# QUERY-BASED VIDEO SUMMARIZATION WITH PSEUDO LABEL SUPERVISION

*Jia-Hong Huang*[1*], *Luka Murn*[2], *Marta Mrak*[2], *Marcel Worring*[1]

[1]University of Amsterdam, Amsterdam, Netherlands ; [2]BBC Research and Development, London, UK

## ABSTRACT

Existing datasets for manually labelled query-based video summarization are costly and thus small, limiting the performance of supervised deep video summarization models. Self-supervision can address the data sparsity challenge by using a pretext task and defining a method to acquire extra data with pseudo labels to pre-train a supervised deep model. In this work, we introduce segment-level pseudo labels from input videos to properly model both the relationship between a pretext task and a target task, and the implicit relationship between the pseudo label and the human-defined label. The pseudo labels are generated based on existing human-defined frame-level labels. To create more accurate query-dependent video summaries, a semantics booster is proposed to generate context-aware query representations. Furthermore, we propose mutual attention to help capture the interactive information between visual and textual modalities. Three commonly-used video summarization benchmarks are used to thoroughly validate the proposed approach. Experimental results show that the proposed video summarization algorithm achieves state-of-the-art performance.

***Index Terms***— Query-based video summarization, semantics, self-supervision, weak supervision, pseudo labels

## 1. INTRODUCTION

Query-based video summarization automatically generates a short video clip to summarize the content of a given video by capturing its query-dependent parts, as shown in Fig. 1. Such a task can be modeled as a fully-supervised machine learning problem [1, 2, 3]. However, creating a large-scale manually-labeled video dataset for a fully-supervised task is costly. Hence, existing datasets, e.g., TVSum [4], SumMe [5], and QueryVS [2], are quite small.

The lack of larger human-annotated datasets is common in fully-supervised deep learning tasks. Self-supervised learning is one of the most successful ways to alleviate this challenge [6, 7, 8, 9]. According to [7, 10], self-supervision is an effective method to balance the cost of data labelling and the performance gain of a fully-supervised deep model. The main idea of self-supervised learning is defining a pretext task and
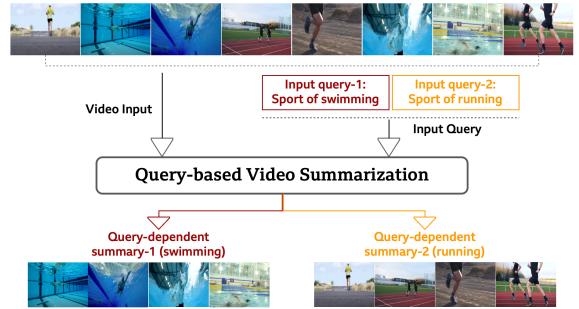


**Fig. 1**: Query-based video summarization. A video is summarized based on textual queries. The summarization algorithm runs independently for each query.

introducing a way to acquire extra data with reliable pseudo labels to pre-train a fully-supervised deep model for performing a target task [6, 7].

Existing self-supervision methods assume that the relation between a target task with human-defined labels and an introduced pretext task with pseudo labels does not exist or exists in a very limited way [7, 10]. However, this assumption may not be accurate for query-based video summarization, where frame-level human-defined labels can be considered as supervision signals of a target task. Segment-level pseudo labels can be considered as supervision signals of a pretext task. Since a video segment is composed of frames, there is an implicit relation between the entire segment and the corresponding frames. The improvement in model performance can hit a bottleneck without modelling these implicit relations.

In this work, a segment-based video summarization pretext task with specially designed pseudo labels is introduced to address this challenge, detailed in Fig. 2. Pseudo labels are generated based on existing human-defined annotations, helping to model the implicit relations between the pretext task and the target task, i.e., frame-based video summarization [2, 4, 5]. In query-based video summarization, we observe that generating accurate query-dependent video summaries can be challenging in practice due to ineffective semantics embedding of textual queries. We address this issue by proposing a semantics booster that generates context-aware query representations which are capable of efficiently capturing the semantics. Furthermore, we noticed that the query

---

input does not always help model performance, most likely due to the interactions between textual and visual modalities not being properly modelled. We address this challenge by introducing mutual attention that helps capture the interactive information between different modalities.

