

# Learning Pixel-Level Distinctions for Video Highlight Detection

Fanyue Wei<sup>1\*</sup> Biao Wang<sup>2</sup> Tiezheng Ge<sup>2</sup> Yuning Jiang<sup>2</sup> Wen Li<sup>1</sup> Lixin Duan<sup>1</sup>

<sup>1</sup>School of Computer Science and Engineering & Shenzhen Institute for Advanced Study, UESTC <sup>2</sup>Alibaba Group  
{wfanyue, liwenbnu, lxduan}@gmail.com, {eric.wb, tiezheng.gtz, mengzhu.jyn}@alibaba-inc.com

## Abstract

The goal of video highlight detection is to select the most attractive segments from a long video to depict the most interesting parts of the video. Existing methods typically focus on modeling relationship between different video segments in order to learning a model that can assign highlight scores to these segments; however, these approaches do not explicitly consider the contextual dependency within individual segments. To this end, we propose to learn pixel-level distinctions to improve the video highlight detection. This pixel-level distinction indicates whether or not each pixel in one video belongs to an interesting section. The advantages of modeling such fine-level distinctions are two-fold. First, it allows us to exploit the temporal and spatial relations of the content in one video, since the distinction of a pixel in one frame is highly dependent on both the content before this frame and the content around this pixel in this frame. Second, learning the pixel-level distinction also gives a good explanation to the video highlight task regarding what contents in a highlight segment will be attractive to people. We design an encoder-decoder network to estimate the pixel-level distinction, in which we leverage the 3D convolutional neural networks to exploit the temporal context information, and further take advantage of the visual saliency to model the spatial distinction. State-of-the-art performance on three public benchmarks clearly validates the effectiveness of our framework for video highlight detection.

## 1. Introduction

Along with the explosive development of mobile devices, a tremendous number of videos are now produced and uploaded to the Internet every day. As a result, picking the most attractive video clips from a lengthy video to create a selection of shining moments is becoming increasingly important, especially for social video platforms such as YouTube and Instagram. As a result, video highlight de-

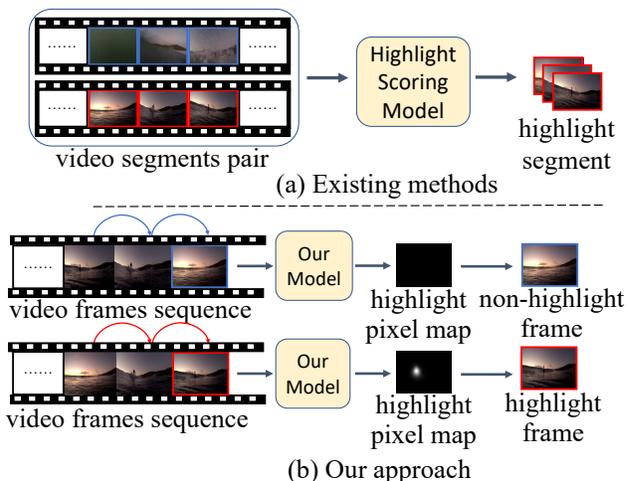


Figure 1. Video highlight detection is highly context-dependent. While previous methods are usually trained to predict the highlight score for a video segment directly, our method takes the temporal and spatial information into account and predicts the fine-level pixel-level distinction as the surrogate task.

tection, which aims to select the most attractive segments from an unedited video, has drawn increasing interest from in the research community.

Most existing works [3, 4, 10] interpret the video highlight detection task as a segment-level ranking problem. These approaches treat each segment as an individual sample and extract the features for video segments. They then compare pairwise segments in order to learn a model that assigns highlight scores to these segments, such that the highlight segments receive higher scores than the non-highlight segments from the raw video. Recently, SL [19] developed a set-based mechanism that is capable of identifying whether or not a video segment is highlight by transformer.

However, these existing methods do utilize both temporal and localized information, but not explicitly considering the contextual dependency within the segment, which is in fact crucial for video highlight detection.

Intuitively, when people watch videos, a specific part is considered to be interesting, usually depends on the previ-

\*Work done during an internship at Alibaba Group

ous parts they have watched. For example, considering a video in which a gymnast performs a somersault, the jumping up before the somersault and jumping down after the somersault are visually quite similar; however, people tend to rate the jumping up as more appealing than the jumping down; because the former contributes to the climax of the somersault, while the latter decays the highlight level after climax. This indicates that predicting the highlight score of one frame highly depends on the context before the current frame.

Similarly, the spatial context is also important for video highlight detection. A dog might not be interesting if it appears together with a group of dogs, while it will definitely be the focus in a dog show scenario. In this case, the context information within one frame would be very helpful for estimating the highlight score.

Accordingly, to exploit the temporal and spatial context for video highlight detection, in this paper, we cast the video highlight detection into a new task: pixel-level distinction estimation. More specifically, rather than assigning highlight scores to video segments (as in the existing works), we aim to predict the attractiveness of each pixel in the video. Such fine-level task offers two benefits. First, as the distinction of a pixel in one frame often depends on the temporal and spatial context, predicting pixel-level distinction allows to exploit such context information in our model, leading to more robust highlight detection results. Second, learning the pixel-level distinction also offers a good explanation for the video highlight task also to what content in a highlight segment might be more appealing to people, making the video highlight detection model more explainable. After estimating the pixel-level distinction, the highlight score of a video segment can be readily obtained by averaging the distinctions of all pixels in the segment.

We develop an encoder-decoder network to estimate the pixel-level distinction. This network is designed to output a distinction map for each frame in the input video. To exploit the temporal context, we employ a 3D convolutional neural networks to incorporate the frames before the current frame in order to predict the distinction map. To model the spatial distinction, we take advantage of the visual saliency to generate pixel-level pseudo-distinction labels for frames in the highlight segments. We demonstrate that the strategies discussed above can be simply integrated into the encoder-decoder network.

Experiments on three challenging benchmarks-YouTube [8], TvSum [15] and CoSum [16]-show that our proposed approach outperforms existing methods by clear margins. We further validate the effectiveness of our proposed modules with ablation studies, and provide qualitative results to show the explainable ability of our proposed model.

