

# Flow Guided Recurrent Neural Encoder for Video Salient Object Detection

Guanbin Li<sup>1</sup> Yuan Xie<sup>1</sup> Tianhao Wei<sup>2</sup> Keze Wang<sup>1</sup> Liang Lin<sup>1,3\*</sup>

<sup>1</sup>Sun Yat-sen University <sup>2</sup>Zhejiang University <sup>3</sup>SenseTime Group Limited

liguanbin@mail.sysu.edu.cn, xiey39@mail2.sysu.edu.cn, thwei@zju.edu.cn,

wangkeze@mail2.sysu.edu.cn, linliang@ieee.org

## Abstract

Image saliency detection has recently witnessed significant progress due to deep convolutional neural networks. However, extending state-of-the-art saliency detectors from image to video is challenging. The performance of salient object detection suffers from object or camera motion and the dramatic change of the appearance contrast in videos. In this paper, we present flow guided recurrent neural encoder (FGRNE), an accurate and end-to-end learning framework for video salient object detection. It works by enhancing the temporal coherence of the per-frame feature by exploiting both motion information in terms of optical flow and sequential feature evolution encoding in terms of LSTM networks. It can be considered as a universal framework to extend any FCN based static saliency detector to video salient object detection. Intensive experimental results verify the effectiveness of each part of FGRNE and confirm that our proposed method significantly outperforms state-of-the-art methods on the public benchmarks of DAVIS and FBMS.

## 1. Introduction

Salient object detection aims at identifying the most visually distinctive objects in an image or video that attract human attention. It has drawn a lot of attention due to the need for solving this problem in many computer vision applications such as image and video compression [12], object segmentation [37], visual tracking [38] and person re-identification [43]. Although image based salient object detection has been extensively studied during the past decade, video based salient object detection is much less explored

\*The first two authors contribute equally to this paper. Corresponding author is Liang Lin. This work was supported by the State Key Development Program under Grant 2016YFB1001004, the National Natural Science Foundation of China under Grant 61702565, Guangdong Natural Science Foundation Project for Research Teams under Grant 2017A030312006, and was also sponsored by CCF-Tencent Open Research Fund.

due to its high complexity and the lack of large-scale annotated video datasets.

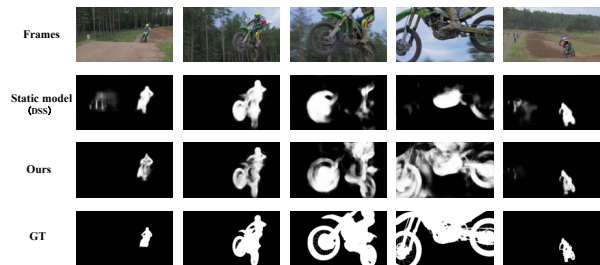


Figure 1. The challenges of still-image saliency detector and the effectiveness of temporal coherence modeling in video based salient object detection.

In recent years, due to the success deployment of deep convolutional neural networks (CNN), the performance of salient object detection in static image has been increased by a significant margin [21, 10, 18, 20]. Nevertheless, directly applying these methods to video salient object detection is non-trivial and challenging. The performance of salient object detection suffers from object or camera motion and the dramatic change of the appearance contrast in videos. As shown in the second row of Fig. 1, state-of-the-art still-image salient object detectors (e.g. DSS [10]) deteriorates drastically from the inability to maintain the visual continuity and temporal correlation of salient objects between consecutive frames.