These novel design choices enable us to improve the model performance of query-based video summarization with self-supervision. Extensive experiments show that the proposed method is effective and achieves state-of-the-art performance. If we examine the problem from the perspective of frame-level label vs. segment-level label, the proposed method can also be considered as a weakly-supervised video summarization approach. Hence, existing weakly-supervised methods are also considered as baselines in this work.

## 2. RELATED WORK

### 2.1. Fully-supervised video summarization

Fully-supervised learning is a common way to model video summarization [5, 11, 12, 13, 14]. In fully-supervised video summarization, labels defined by human experts are used to supervise a model in the training phase. In [5], a video summarization approach is proposed to automatically summarize user videos that contain a set of interesting events. The authors start by dividing a video based on a superframe segmentation, tailored to raw videos. Then, various levels of features are used to predict the score of visual interestingness per superframe. Finally, a video summary is produced by selecting a set of superframes in an optimized way. In [12, 13], a Recurrent Neural Network (RNN) is used in a hierarchical way to model the temporal structure in video data. The authors of [11] consider video summarization as a problem of structured prediction. A deep-learning-based method is proposed to estimate the importance of video frames based on modelling their temporal dependency. The authors of [14] propose an importance propagation-based collaborative teaching network (iPT-Net) for video summarization by transferring samples from a video moment localization correlated task equipped with a lot of training data. In [2, 3, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34], the model learning process expands beyond solely utilizing visual inputs and incorporates an additional modality, such as viewers' comments, video captions, or any other contextual data available.

The aforementioned fully-supervised methods exploit a full set of human expert annotations to supervise the model in the training phase. Although such a method performs well, it is costly. Therefore, a better solution should be developed for video summarization.

### 2.2. Weakly-supervised video summarization

In [35, 36, 37, 38, 39], video summarization is considered as a weakly-supervised learning task. Weakly-supervised learning can mitigate the need for extensive datasets with human expert annotations. Instead of using a full set of data with human expert labels, such as frame-level annotations, weakly-supervised approaches exploit less-expensive weak labels, such as video-level annotations from human experts. Although weak labels are imperfect compared to a full set of human expert annotations, they still can be used to train video summarization models effectively.

### 2.3. Self-supervision in video summarization

In [40, 41], image pretext tasks [7] are extended to video for self-supervision in video summarization. In [40], the keyframes of a video are defined as those which are very different in their optical flow features and appearance from the rest of the frames of the video. The authors of [41] claim that a good video sequence encoder should have the ability to model the correct order of video segments. Segments are selected from a given video based on a fixed proportion before feeding it into a neural network. They are randomly shuffled and used to train the neural network and distinguish the odd-position segments to control the difficulty of the auxiliary self-supervision task.

Existing work related to self-supervision in video summarization is very limited, and they do not focus on query-based video summarization. To the best of our knowledge, our proposed method is one of the pioneer works of self-supervision in query-based video summarization.

### 2.4. Word embedding methods

According to [42], static word embeddings and contextualized word representations are commonly used to encode textual data. Both of them are more effective than the Bag of Words (BoW) method. Skip-gram with negative sampling (SGNS) [43] and GloVe [44] are well-known models for generating static word embeddings. According to [45, 46], these models learn word embeddings iteratively in practice. However, it has been proven that both of them implicitly factorize a word-context matrix containing a co-occurrence statistic.

The authors of [42] mention that in static word embeddings methods, all meanings of a polysemous word must share a single vector because a single representation for each word is created. Hence, the contextualized word representations method is more effective than static word embeddings because of its context-sensitive word representations. In [47, 48, 49], the proposed neural language models are fine-tuned to create deep learning-based models for a wide range of downstream natural language processing tasks.

In this work, a contextualized word representation-based method is used to encode the text-based input query.

## 3. METHODOLOGY

In this section, the proposed query-based video summarization method is described in detail, and illustrated in Fig. 2. The approach is based on contextualized query representations, attentive convolutional 2D and 3D features, interactive