In summary, the main contributions of this paper are as

follows:

- We propose a new pixel-level distinction estimation task for video highlight detection, which is able to explore the fine-level context in order to predict the attractiveness of specific segments.
- We design an encoder-decoder network for estimating the pixel-level distinction, which takes advantage of the 3D convolutional neural networks and the visual saliency map to exploit the temporal and spatial context, respectively.
- We achieve new state-of-the-art performance on three public benchmarks. Moreover, our model also exhibits good explainable ability, and is able to directly output the most appealing regions in the highlighted video segments.

## 2. Related Work

### 2.1. Video Highlight Detection

The goal of video highlight detection is to find the most attractive parts of videos. Prior methods have largely focused on generating highlight from sports videos [11–13, 17]. More recent approaches focus on addressing the Internet videos and first-person videos. These recent methods can be divided into two aspects: supervised and unsupervised (or weakly supervised).

Supervised methods mainly treat video highlight detection as a segment-level ranking or scoring task [3, 4, 10, 19, 20]. These methods generally construct a pair-wise ranking constraint for two video segments, highlight and non-highlight segments. Video2GIF [3] proposes a method to learn from manually generated video-GIF pairs. By utilizing an adaptive Huber loss to overcome noisy data, a robust deep RankNet can generate a ranked list of video segments.

Additionally, GNN [5] introduces object semantics to video highlight task and further model the relationships between objects via graph neural network.

SL [19] utilizes transformer structure to capture the multiple segments that contributes to the target segment. Moreover, SA [20] proposes that audio and visual information are highly related to highlight detection. They fuse audio and visual information by attention for video highlight detection.

Unsupervised or weakly supervised methods often introduces some prior information as a supervised signal rather than using the highlight annotations for training.

LIM-s [2] leverages the video duration as the implicit supervised signal. They contend that user-generated videos have the relationship that video segments from shorter videos are more likely to be highlights than those from longer videos. Therefore, they propose a model that learns

to score the highlight segments higher than non-highlight segments. More recently, MINI-Net [1] casts video highlight detection as multiple instance learning problem. They characterize each video as a bag of segments, aiming to score a positive bag about specific event higher than a negative bag that events are irrelevant.

Most methods generate the highlight clips by ranking the highlight and non-highlight segments based on segment-level feature representation. In a departure from the existing methods, our work captures the visual temporal distinction via sliding window and introduces visual saliency to model the highlight spatially with pixel-level loss.

## 2.2. Video Summary

The goal of video summarization, which is highly related to video highlights, is to produce the most informative clip that incorporates the complete plot of an entire given video [30–33]. Video summary models often learn to score a sequence of selected frames [23] or clips [24]. Additionally, some video summary methods consider not only importance but also representation [22], diversity [21] and coherency [25]. [26] aims to select a subset of frames that optimally represent a given video, performing unsupervised video summarization with adversarial LSTM networks. The model comprises a summarizer, which aims to obtain an optimal summarization of a new video, and a discriminator, to distinguish between the original video and its reconstruction obtained from the summarizer. [22] formalizes video summarization into a sequential decision-making process. By training an end-to-end reinforcement learning framework, the proposed model predicts frame-level probability to be chosen to form summary. Moreover, some methods [27, 28] take a hierarchical recurrent neural network to exploit the long temporal dependency among frames for video summarization. [29] captures the temporal dependencies with LSTM and GCN hierarchically.

## 2.3. Visual Saliency

Visual saliency aims to model the gaze fixation. Previous methods have utilized optical flow to make use of temporal information [45, 46]. Moreover, some methods aggregate temporal information with LSTM [47]. ACLNet [48] enhances the ability of LSTM to capture the dynamic saliency through the use of a frame-wise attention mask. Recently, TASED-Net [38] aggregates temporal information and takes an encoder-decoder structure spatially to predict frame-wise saliency maps in a sliding-window fashion for a given video. STAViS [49] combines spatiotemporal auditory and visual information to address video saliency.

## 3. Approach

In this paper, we propose to leverage the temporal and spatial relations within a segment to improve video high-

light detection. Our motivation derives from the fact that video highlights is highly context-dependent; *i.e.*, whether or not the content in a video segment should be highlighted depends on the content that comes before it in temporal dimension and the content surrounding it in the spatial dimension.

Paradigms in previous methods [3, 5, 8, 19, 20] have proposed different strategies for learning the scoring function  $f(s_i)$ , which largely involve assigning higher scores to highlight segments and lower scores to non-highlight segments. These methods usually obtain the whole feature representation from each video segment, and learn score function differently. However, these methods tend to ignore the spatial-temporal relation of the content among frames within each segment, which is in fact crucial to video highlight detection.

To capture this context-dependent property, rather than scoring each video segment on its global feature representation, we instead propose to predict the highlight score for each pixel per frame, an approach referred to as *Pixel-Level Distinction Video Highlight Detection (PLD-VHD)* in the paper.

## 3.1. Modeling Temporal Dependency

Formally, for one video  $V$ , denote  $S = \{s_1, \dots, s_n\}$  as the set of video segments after division, where each  $s_i$  is a segment for  $i = 1, \dots, n$ . Each  $s_i$  is accompanied by a label  $y_i$ , where  $y_i = 1$  indicates that  $s_i$  is a highlight segment while  $y_i = 0$  denotes opposite.

