.12763*, 2019.
- [33] J. Lei, L. Yu, T. L. Berg, and M. Bansal, "Tvr: A large-scale dataset for video-subtitle moment retrieval," in *ECCV*. Springer, 2020, pp. 447–463.
- [34] Y. Liu, S. Li, Y. Wu, C.-W. Chen, Y. Shan, and X. Qie, "UMT: Unified Multi-modal Transformers for Joint Video Moment Retrieval and Highlight Detection," in *CVPR*, 2022, pp. 3042–3051.