**Fig. 2**: Flowchart of the proposed self-supervision method for query-based video summarization. The model is pre-trained by the textual-spatial features from the Mutual Attention Mechanism and pseudo segment-level labels. The completely trained video summary generator exploits the fully-connected layer to produce a frame-level score vector for the given input video and outputs the final query-dependent video summary.
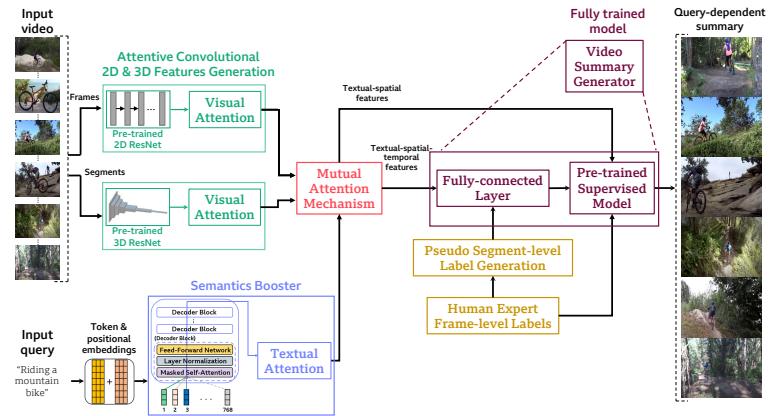


attention mechanism, mean-based pseudo shot label generation, and video summary generation.

### 3.1. Semantics Booster

Generating an accurate query-dependent video summary is challenging because of the ineffective semantics embedding of input textual queries. In this work, a semantics booster is introduced to capture the semantics of the input query effectively. The transformer-based model architecture has been firmly established as one of the state-of-the-art approaches in language modeling and machine translation [50]. Hence, the proposed semantics booster is built on top of the transformer architecture to generate context-aware query representations, described as follows.

For an input token $k_n$, its embedding $x_n$ is defined as: $x_n = W_e * k_n + P_{k_n}, n \in \{1, ..., N\}$, where $W_e \in \mathbb{R}^{E_s \times V_s}$ is the input text-based query token embedding matrix with the vocabulary size $V_s$ and the word embedding size $E_s$, the positional encoding of $k_n$ is $P_{k_n}$, and $N$ denotes the number of input tokens. The subscripts $s$ and $e$ denote size and embedding, respectively. The representation of the current word $Q$ is generated by one linear layer defined as: $Q = W_q * x_n + b_q$, where $b_q$ and $W_q \in \mathbb{R}^{H_s \times E_s}$ are learnable parameters of the linear layer, the output size of the linear layer is $H_s$ and the subscript $q$ denotes query. The key vector $K$ is calculated by the other linear layer defined as: $K = W_k * x_n + b_k$, where $b_k$ and $W_k \in \mathbb{R}^{H_s \times E_s}$ are learnable parameters of the linear layer. The subscript $k$ denotes key. The value vector $V$ is generated by another linear layer defined as: $V = W_v * x_n + b_v$, where $b_v$ and $W_v \in \mathbb{R}^{H_s \times E_s}$ are learnable parameters of the linear layer. The subscript $v$ denotes value.

After $Q$, $K$, and $V$ are calculated, the masked self-attention is generated as: $\text{MaskAtten}(Q, K, V) = \text{softmax}(m(\frac{QK^T}{\sqrt{d_k}}))V$, where $m(\cdot)$ and $d_k$ denote a masked self-attention function and a scaling factor, respectively. The layer normalization is calculated as: $Z_{\text{Norm}} = \text{LayerNorm}(\text{MaskAtten}(Q, K, V))$, where $\text{LayerNorm}(\cdot)$ de-

notes a layer normalization function. Then, the introduced context-aware representation $\mathcal{R}_{\text{context}}$ of the input text-based query is derived as: $\mathcal{R}_{\text{context}} = \sigma(W_1 Z_{\text{Norm}} + b_1)W_2 + b_2$, where $\sigma$ is an activation function, $W_1$, $W_2$, $b_1$, and $b_2$ are learnable parameters of a position-wise feed-forward network. To have even better textual representations, a textual attention function $\text{TextAtten}(\cdot)$ is introduced to reinforce the context-aware representation. The function takes $\mathcal{R}_{\text{context}}$ as input and calculates the attention and textual representation in an element-wise way. The attentive context-aware representation is calculated as $Z_{ta} = \text{TextAtten}(\mathcal{R}_{\text{context}})$, where $ta$ indicates textual attention.

### 3.2. Visual Attention

A 2D ConvNet and a 3D ConvNet are exploited to distill the video frame and video segment information, respectively. To reinforce the generated 2D and 3D features, a visual attention function $\text{AttenVisual}(\cdot)$ is introduced to improve the quality of features.