We begin from a basic model for pixel-level distinction estimation. As the ground-truth pixel-level distinction  $d_t(i, j)$  is unknown, we need to construct a pseudo-distinction label for each pixel using the segment-level highlight label  $y_s$ . The basic concept of our approach is simple. For frames from non-highlight segments  $s_n$ , we set the distinction labels of all pixels as zeros; for those from highlight segments  $s_h$ , their distinction labels are set as ones. The pseudo-distinction label can be defined as follows:

$$d_t(i, j) = \begin{cases} 1, & I_t \in s_h \\ 0, & I_t \in s_n \end{cases} \quad (1)$$

For simplicity, we use  $D_t$  to represent the distinction map for a frame  $I_t$ , where  $d_t(i, j)$  is the pseudo-distinction label defined in Eq. (1). For ease of presentation, we also use  $f(I_t)$  to denote the pixel-level distinction estimation function for the entire frame  $I_t$ . We then take a simple mean squared error (MSE) as the loss, and the problem for learning pixel-level distinctions can be formulated as:

$$\min \mathcal{L}(f(I_t), D_t) = \frac{1}{W \cdot H} \sum_{i=1}^W \sum_{j=1}^H (p_t(i, j) - d_t(i, j))^2 \quad (2)$$

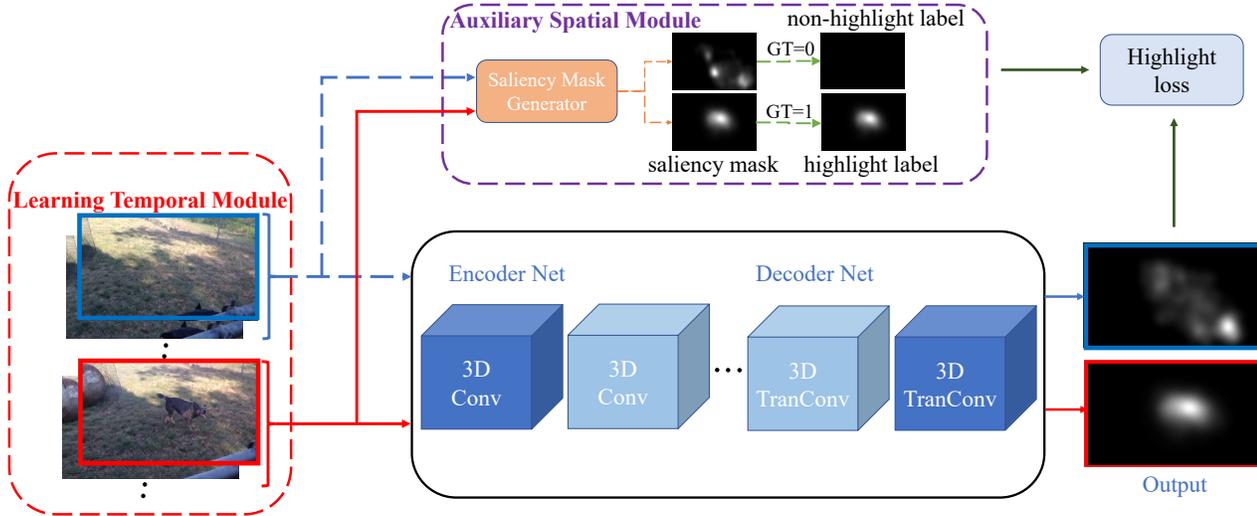


Figure 2. Our network follows an encoder-decoder structure. The encoder net is a 3D ConvNet for extracting features of input frames, while the decoder net aims to obtain a map that is the same size as the input frames for pixel-level distinction. The learning temporal module incorporates the previous frames before the target frame, while the auxiliary spatial module generates the pseudo labels. The frames in blue and red boxes represent targets with two kinds of labels: highlight and non-highlight, respectively.

where  $f(\cdot)$  is the distinction estimation function,  $p_t(i, j)$  denotes the pixel-level distinctions obtained by  $f(\cdot)$ ,  $W$  and  $H$  denote the width and height of the frame.

The distinction estimation function can be implemented with an encoder-decoder network. The input frame  $I_t$  are first fed into an encoder to obtain the latent feature representation, after which the feature map is upsampled by the decoder network for pixel-level distinction prediction.

However, as discussed above, creating video highlights is highly context-dependent. When people watch videos, the current frame becomes interesting because people have watched the previous frames. This means that distinctions of the current frame should depend on the frames that came before it.

Accordingly, to take the temporal dependency into account, rather than directly using each frame as input, we use a video clip to predict the pixel-level distinction. This clip contains both the current frame and a number of frames before it. Given the  $t$ -th frame  $I_t$ , let us denote the corresponding video clip as  $C_t = \{I_{t-L+1}, I_{t-L+2}, \dots, I_t\}$  where  $L$  is the total length of the video clip. The model for predicting the pixel-level distinctions can thus be updated as follows:

$$\min \mathcal{L}(f(C_t), D_t) \quad (3)$$

where  $\mathcal{L}$  is the MSE loss defined as in Eq. (2).

In our implementation, we apply an  $L$ -length sliding window to the video to generate the video clips. For the first  $T - 1$  frames in each video, we reverse the order of these frames and pad them to the beginning of the video to

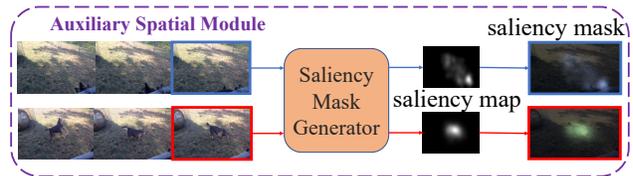


Figure 3. The spatial distinction of our approach. With the mask generated by our Saliency Mask Generator, the pseudo-label can eliminate the noise (such as the background) that makes no contribution to the highlight.

ensure the sliding window works. Each video clip is then fed into a 3D convolutional neural network (*e.g.*, C3D [34] or TASED [38]) to automatically exploit the temporal relation among frames within each video clip, as illustrated in Figure 1.

### 3.2. Spatial Highlight with Visual Saliency

In addition to the contextual temporal dependency within video clips, the attractiveness of an object is also often dependent on the surrounding context. For example, a single dog might not be particularly appealing when it appears in a group of dogs in one image, but it is the shining star in a dog show scenario.

We therefore further consider the spatial relationship for learning pixel-level distinctions. In particular, for non-highlight segments, their pseudo-labels are still all-zeros,

Table 1. Video highlight detection results of different methods on the YouTube Highlights dataset.