Cognitive studies have shown that visual contrast is the key factor that leads to a specific region becoming salient in static images. For dynamic videos, the difference between consecutive frames caused by object motion are more attractive to people's attention [13]. Such temporal information has been exploited in existing video salient object detection methods either in the form of graphics model [35, 3] or simply embedded in a convolutional neural network [36]. Graphics model based methods generally employ generative framework which first infers an initial saliency map from either intra-frame appearance contrast

information [3] or inter-frame gradient flow field [35], and further incorporates an energy function with some heuristic spatio-temporal modeling to encourage the cross-frame consistency of the output saliency maps. Due to their independence from training data and the use of handcrafted low-level features, it is arduous for graphics model based methods to cope with videos with complex semantic contrast and objects motion. Although optical flow has been exploited in these methods, it is only used in an off-the-shelf mode for heuristic post-processing. Recently, with the thriving application of deep CNN in salient object detection of static images, there are also attempts to extend CNN to video salient object detection [36, 16]. They simply concatenate consecutive frame images and feed to convolutional neural networks for temporal coherence modeling. However, since convolutional neural network does not have the memory function, this naive aggregation of raw frame images followed by a series of convolution operations can not well characterize the continuous dynamic evolution of video frames in the temporal domain. Moreover, this simple spatio-temporal modeling strategy lacks explicit compensation for object's motion, making it hard to detect the salient objects with strenuous movement while maintaining temporal coherence (e.g. object moves beyond the receptive field of the neural network).

In this work, we present flow guided recurrent neural encoder (FGRNE), an end-to-end learning framework to extend any FCN based still-image saliency detectors to video salient object detection. It works by enhancing the temporal coherence of the per-frame feature by exploiting both motion information in terms of optical flow and sequential feature evolution encoding in terms of LSTM networks. Specifically, we employ an off-the-shelf FCN based image saliency detector (e.g. DSS [10]) as our host network for feature extraction and ultimate saliency inference, and a pre-trained FlowNet [7] for motion estimation between a frame pair. Our FGRNE learns to improve the per-frame feature by incorporating a flow guided feature warping followed by a LSTM based temporal coherence feature encoding. The output feature map at the last time-step is considered as our encoded feature and is fed to the upper part of the host network for saliency inference. Moreover, our FGRNE also involves another LSTM module to improve the estimated optical flow from frame pairs with large time interval. All the three modules of FGRNE including motion computing and updating, flow guided feature warping as well as temporal coherence feature encoding are trained end-to-end with the host network.

In summary, this paper has the following contributions:

- We introduce a flow guided recurrent neural encoder framework to enhance the temporal coherence modeling of the per-frame feature representation, which can be exploited to extend any FCN based still-image

saliency detector to video salient object detection.

- We propose to incorporate an optical flow network in FGRNE framework to estimate the motion of each frame, which is further used in feature warping to explicitly compensate for object's motion.
- We proposed to exploit a ConvLSTM in our FGRNE for sequential feature encoding, which can capture the evolution of appearance contrast in temporal domain and is complementary to feature warping towards an improved performance for video salient object detection.

## 2. Related Work

### 2.1. Still-Image Salient Object Detection

Image salient object detection has been extensively studied for decades. Conventional methods can be divided into bottom-up approaches based on low-level features [8, 15, 5] and top-down models guided by high-level knowledge [14, 40, 22]. In recent years, the profound deep CNN has pushed the research on salient object detection into a new phase and has become the dominant research direction in this field. Deep CNN based methods can be further divided into two categories, including region based deep feature learning [19, 42, 32] and end-to-end fully convolutional network based methods [20, 10, 18, 33, 17]. Methods in the first category separate an image into regions, and treat each region as an independent unit for deep feature extraction and saliency inference. They are generally space and time consuming due to significant redundancy in feature extraction and storage. To overcome this deficiency, deep FCN based models have been developed to directly map a raw input image to its corresponding saliency map in an end-to-end trainable way. These kind of methods can make the best of feature sharing mechanism and generate the hierarchical feature of each region in a single network forward operation. They can produce superior saliency maps and have become the fundamental component in state-of-the-art methods of this field.

In contrast to these still-image based salient object detection methods, we focus on video salient object detection, which incorporates both temporal and motion information to improve the feature map representation for saliency inference. It can be considered as a universal framework to extend any FCN based models to video salient object detection, and can easily benefit from the improvement of still-image salient object detectors.