Let $E$ and $X$ be a feature generator and a set of video clips, respectively. A feature generator $E$ maps an input $x \in X$ to a feature vector $f \in \mathbb{R}^d$. $F = \{f = E(x) \in \mathbb{R}^d \mid x \in X\}$ denotes a set of features produced by the feature generator $E$. Let $F_s$ be the generated features from the video spatial feature generator $E_s$. $F_{st}$ denotes the generated features from the video spatio-temporal feature generator $E_{st}$. Frame-level and segment-level data both are exploited to train the proposed query-based video summarization model, meaning $F = F_s \cup F_{st}$. In the frame-level case, the attentive feature generator $\text{AttenVisual}(\cdot)$ learns attention weights and produces attentive spatial features $Z_{as} = \{f_{as} = \text{AttenVisual}(f) \in \mathbb{R}^d \mid f \in F_s\}$, i.e., attentive convolutional 2D features. In the segment-level case, the attentive feature generator learns attention weights and produces attentive spatio-temporal features $Z_{ast} = \{f_{ast} = \text{AttenVisual}(f) \in \mathbb{R}^d \mid f \in F_{st}\}$, i.e., attentive convolutional 3D features.

## 3.3. Mutual Attention

We observe that textual queries do not always help the model performance due to the interactions between the video and query inputs not being modelled effectively. In this work, a mutual attention mechanism MutualAtten$(\cdot)$ is introduced to address this issue and model the interactive information between the video and query. The mutual attention $Z_{ma}$ performs one by one convolution, i.e., convolutional attention. $Z_{ma} = \text{MutualAtten}(Z_{ta} \odot Z_{as} \odot Z_{ast})$, where $Z_{ta}$ indicates textual attention and $\odot$ denotes Hadamard product.

## 3.4. Pseudo Segment-level Label Generation

Let $S_f$ be a set of human experts' frame-level score annotations and $P$ a pseudo score annotation generator that maps frame-level human expert scores to a segment-level pseudo score.

In [4], the authors empirically find that a two-second segment is suitable for capturing local context of a video as it achieves good visual coherence. Based on this observation, in this work the proposed pseudo label generator $P$ is designed to generate a segment-level score every two seconds. In practice, since the generated pseudo score annotations are not validated by human experts, they might contain noisy or biased information. Based on [51], the Mean function is one of the effective ways to reduce the noise contained in the segment-level pseudo label. Hence, Mean function is used to design the proposed pseudo label generator $P$ to produce the mean score $S_{\text{mean}} = P(S_f) = \text{Mean}(S_f)$, i.e., the two-second segment-level pseudo score label. In the training phase, compared with the frame-level label, the mean-based pseudo segment label $S_{\text{mean}}$ is used not only for spatial supervision but also for temporal supervision. The temporal supervision with the segment-level pseudo annotations improves the query-based video summarization model performance.

## 3.5. Loss Function

According to [2], query-based video summarization can be modeled as a classification problem. Thus, in this work, the categorical cross-entropy loss function is adopted to build the proposed approach:

$$\text{Loss} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} \mathbf{1}_{y_i \in C_c} \log(P_{\text{model}}[y_i \in C_c]), \quad (1)$$

where $N$ indicates the number of observations, $C$ denotes the number of categories, $\mathbf{1}_{y_i \in C_c}$ is an indicator function of the $i$-th observation belonging to the $c$-th category, and $P_{\text{model}}[y_i \in C_c]$ is the probability predicted by the model for the $i$-th observation to belong to the $c$-th category.

## 4. EXPERIMENTS AND ANALYSIS

### 4.1. Datasets and evaluation metrics

**Datasets.** TVSum [4] is a commonly used dataset for traditional video summarization, containing only the video in-

| Pseudo label pre-training | Mutual attention | Semantics booster | TVSum | QueryVS |
|---|---|---|---|---|
| - | - | - | 47.5 | 50.8 |
| ✓ | - | - | 61.3 | 52.9 |
| - | ✓ | - | 58.9 | 52.0 |
| - | - | ✓ | 56.4 | 52.3 |
| ✓ | ✓ | ✓ | **68.4** | **55.3** |

**Table 1**: Ablation study of the pseudo segment-level label pre-training, semantics booster, and mutual attention mechanism using $F_1$-score.

put. However, authors of [17, 18] consider TVSum metadata, e.g., video title, as a text-based query input to generate the query-dependent video summary. In our experiments, the TVSum dataset is randomly divided into 40/5/5 videos for training/validation/testing, respectively. The video length is ranging from 2 to 10 minutes. The human expert score labels range from 1 to 5, and are annotated with 20 frame-level responses per video [18].