	LSVM	RRAE	Video2gif	LIM-s	MINI-Net	AFM-F-M	GNN	SL	SA	PLD-VHD
dog	0.60	0.49	0.308	0.579	0.582	0.72	0.67	0.708	0.649	<b>0.749</b>
gymnastics	0.41	0.35	0.335	0.417	0.617	0.56	0.66	0.532	<b>0.715</b>	0.702
parkour	0.61	0.50	0.540	0.670	0.702	0.75	<b>0.83</b>	0.772	0.766	0.779
skating	0.62	0.25	0.554	0.578	0.722	0.68	0.70	<b>0.725</b>	0.606	0.575
skiing	0.36	0.22	0.328	0.486	0.587	0.64	0.69	0.661	<b>0.712</b>	0.707
surfing	0.61	0.49	0.541	0.651	0.651	0.78	0.69	0.762	0.782	<b>0.790</b>
Average	0.536	0.412	0.464	0.564	0.644	0.68	0.69	0.693	0.705	<b>0.730</b>

Table 2. Experimental results (top-5 mAP score) of compared methods on the TVsum dataset.

	KVS	DPP	sLstm	SM	Quasi	MBF	CVS	SG	LIM-s	VESD	DSN	MINI-Net	SL	SA	PLD-VHD
BK	0.342	0.395	0.406	0.407	0.295	0.313	0.326	0.417	0.663	0.441	0.368	0.717	0.726	0.681	<b>0.845</b>
BT	0.419	0.464	0.471	0.473	0.327	0.365	0.402	0.483	0.691	0.492	0.435	0.769	0.789	<b>0.950</b>	0.809
DS	0.394	0.449	0.455	0.453	0.309	0.357	0.378	0.466	0.626	0.488	0.416	0.591	0.640	0.608	<b>0.703</b>
FM	0.397	0.442	0.452	0.451	0.318	0.365	0.365	0.464	0.432	0.487	0.412	0.559	0.589	0.669	<b>0.725</b>
GA	0.402	0.457	0.463	0.469	0.342	0.325	0.379	0.475	0.612	0.496	0.428	0.754	0.749	<b>0.844</b>	0.764
MS	0.417	0.462	0.477	0.478	0.375	0.412	0.398	0.489	0.54	0.503	0.436	0.813	0.862	0.865	<b>0.872</b>
PK	0.382	0.437	0.448	0.445	0.324	0.318	0.354	0.456	0.604	0.478	0.411	0.780	<b>0.790</b>	0.703	0.719
PR	0.403	0.446	0.461	0.458	0.301	0.334	0.381	0.473	0.475	0.485	0.417	0.545	0.632	0.675	<b>0.740</b>
VT	0.353	0.399	0.411	0.415	0.336	0.295	0.328	0.423	0.559	0.447	0.373	0.803	<b>0.865</b>	0.834	0.744
VU	0.441	0.453	0.462	0.467	0.369	0.357	0.413	0.472	0.429	0.493	0.441	0.653	0.687	0.647	<b>0.791</b>
Average	0.398	0.447	0.451	0.461	0.329	0.345	0.372	0.462	0.563	0.481	0.424	0.698	0.733	0.748	<b>0.771</b>

since none of the pixels in these segments are of interest. For highlight segments, we take advantage of the visual saliency to exploit the spatial context in each frame.

On one hand, as Figure 3 shows, visual saliency, which can be seen as robust general visual signals, aims to model the gaze fixation people display when they are watching videos, which is in line with the goal of video highlight detection. Using saliency helps us to identify the fine-level regions that attract people. On the other hand, while we annotate the pixel-level distinctions for all pixels in the highlight segments in Eq. (1), not all regions in the highlight segments are truly attractive, which yields a considerable amount of noise when optimizing the learning problem in Eq. (3). Using saliency information to eliminate the background noise facilitates the learning of a more robust pixel-level distinction estimation model.

More specifically, we use the saliency mask as the pseudo-labels to annotate the pixel-level distinctions for pixels in the highlight segments. Given any frame  $I_t$  in the highlight segment, we denote its saliency mask as  $M_t$ . The pixel-level distinctions can thus be defined as follows:

$$\hat{d}_t(i, j) = \begin{cases} 0, & M_t(i, j) \leq \beta \\ 1, & M_t(i, j) > \beta \end{cases} \quad (4)$$

where  $\beta$  is a hyper-parameter threshold we simply set as 0.0005 in most cases.

Note that, by using the above definition of pixel-level distinctions  $\hat{d}_t(i, j)$  to replace the original pixel-level dis-

tinctions  $d_t(i, j)$  in Eq. (1), the spatial distinction can be seamlessly integrated with the temporal dependency learning framework. We are also able to jointly exploit the spatial and temporal dependencies in order to estimate the pixel-level distinctions. Denoting  $\hat{D}_t$  as the new distinction map for frame  $I_t$ , the learning objective can be updated as follows:

$$\min \mathcal{L}(f(C_t), \hat{D}_t) \quad (5)$$

where  $\mathcal{L}$  is the MSE loss defined as in Eq.(2).

After learning the pixel-level distinction estimation model, given a segment from any video, the highlight score can be calculated by averaging all pixel-level distinctions in the segments, as follows:

$$f(s_d) = \sum_{t=1}^N \sum_{i=1}^H \sum_{j=1}^W \frac{1}{N \cdot H \cdot W} (f(C_t)^{(i,j)}) \quad (6)$$

where  $s_d$  denotes the  $d$ -th segment in a video, while  $f(C_t)^{(i,j)}$  is the  $(i, j)$ -th element of the estimated distinction map for  $I_t$ ; moreover,  $N$ ,  $H$  and  $W$  respectively denote the number of frames in  $s_d$ , the height and the width of the frames. The highlight score of a video segment can be estimated by using the mean of all pixel-level distinctions in this segment, while the highlighted video can be obtained by ensembling the video segments with highest scores similarly as in existing video highlight detection works.

Table 3. Results (top-5 mAP score) on the CoSum Dataset. Our method outperforms all comparison methods by a large margin.