### 2.2. Video Salient Object Detection

Compared with saliency detection in still images, detecting video salient objects is much more challenging due to the high complexity in effective spatio-temporal modeling

and the lack of large-scale annotated video datasets. It is far less explored in the research community. Earlier methods to this problem can be considered as simple extensions of some static saliency models with extra crafted temporal features [24, 9]. More recent and noteworthy works generally formulate video saliency detection as a spatio-temporal context modeling problem over consecutive frames, and incorporates energy functions with handcrafted rules to encourage both the spatial smoothness and temporal consistency of the output saliency maps [3, 35, 6]. However, these approaches all belong to unsupervised generative models and depend on handcrafted low-level features for heuristic saliency inference, and hence are not able to handle complex videos that require knowledge and semantic reasoning. Though recently an unpublished work by *Le et al.* [16] proposes to incorporate deep CNN feature in a spatio-temporal CRF framework for temporal consistency enhancement, it still suffers from the deficiency of multi-stage pipeline and its high-computational costs. The most relevant work with us is [33], which exploits a second FCN to improve the temporal coherence of the saliency map generated from an initial static FCN based saliency network, by taking as input the concatenation of successive frame pair as well as the initial saliency map and directly mapping to a refined saliency map in a forward network operation. Since convolutional neural network does not have the memory function, it is not able to well model the continuous evolution of video frames in the temporal domain. Moreover, this rough strategy of spatio-temporal modeling lacks explicit compensation for objects motion, making it hard to detect the salient objects with strenuous movement.

By contrast, our method considers temporal information in the feature level instead of the raw input frames and incorporates LSTM network to naturally encode the sequential feature evolution. The entire framework is trained end-to-end and the inference process is highly efficient. Besides, our method can further incorporate such graphics model based post-processing technique (e.g. CRF) to improve the performance.

### 2.3. Optical Flow based Motion Estimation

Optical flow estimates the per-pixel motion between two consecutive frames and it is widely used in a variety of video analysis tasks. Traditional methods are mostly based on variational formulation, which mainly tackle small displacements and are limited by their high-computational costs for efficient video applications. Recently, deep learning based methods have been employed to optical flow computation [7, 28, 11]. The most representative work is FlowNet [7] which shows that CNN can be applied to highly effective optical flow inference. There are also attempts to incorporate FlowNet in contemporary deep learning framework to enhance the temporal continuity of the represen-

tation of video features, which has brought performance improvements to various of video comprehension tasks, including video recognition [45], object detection [44] and video object segmentation [29].

Optical flow has been exploited in existing video salient object detection models, however, it is either used as an auxiliary motion feature or a handcrafted rule in post-processing for temporal coherence improvement. Inspired by [45, 44], we incorporate optical flow to enable feature warping across frames and compensate for the change caused by object motion. However, unlike these efforts, the motion flow is dynamically updated in our framework and the result of feature warping is exploited to temporal feature encoding instead of feature aggregation. Moreover, we are first to integrate optical flow in recurrent neural encoder for effective spatio-temporal feature learning and have demonstrated its superior performance on the task of video salient object detection.

### 3. Flow Guided Recurrent Neural Encoder

Given the a video frame sequence  $I_i, i = 1, 2, \dots, N$ , the objective of video salient object detection is to output the saliency maps of all frames,  $S_i, i = 1, 2, \dots, N$ . State-of-the-art salient object detectors for static image are mostly based on FCN structure [20, 23, 18, 10]. Given a pre-trained static model  $\mathcal{N}$  (e.g. DSS [10] model), it can be considered as a feature extraction module  $\mathcal{N}_{\text{fea}}$  followed by a pixel-wise saliency regression module  $\mathcal{N}_{\text{reg}}$ . The output saliency map  $S$  of a given image  $I$  can be computed as  $S = \mathcal{N}_{\text{reg}}(\mathcal{N}_{\text{fea}}(I))$ . Directly applying this model to each individual frame usually generates unstable and temporally inconsistent saliency maps due to the lack of temporal coherence modeling in feature representation.