The SumMe [5] dataset is randomly divided into 19 videos for training, 3 videos for validation, and 3 videos for testing. The video duration in SumMe is ranging from 1 to 6 minutes. In SumMe, the human expert annotation score ranges from 0 to 1. SumMe is not used for query-based video summarization and we do not have a query input when a model is evaluated on this dataset.

QueryVS [2] is an existing dataset designed for query-based video summarization. In our experiments, the QueryVS dataset is separated into 114/38/38 videos for training/validation/testing, respectively. The video length in QueryVS is ranging from 2 to 3 minutes, and every video is retrieved based on a given text-based query.

To validate the proposed query-based video summarization method, three segment-level datasets are created based on the above frame-level datasets. Both the segment-level dataset, i.e., for pre-training, and the frame-level dataset, i.e., the target dataset, are used to conduct our experiments.

**Evaluation metric.** Based on [4, 5, 17, 18, 52], the $F_\beta$-score with the hyper-parameter $\beta = 1$ is a commonly used metric for assessing the performance of supervised video summarization approaches. It is based on measuring the agreement between the predicted score and ground truth score provided by the human expert. The $F_\beta$-score is defined as: $F_\beta = \frac{1}{N} \sum_{i=1}^{N} \frac{(1+\beta^2) \times p_i \times r_i}{(\beta^2 \times p_i) + r_i}$, where $r_i$ indicates $i$-th recall, $p_i$ indicates $i$-th precision, $N$ indicates number of $(r_i, p_i)$ pairs, "$\times$" denotes scalar product, and $\beta$ is used to balance the relative importance between recall and precision.

### 4.2. Experimental settings

In the experiments, a 2D ResNet-34 network pre-trained on the ImageNet database [55] is adopted to generate frame-level features for each input video. The 512 features are extracted from the visual layer one layer below the classification layer. A 3D ResNet-34 pre-trained on the Kinetics benchmark [56] is used in the experiments to generate segment-level features

| Model | Method | TVSum | SumMe | QueryVS |
|-------|--------|-------|-------|---------|
| vsLSTM [11] | Fully supervised | 54.2 | 37.6 | - |
| H-RNN [12] | | 57.7 | 41.1 | - |
| HSA-RNN [13] | | 59.8 | 44.1 | - |
| iPTNet [14] | | 63.4 | 54.5 | - |
| SMLD [53] | | 61.0 | 47.6 | - |
| SMN [54] | | 64.5 | **58.3** | - |
| FPVSF [36] | Weakly supervised | - | 41.9 | - |
| WS-HRL [39] | | 58.4 | 43.6 | - |
| DSSE [17] | Query based | 57.0 | - | - |
| DQSN [18] | | 58.6 | - | - |
| QueryVS [2] | | - | - | 41.4 |
| GPT2MVS [3] | | - | - | 54.8 |
| Ours | | **68.4** | 52.4 | **55.3** |

**Table 2**: Comparison with state-of-the-art video summarization methods based on the $F_1$-score, best highlighted in bold. '-' denotes unavailability from previous work.



**Fig. 3**: Randomly selected qualitative results of the proposed method. Selected frames from the ground truth frame-based score annotations for the input video are highlighted in gray, with red representing the frames not selected. Frames selected for the query-dependent video summary are highlighted in green. 217 denotes the video length before video preprocessing and 647 denotes the video length after the video preprocessing.

for each input video. The features with 512 dimensions are located in the visual layer which is right after the global average pooling layer.

The video lengths in the SumMe, TVSum and QueryVS datasets vary, with the maximum number of frames in a video being 388 for SumMe, 199 for QueryVS, and 647 for TVSum. A frame-repeating preprocessing technique [2] is followed to make all the videos in each dataset the same length.

The input size of the CNN is 224 by 224 with RGB channels. Every channel is normalized by standard deviation $= (0.2737, 0.2631, 0.2601)$ and mean $= (0.4280, 0.4106, 0.3589)$. PyTorch is used for the implementation and to train models for 100 epochs with $1e - 7$ learning rate. The Adam optimizer is used, with hyper-parameters set as $\epsilon = 1e - 8$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$.