	KVS	DPP	sLstm	SM	SMRS	Quasi	MBF	CVS	SG	VESD	DSN	MINI-Net	PLD-VHD
BJ	0.662	0.672	0.683	0.692	0.504	0.561	0.631	0.658	0.698	0.685	0.715	0.776	<b>0.900</b>
BP	0.674	0.682	0.701	0.722	0.492	0.625	0.592	0.675	0.713	0.714	0.746	0.963	<b>0.970</b>
ET	0.731	0.744	0.749	0.789	0.556	0.575	0.618	0.722	0.759	0.783	0.813	0.786	<b>0.817</b>
ERC	0.685	0.694	0.717	0.728	0.525	0.563	0.575	0.693	0.729	0.721	0.756	0.953	<b>1.000</b>
KP	0.701	0.705	0.714	0.745	0.521	0.557	0.594	0.707	0.729	0.742	0.772	0.959	<b>1.000</b>
MLB	0.668	0.677	0.714	0.693	0.543	0.563	0.624	0.679	0.721	0.687	0.727	0.869	<b>1.000</b>
NFL	0.671	0.681	0.681	0.727	0.558	0.587	0.603	0.674	0.693	0.724	0.737	0.897	<b>0.970</b>
NDC	0.698	0.704	0.722	0.759	0.496	0.617	0.694	0.702	0.738	0.751	0.782	0.890	<b>0.958</b>
SL	0.713	0.722	0.721	0.766	0.525	0.551	0.624	0.715	0.743	0.763	0.794	0.787	<b>0.844</b>
SF	0.642	0.648	0.653	0.653	0.533	0.562	0.603	0.647	0.681	0.674	0.709	0.727	<b>1.000</b>
Average	0.684	0.692	0.705	0.735	0.525	0.576	0.602	0.687	0.720	0.721	0.755	0.861	<b>0.946</b>

### 3.3. Network Architecture

As shown in Figure 2, our model is constituted by an encoder network (used to extract the features of the input video clip) and a decoder network (used to generate a distinction map corresponding to the target frame). The *Temporal Module* is designed to obtain the auxiliary past information for the current frame to be predicted. It utilizes 3D convolution neural network input consisting of past consequent frames and the current target frame. Moreover, the *Auxiliary Spatial Module* is a visual saliency model that works as an encoder network to generate a saliency mask. For this purpose, we here adopt TASED-Net [38] pretrained on DHF1K [48]. The final output of our whole framework is a highlight map as same size as the input target frame.

## 4. Experiments

In this section, we validate our model on several challenging public benchmarks-YouTube [8], TvSum [15] and CoSum [16]-and compare the results with those of several state-of-the-art video highlight detection methods. More experimental details are reported in the Supplementary Materials.

### 4.1. Experimental Setup

#### 4.1.1 Dataset and Evaluation Metric

- *YouTube Highlight* [8] is a popular video highlight detection dataset that collects videos from six different domains. Each domain contains 50 to 90 videos with varying duration. Each video is divided into several segments that contain approximately 100 frames, each of which are annotated with three different kinds of labels: 1-selected by users as highlight; 0-borderline cases, and -1-non-highlight. We treat the borderline cases as non-highlight.
- *TvSum* [15] contains 50 videos of 10 classes. Following [1, 2] we select top 50% of shots in terms of the

scores provided by annotators for each video as the human-created ground truth.

- *CoSum* [16] consists of 51 videos of 10 events. In this work, following [1], we compare each generated highlight with three human-created ground truth.

Like most existing methods [1, 2], we follow Video2gif [3], using mean average precision(mAP) as the evaluation metric.

#### 4.1.2 Comparison Methods

We compare our methods (*PLD-VHD*) with the following state-of-the-art video highlight detection baselines on three datasets.

- *Weakly supervised methods.* The comparison methods in this category are RRAE [9], MBF [16], SMRS [50], Quasi [51], CVS [52], SG [26], VESD [53], DSN [54], LIM-s [2] and MINI-Net [1].
- *Supervised methods.* There are also several supervised methods selected for comparison, *i.e.*, LSVM [8], Video2gif [3], KVS [55], DPP [21], sLstm [56], SM [24], AFM-F-M [10], GNN [5], SL [19] and SA [20].

Although some of these methods are used for video summarization, following [1, 2], their performance is evaluated using the same metrics as those used in this study.

### 4.2. Video Highlight Detection Results

The public datasets contain videos under different situations, such as camera view changes. In terms of the overall experimental results, our proposed method with pseudo-distinction labels achieves the best performance. Table 1 presents the results of different methods for video highlight detection on the *YouTube Highlight dataset* [8]. We report the results of our proposed approach using TASED-Net [38] as the backbone networks, denoted by “*PLD-VHD*”. For



Figure 4. Showcase examples from different domain. Red means a higher highlight score, light green indicates a lower highlight score, and blue represents a medium highlight score.

Table 4. Results of ablation studies on YouTube Highlights dataset.

	C3D w/o temporal	C3D w/o spatial	C3D full	TASED w/o temporal	TASED w/o spatial	TASED full
dog	0.594	0.700	0.718	0.668	0.734	0.749
gymnastics	0.707	0.716	0.730	0.691	0.701	0.702
parkour	0.578	0.677	0.746	0.658	0.756	0.779
skating	0.360	0.405	0.490	0.411	0.521	0.575
skiing	0.667	0.670	0.696	0.654	0.705	0.707
surfing	0.725	0.756	0.758	0.736	0.779	0.790
Average	0.651	0.664	0.712	0.667	0.702	0.730

the baseline methods, their results are copied from their original papers or borrowed from [1, 2].

Moreover, *PLD-VHD* improves the SOTA methods (*i.e.* SA [20] and MINI-Net [1]) on *TvSum* and *CoSum* by 3.1% and 9.9%, respectively, as shown in Tables 2 and 3, respectively. This clearly demonstrates the effectiveness of our approach by learning pixel-level distinctions. For limitations, our method mainly fails in first-person videos shot by hand-held cameras, especially for skating in YouTube Highlight, which contain a lot of cluttered background due to uncontrollable camera motions

We also present some visual examples for video highlight detection in different domains.