### 4.3. Ablation Study

The ablation study of the proposed method is presented in Table 1. The baseline model without the mutual attention mechanism and pseudo segment label pre-training and no semantics booster performs significantly worse than approaches utilising any or all of the proposed improvements. Note that when the semantics booster is not adopted, the BoW embedding method is used.

The mutual attention mechanism helps capture the interaction between the input query and video more effectively. The pseudo segment-level label pre-training helps the proposed model have better initialization. The semantics booster captures the semantic meaning of the text-based query.

### 4.4. Comparison with state-of-the-art models

The comparison with existing fully-supervised, weakly-supervised and query-based approaches is presented in Table 2. The results show the performance of our proposed method is the best on TVSum and QueryVS datasets, with a competitive performance on the SumMe dataset.

The correctness of the generated segment-level pseudo labels is not guaranteed by human experts, but it still contains useful information, e.g., better temporal information, to supervise the proposed model during pre-training. In weakly-supervised methods, although the correctness of the coarse labels, e.g., video-level label, is guaranteed by human experts, it is still not good enough to boost the model performance better than our proposed method. In query-based summarization methods, although the other modality is used to help the model performance, the effectiveness of the multi-modal feature fusion could limit the performance improvement.

Randomly selected qualitative results are shown in Fig. 3.

## 5. CONCLUSION

In this work, a new query-based video summarization approach is proposed. The method is based on the self-supervision of segment-level pseudo scores, semantics booster, and a mutual attention mechanism. Additionally, three segment-level video summarization datasets for self-supervision are proposed based on existing small-scale query-based video summarization datasets. Experimental results show the mean-based segment-level pseudo labels provide effective temporal supervision. The proposed approach achieves state-of-the-art performance in terms of the $F_1$-score. Nowadays, video content is growing at an ever-increasing speed and beyond the capacity of an individual for full comprehension. In such cases, the proposed query-based video summarization method has the potential to improve the efficiency of video exploration.

## 6. ACKNOWLEDGMENTS

# 7. REFERENCES

[1] Arun Balajee Vasudevan, Michael Gygli, Anna Volokitin, and Luc Van Gool, "Query-adaptive video summarization via quality-aware relevance estimation," in *MM*, 2017, pp. 582–590.

[2] Jia-Hong Huang and Marcel Worring, "Query-controllable video summarization," in *ICMR*, 2020, pp. 242–250.

[3] Jia-Hong Huang, Luka Murn, Marta Mrak, and Marcel Worring, "Gpt2mvs: Generative pre-trained transformer-2 for multi-modal video summarization," in *ICMR*, 2021, pp. 580–589.

[4] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes, "Tvsum: Summarizing web videos using titles," in *CVPR*, 2015, pp. 5179–5187.

[5] Michael Gygli, Helmut Grabner, Hayko Riemenschneider, and Luc Van Gool, "Creating summaries from user videos," in *ECCV*. Springer, 2014, pp. 505–520.

[6] Carl Doersch, Abhinav Gupta, and Alexei A Efros, "Unsupervised visual representation learning by context prediction," in *ICCV*, 2015, pp. 1422–1430.

[7] Humam Alwassel, Dhruv Mahajan, Bruno Korbar, Lorenzo Torresani, Bernard Ghanem, and Du Tran, "Self-supervised learning by cross-modal audio-video clustering," *arXiv preprint arXiv:1911.12667*, 2019.

[8] Zihang Lai and Weidi Xie, "Self-supervised learning for video correspondence flow," *arXiv preprint arXiv:1905.00875*, 2019.

[9] Dahun Kim, Donghyeon Cho, and In So Kweon, "Self-supervised video representation learning with space-time cubic puzzles," in *AAAI*, 2019, vol. 33, pp. 8545–8552.

[10] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze, "Deep clustering for unsupervised learning of visual features," in *ECCV*, 2018, pp. 132–149.

[11] Ke Zhang, Wei-Lun Chao, Fei Sha, and Kristen Grauman, "Video summarization with long short-term memory," in *ECCV*. Springer, 2016, pp. 766–782.

[12] Bin Zhao, Xuelong Li, and Xiaoqiang Lu, "Hierarchical recurrent neural network for video summarization," in *MM*, 2017, pp. 863–871.

[13] Bin Zhao, Xuelong Li, and Xiaoqiang Lu, "Hsa-rnn: Hierarchical structure-adaptive rnn for video summarization," in *CVPR*, 2018, pp. 7405–7414.

[14] Hao Jiang and Yadong Mu, "Joint video summarization and moment localization by cross-task sample transfer," in *CVPR*, 2022, pp. 16388–16398.

[15] Jia-Hong Huang, Chao-Han Huck Yang, Pin-Yu Chen, Andrew Brown, and Marcel Worring, "Causal video summarizer for video exploration," in *2022 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2022, pp. 1–6.

[16] Jia-Hong Huang, Chao-Han Huck Yang, Pin-Yu Chen, Min-Hung Chen, and Marcel Worring, "Causalainer: Causal explainer for automatic video summarization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 2629–2635.

[17] Yitian Yuan, Tao Mei, Peng Cui, and Wenwu Zhu, "Video summarization by learning deep side semantic embedding," *Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 1, pp. 226–237, 2017.

[18] Kaiyang Zhou, Tao Xiang, and Andrea Cavallaro, "Video summarisation by classification with deep reinforcement learning," *arXiv:1807.03089*, 2018.

[19] Ting-Wei Wu, Jia-Hong Huang, Joseph Lin, and Marcel Worring, "Expert-defined keywords improve interpretability of retinal image captioning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 1859–1868.

[20] Riccardo Di Sipio, Jia-Hong Huang, Samuel Yen-Chi Chen, Stefano Mangini, and Marcel Worring, "The dawn of quantum natural language processing," *ICASSP*, 2022.

[21] C-H Huck Yang, Jia-Hong Huang, Fangyu Liu, Fang-Yi Chiu, Mengya Gao, Weifeng Lyu, Jesper Tegner, et al., "A novel hybrid machine learning model for auto-classification of retinal diseases," *Workshop on Computational Biology, ICML*, 2018.

[22] Yi-Chieh Liu, Hao-Hsiang Yang, C-H Huck Yang, Jia-Hong Huang, Meng Tian, Hiromasa Morikawa, Yi-Chang James Tsai, and Jesper Tegner, "Synthesizing new retinal symptom images by multiple generative models," in *ACCV*. Springer, 2018, pp. 235–250.

[23] C-H Huck Yang, Fangyu Liu, Jia-Hong Huang, Meng Tian, I-Hung Lin, Yi Chieh Liu, Hiromasa Morikawa, Hao-Hsiang Yang, and Jesper Tegner, "Auto-classification of retinal diseases in the limit of sparse data using a two-streams machine learning model," in *ACCV*. Springer, 2018, pp. 323–338.

[24] Jia-Hong Huang, Ting-Wei Wu, Chao-Han Huck Yang, and Marcel Worring, "Longer version for" deep context-encoding network for retinal image captioning"," *arXiv preprint arXiv:2105.14538*, 2021.

[25] Jia-Hong Huang, Ting-Wei Wu, Chao-Han Huck Yang, and Marcel Worring, "Deep context-encoding network for retinal image captioning," in *ICIP*. IEEE, 2021, pp. 3762–3766.

[26] Jia-Hong Huang, Ting-Wei Wu, and Marcel Worring, "Contextualized keyword representations for multi-modal retinal image captioning," in *ICMR*, 2021, pp. 645–652.

[27] Jia-Hong Huang, Ting-Wei Wu, C-H Huck Yang, Zenglin Shi, I Lin, Jesper Tegner, Marcel Worring, et al., "Non-local attention improves description generation for retinal images," in *WACV*, 2022, pp. 1606–1615.

[28] Jia-Hong Huang, C-H Huck Yang, Fangyu Liu, Meng Tian, Yi-Chieh Liu, Ting-Wei Wu, I Lin, Kang

Wang, Hiromasa Morikawa, Hernghua Chang, et al., "Deepopht: medical report generation for retinal images via deep models and visual explanation," in *WACV*, 2021, pp. 2442–2452.

[29] Jia-Hong Huang, Modar Alfadly, Bernard Ghanem, and Marcel Worring, "Assessing the robustness of visual question answering," *arXiv preprint arXiv:1912.01452*, 2019.