As it is shown in Figure 4 and the supplementary material, by learning pixel-level distinctions, our framework can perform video highlight detection effectively.

### 4.3. Ablation Studies

We conducted an additional experiment by changing the backbone network to C3D [34] pretrained on Sports1M [43], which is the same setting using in video2gif, and is simpler than that [36] used in MINI-Net. As is shown in

Table 5. Results of ablation studies on TvSum and CoSum.

	TASED w/o temporal	TASED w/o spatial	TASED full
TvSum	0.729	0.741	0.771
CoSum	0.888	0.915	0.946

Table 4, although the performance of Ours(C3D) is slightly worse than Ours(TASED) due to the use of a weaker backbone networks, it still outperforms all other existing methods. This again confirms that our proposed approach is effective even when different backbone networks are used.

We further validate the effectiveness of different components in our proposed approach. Specifically, two kinds of cues are used in our approach: the temporal cues and spatial cues. To validate their effects, we conduct ablation experiments by respectively removing those two types of cues as outlined below, with a total of four variants:

- **C3D w/o spatial** takes C3D [34] as the backbone network, but does not use the saliency mask to generate the pseudo-distinction labels. In other words, we use the distinction label defined in Eq. (1) in this case.

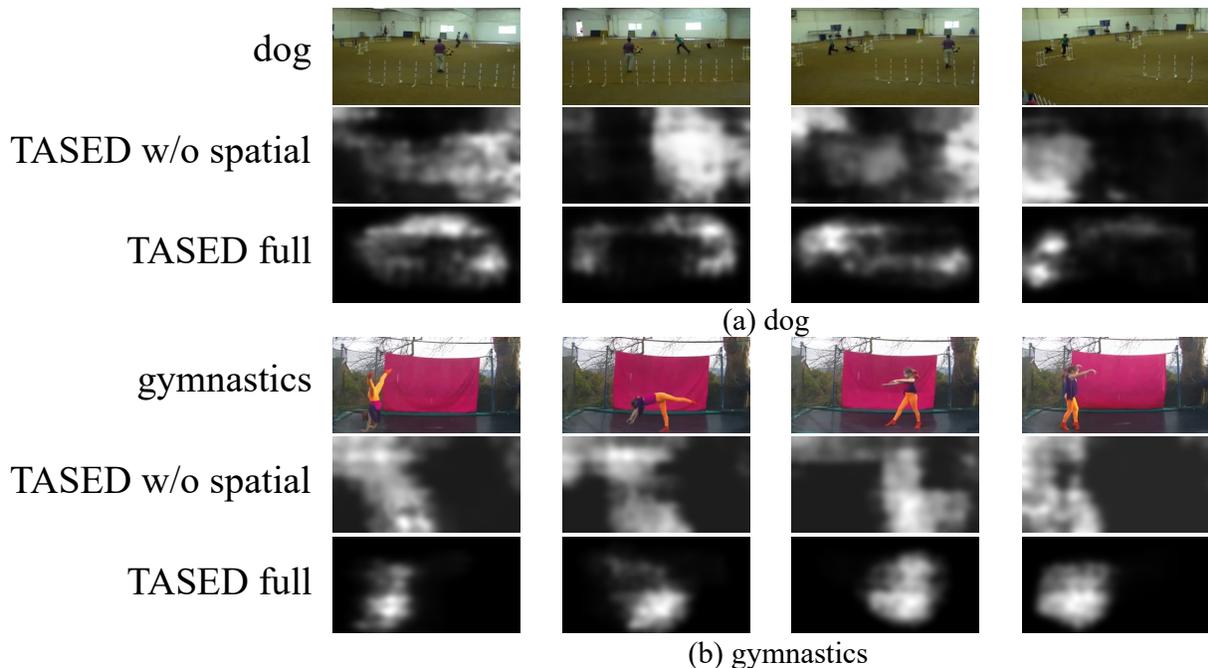


Figure 5. The frames on the first line of each subfigure are sampled from dog(gymnastics) in YouTube Highlight dataset. In the third line, the white regions inferred by our full model present the trajectory of the dog show and the action of the actress in the gymnastics clip, while the second line inferred without the spatial module may contain some background noise and cannot provide a clear highlight cues.

- **C3D w/o temporal** removes the effect of the temporal context by duplicating the target frame  $I_t$  to fill the video clip  $C_t$ . It can be determined using only  $I_t$  for distinction estimation as described in Eq. (4).
- **TASED w/o spatial** follows the same setting as *C3D w/o spatial*, but uses the TASED-Net as the backbone.
- **TASED w/o temporal** follows the same setting as *C3D w/o temporal*, but uses the TASED-Net as the backbone.

The results of different variants and the full models are shown in Table 4. We can observe from these results that both temporal and spatial cues are important. In particular, when spatial cues are removed, the performance of our model using C3D (TASED) drops from 0.712 to 0.664 (0.730 to 0.702). Benefiting from the spatial highlight cues, our model can be more robust to noise and exhibit improved learning of pixel-level distinction for video highlight detection, as shown in Figure 5. Similarly, when the temporal context is eliminated, our model using C3D (TASED) drops from 0.712 to 0.651 (0.730 to 0.667). Similar ablation studies of our temporal and spatial module on *TvSum* and *CoSum* are presented in Table 5.

More detailed results of each domain on *TvSum* and *CoSum* are included in the supplementary materials. This con-

firms our analysis suggesting that the video highlight task is highly dependent on the context preceding the current frame, both temporally and spatially.

## 5. Conclusion

In this work, we make pixel-level distinctions for video highlight detection by exploiting the temporal and spatial relations within video segments. For temporal relations, we utilize a 3D convolutional neural network to capture the distinctions by incorporating frames prior to the current frame while also making use of visual saliency to model the distinctions for spatial relations. We further adopt an encoder-decoder structure to predict pixel-level distinctions for highlight detection. In addition to achieving state-of-the-art performance, our proposed approach also has the advantage of explainability.

## Acknowledgements

This work is supported by the Major Project for New Generation of AI under Grant No. 2018AAA0100400, the National Natural Science Foundation of China (Grant No. 62176047), Beijing Natural Science Foundation (Z190023), and Alibaba Group through Alibaba Innovation Research Program.

## References

- [1] F.-T. Hong, X. Huang, W.-H. Li, and W.-S. Zheng, “Mininet: Multiple instance ranking network for video highlight detection,” in *ECCV*, 2020, pp. 345–360. [3](#), [6](#), [7](#)
- [2] B. Xiong, Y. Kalantidis, D. Ghadiyaram, and K. Grauman, “Less is more: Learning highlight detection from video duration,” in *CVPR*, 2019, pp. 1258–1267. [2](#), [6](#), [7](#)
- [3] M. Gygli, Y. Song, and L. Cao, “Video2gif: Automatic generation of animated gifs from video,” in *CVPR*, 2016, pp. 1001–1009. [1](#), [2](#), [3](#), [6](#)
- [4] T. Yao, T. Mei, and Y. Rui, “Highlight detection with pairwise deep ranking for first-person video summarization,” in *CVPR*, 2016, pp. 982–990. [1](#), [2](#)
- [5] Y. Zhang, J. Gao, X. Yang, C. Liu, Y. Li, and C. Xu, “Find objects and focus on highlights: Mining object semantics for video highlight detection via graph neural networks,” in *AAAI*, 2020, pp. 12902–12909. [2](#), [3](#), [6](#)
- [6] A. Garcia del Molino and M. Gygli, “Phd-gifs: personalized highlight detection for automatic gif creation,” in *ACM MM*, 2018, pp. 600–608.
- [7] M. Rochan, M. K. K. Reddy, L. Ye, and Y. Wang, “Adaptive video highlight detection by learning from user history,” in *ECCV*, 2020, pp. 261–278.
- [8] M. Sun, A. Farhadi, and S. Seitz, “Ranking domain-specific highlights by analyzing edited videos,” in *ECCV*, 2014, pp. 787–802. [2](#), [3](#), [6](#)
- [9] H. Yang, B. Wang, S. Lin, D. Wipf, M. Guo, and B. Guo, “Unsupervised extraction of video highlights via robust recurrent auto-encoders,” in *ICCV*, 2015, pp. 4633–4641. [6](#)
- [10] Y. Jiao, Z. Li, S. Huang, X. Yang, B. Liu, and T. Zhang, “Three-dimensional attention-based deep ranking model for video highlight detection,” *IEEE Transactions on Multimedia*, vol. 20, no. 10, pp. 2693–2705, 2018. [1](#), [2](#), [6](#)
- [11] J. Wang, C. Xu, E. Chng, and Q. Tian, “Sports highlight detection from keyword sequences using hmm,” in *ICME*, vol. 1, 2004, pp. 599–602. [2](#)
- [12] C. Xu, J. Wang, K. Wan, Y. Li, and L. Duan, “Live sports event detection based on broadcast video and web-casting text,” in *ACM MM*, 2006, pp. 221–230. [2](#)
- [13] C. Xu, J. Wang, H. Lu, and Y. Zhang, “A novel framework for semantic annotation and personalized retrieval of sports video,” *IEEE transactions on multimedia*, vol. 10, no. 3, pp. 421–436, 2008. [2](#)
- [14] F. Qi, X. Yang, and C. Xu, “Emotion knowledge driven video highlight detection,” *IEEE Transactions on Multimedia*, 2020.
- [15] Y. Song, J. Vallmitjana, A. Stent, and A. Jaimes, “Tvsum: Summarizing web videos using titles,” in *CVPR*, 2015, pp. 5179–5187. [2](#), [6](#)
- [16] W.-S. Chu, Y. Song, and A. Jaimes, “Video co-summarization: Video summarization by visual co-occurrence,” in *CVPR*, 2015, pp. 3584–3592. [2](#), [6](#)
- [17] G. Zhu, Q. Huang, C. Xu, L. Xing, W. Gao, and H. Yao, “Human behavior analysis for highlight ranking in broadcast racket sports video,” *IEEE Transactions on Multimedia*, vol. 9, no. 6, pp. 1167–1182, 2007. [2](#)
- [18] Z. Guo, Z. Zhao, W. Jin, W. Dazhou, L. Ruitao, and J. Yu, “Taohighlight: Commodity-aware multi-modal video highlight detection in e-commerce,” *IEEE Transactions on Multimedia*, 2021.
- [19] M. Xu, H. Wang, B. Ni, R. Zhu, Z. Sun, and C. Wang, “Cross-category video highlight detection via set-based learning,” in *ICCV*, 2021, pp. 7970–7979. [1](#), [2](#), [3](#), [6](#)
- [20] T. Badamdorj, M. Rochan, Y. Wang, and L. Cheng, “Joint visual and audio learning for video highlight detection,” in *ICCV*, 2021, pp. 8127–8137. [2](#), [3](#), [6](#), [7](#)
- [21] B. Gong, W.-L. Chao, K. Grauman, and F. Sha, “Diverse sequential subset selection for supervised video summarization,” *NIPS*, vol. 27, pp. 2069–2077, 2014. [3](#), [6](#)
- [22] K. Zhou, Y. Qiao, and T. Xiang, “Deep reinforcement learning for unsupervised video summarization with diversity-representativeness reward,” in *AAAI*, 2018. [3](#)
- [23] Y. J. Lee, J. Ghosh, and K. Grauman, “Discovering important people and objects for egocentric video summarization,” in *CVPR*, 2012, pp. 1346–1353. [3](#)
- [24] M. Gygli, H. Grabner, and L. Van Gool, “Video summarization by learning submodular mixtures of objectives,” in *CVPR*, 2015, pp. 3090–3098. [3](#), [6](#)
- [25] Z. Lu and K. Grauman, “Story-driven summarization for egocentric video,” in *CVPR*, 2013, pp. 2714–2721. [3](#)
- [26] B. Mahasseni, M. Lam, and S. Todorovic, “Unsupervised video summarization with adversarial lstm networks,” in *CVPR*, July 2017. [3](#), [6](#)
- [27] B. Zhao, X. Li, and X. Lu, “Hierarchical recurrent neural network for video summarization,” in *ACM MM*, 2017, pp. 863–871. [3](#)
- [28] —, “Hsa-rnn: Hierarchical structure-adaptive rnn for video summarization,” in *CVPR*, 2018, pp. 7405–7414. [3](#)
- [29] B. Zhao, H. Li, X. Lu, and X. Li, “Reconstructive sequence-graph network for video summarization,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [3](#)
- [30] K. Zhang, K. Grauman, and F. Sha, “Retrospective encoders for video summarization,” in *ECCV*, 2018, pp. 383–399. [3](#)
- [31] G. Yi, D. Yang, A. Bentalab, W. Li, Y. Li, K. Zheng, J. Liu, W. T. Ooi, and Y. Cui, “The acm multimedia 2019 live video streaming grand challenge,” in *ACM MM*, 2019, pp. 2622–2626. [3](#)
- [32] M. Otani, Y. Nakashima, E. Rahtu, and J. Heikkila, “Re-thinking the evaluation of video summaries,” in *CVPR*, 2019, pp. 7596–7604. [3](#)
- [33] A. B. Vasudevan, M. Gygli, A. Volokitin, and L. Van Gool, “Query-adaptive video summarization via quality-aware relevance estimation,” in *ACM MM*, 2017, pp. 582–590. [3](#)