[30] Jia-Hong Huang, Modar Alfadly, Bernard Ghanem, and Marcel Worring, "Improving visual question answering models through robustness analysis and in-context learning with a chain of basic questions," *arXiv preprint arXiv:2304.03147*, 2023.

[31] Jia-Hong Huang, "Robustness analysis of visual question answering models by basic questions," *King Abdullah University of Science and Technology, Master Thesis*, 2017.

[32] Jia-Hong Huang, Modar Alfadly, and Bernard Ghanem, "Robustness analysis of visual qa models by basic questions," *VQA Challenge and Visual Dialog Workshop, CVPR*, 2018.

[33] Jia-Hong Huang, Modar Alfadly, and Bernard Ghanem, "Vqabq: Visual question answering by basic questions," *VQA Challenge Workshop, CVPR*, 2017.

[34] Jia-Hong Huang, Cuong Duc Dao, Modar Alfadly, and Bernard Ghanem, "A novel framework for robustness analysis of visual qa models," in *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 8449–8456.

[35] Rameswar Panda, Abir Das, Ziyan Wu, Jan Ernst, and Amit K Roy-Chowdhury, "Weakly supervised summarization of web videos," in *ICCV*, 2017, pp. 3657–3666.

[36] Hsuan-I Ho, Wei-Chen Chiu, and Yu-Chiang Frank Wang, "Summarizing first-person videos from third persons' points of view," in *ECCV*, 2018, pp. 70–85.

[37] Sijia Cai, Wangmeng Zuo, Larry S Davis, and Lei Zhang, "Weakly-supervised video summarization using variational encoder-decoder and web prior," in *ECCV*, 2018, pp. 184–200.

[38] Tao Hu, Pascal Mettes, Jia-Hong Huang, and Cees GM Snoek, "Silco: Show a few images, localize the common object," in *ICCV*, 2019, pp. 5067–5076.

[39] Yiyan Chen, Li Tao, Xueting Wang, and Toshihiko Yamasaki, "Weakly supervised video summarization by hierarchical reinforcement learning," in *MM Asia*, pp. 1–6. 2019.

[40] Xiang Yan, Syed Zulqarnain Gilani, Mingtao Feng, Liang Zhang, Hanlin Qin, and Ajmal Mian, "Self-supervised learning to detect key frames in videos," *Sensors*, vol. 20, no. 23, pp. 6941, 2020.

[41] Yudong Jiang, Kaixu Cui, Bo Peng, and Changliang Xu, "Comprehensive video understanding: Video summarization with content-based video recommender design," in *ICCVW*, 2019, pp. 0–0.

[42] Kawin Ethayarajh, "How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings," *arXiv preprint arXiv:1909.00512*, 2019.

[43] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean, "Distributed representations of words and phrases and their compositionality," *arXiv preprint arXiv:1310.4546*, 2013.

[44] Jeffrey Pennington, Richard Socher, and Christopher D Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.

[45] Omer Levy and Yoav Goldberg, "Linguistic regularities in sparse and explicit word representations," in *Proceedings of the eighteenth conference on computational natural language learning*, 2014, pp. 171–180.

[46] Omer Levy and Yoav Goldberg, "Neural word embedding as implicit matrix factorization," *NIPS*, vol. 27, pp. 2177–2185, 2014.

[47] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.

[48] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer, "Deep contextualized word representations," *arXiv preprint arXiv:1802.05365*, 2018.

[49] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, pp. 9, 2019.

[50] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.

[51] Tang Zhi, Li Jingwen, and Zhou Yinqing, "Analysis on noise reduction method for interferometric sar image," in *IGARSS*, 2004, vol. 6, pp. 4243–4246.

[52] George Hripcsak and Adam S Rothschild, "Agreement, the f-measure, and reliability in information retrieval," *Journal of the American medical informatics association*, vol. 12, no. 3, pp. 296–298, 2005.

[53] Wei-Ta Chu and Yu-Hsin Liu, "Spatiotemporal modeling and label distribution learning for video summarization," in *MMSP*, 2019, pp. 1–6.

[54] Junbo Wang, Wei Wang, Zhiyong Wang, Liang Wang, Dagan Feng, and Tieniu Tan, "Stacked memory network for video summarization," in *MM*, 2019, pp. 836–844.

[55] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.

[56] Joao Carreira and Andrew Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *CVPR*, 2017, pp. 6299–6308.