- [34] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, “Learning spatiotemporal features with 3d convolutional networks,” in *ICCV*, 2015, pp. 4489–4497. 4, 7
- [35] J. Carreira and A. Zisserman, “Quo vadis, action recognition? a new model and the kinetics dataset,” in *CVPR*, 2017, pp. 6299–6308.
- [36] K. Hara, H. Kataoka, and Y. Satoh, “Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?” in *CVPR*, 2018, pp. 6546–6555. 7
- [37] S. Xie, C. Sun, J. Huang, Z. Tu, and K. Murphy, “Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification,” in *ECCV*, 2018, pp. 305–321.
- [38] K. Min and J. J. Corso, “Tased-net: Temporally-aggregating spatial encoder-decoder network for video saliency detection,” in *ICCV*, 2019, pp. 2394–2403. 3, 4, 6
- [39] J. Li, H. Zhang, W. Wan, and J. Sun, “Two-class 3d-cnn classifiers combination for video copy detection,” *Multimedia Tools and Applications*, vol. 79, no. 7, pp. 4749–4761, 2020.
- [40] Z. Chen, D. Huang, Y. Wang, and L. Chen, “Fast and light manifold cnn based 3d facial expression recognition across pose variations,” in *ACM MM*, 2018, pp. 229–238.
- [41] T. T. Niemi, M. Viitanen, and J. Vanne, “Binocular multi-cnn system for real-time 3d pose estimation,” in *ACM MM*, 2020, pp. 4553–4555.
- [42] Y. Hao, Z.-N. Liu, H. Zhang, B. Zhu, J. Chen, Y.-G. Jiang, and C.-W. Ngo, “Person-level action recognition in complex events via tsd-tsm networks,” in *ACM MM*, 2020, pp. 4699–4702.
- [43] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei, “Large-scale video classification with convolutional neural networks,” in *CVPR*, 2014, pp. 1725–1732. 7
- [44] L. Wu, Y. Yang, L. Chen, D. Lian, R. Hong, and M. Wang, “Learning to transfer graph embeddings for inductive graph based recommendation,” in *SIGIR*, 2020, pp. 1211–1220.
- [45] W. Wang, J. Shen, and L. Shao, “Consistent video saliency using local gradient flow optimization and global refinement,” *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 4185–4196, 2015. 3
- [46] C. Bak, A. Kocak, E. Erdem, and A. Erdem, “Spatio-temporal saliency networks for dynamic saliency prediction,” *IEEE Transactions on Multimedia*, vol. 20, no. 7, pp. 1688–1698, 2017. 3
- [47] L. Bazzani, H. Larochelle, and L. Torresani, “Recurrent mixture density network for spatiotemporal visual attention,” in *ICLR*, 2017. 3
- [48] W. Wang, J. Shen, F. Guo, M.-M. Cheng, and A. Borji, “Revisiting video saliency: A large-scale benchmark and a new model,” in *CVPR*, 2018, pp. 4894–4903. 3, 6
- [49] A. Tsiami, P. Koutras, and P. Maragos, “Stavis: Spatio-temporal audiovisual saliency network,” in *CVPR*, 2020, pp. 4766–4776. 3
- [50] E. Elhamifar, G. Sapiro, and R. Vidal, “See all by looking at a few: Sparse modeling for finding representative objects,” in *CVPR*, 2012, pp. 1600–1607. 6
- [51] G. Kim, L. Sigal, and E. P. Xing, “Joint summarization of large-scale collections of web images and videos for storyline reconstruction,” in *CVPR*, 2014, pp. 4225–4232. 6
- [52] R. Panda and A. K. Roy-Chowdhury, “Collaborative summarization of topic-related videos,” in *CVPR*, 2017, pp. 7083–7092. 6
- [53] S. Cai, W. Zuo, L. S. Davis, and L. Zhang, “Weakly-supervised video summarization using variational encoder-decoder and web prior,” in *ECCV*, 2018, pp. 184–200. 6
- [54] R. Panda, A. Das, Z. Wu, J. Ernst, and A. K. Roy-Chowdhury, “Weakly supervised summarization of web videos,” in *ICCV*, 2017, pp. 3657–3666. 6
- [55] D. Potapov, M. Douze, Z. Harchaoui, and C. Schmid, “Category-specific video summarization,” in *ECCV*, 2014, pp. 540–555. 6
- [56] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, “Video summarization with long short-term memory,” in *ECCV*, 2016, pp. 766–782. 6