Our proposed FGRNE  $\mathcal{E}$  aims at enhancing the temporal consistency of feature representation by extra looking at a segment of  $k$  former frames. Given a reference frame  $I_i$ , the encoded feature is denoted as  $F_i = \mathcal{E}(\mathcal{N}_{\text{fea}}(I_i), \mathcal{N}_{\text{fea}}(I_{i-1}), \dots, \mathcal{N}_{\text{fea}}(I_{i-k}))$ . As object motion and the change of its appearance contrast are two core influencing factors to video saliency, the proposed FGRNE incorporates an off-the-shelf FlowNet model [7] and a LSTM based feature encoder to respectively take care of these two factors.

As shown in Fig. 2, the architecture of our FGRNE consists of three modules, including motion computing and updating, motion guided feature warping, and temporal coherence feature encoding. Specifically, we first compute an optical flow map for each of the  $k$  former frames relative to the reference frame. Each of the flow map is further fed to a LSTM in reverse order for motion refinement. Secondly, the updated flow map at each time step is applied to warp the feature map accordingly. And finally, each warped feature is consecutively fed to another LSTM for temporal

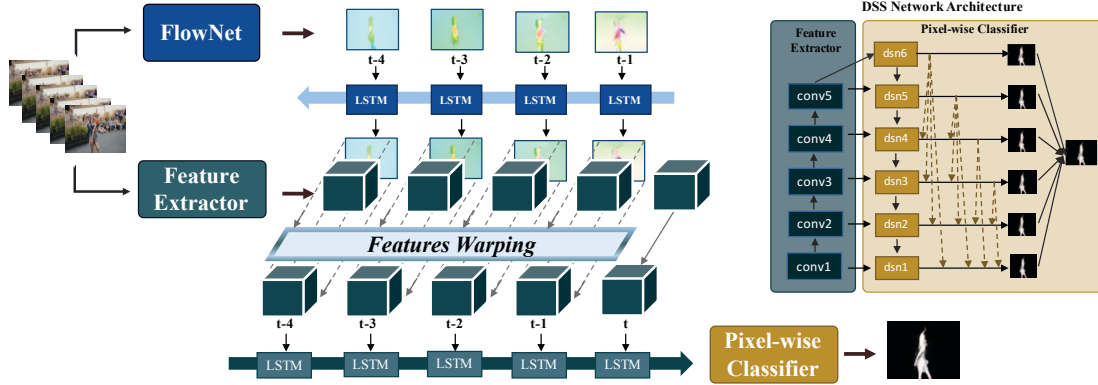


Figure 2. Our overall framework for flow guided recurrent neural encoder. It incorporates a LSTM with reverse sequential input for motion flow update, a flow guided feature warping module and another LSTM for temporal coherence feature encoding.

coherence feature encoding, which produces the resulted feature  $F_i$ . The output saliency map is thus computed as  $S_i = \mathcal{N}_{\text{reg}}(F_i)$ .

### 3.1. Motion Computing and Updating

Given a reference  $I_i$  and a window of  $k$  former frames, we first apply the embedded FlowNet  $\mathcal{F}$  [7] to individually estimate  $k$  initial flow fields  $\{O_{i \rightarrow j} = \mathcal{F}(I_i, I_j) | j = i-1, i-2, \dots, i-k\}$  relative to the reference frame. The resulted flow field  $O_{i \rightarrow j}$  is a position offset map of two channels. It computes the pixel displacement  $(u, v)$  for every pixel location  $(x, y)$  in  $I_i$  to the spatial location  $(x', y')$  in  $I_j$ , i.e.,  $(x', y') = (x + u, y + v)$ , where  $u$  and  $v$  respectively represent the pixel offsets in horizontal and vertical directions.

As the FlowNet is originally trained from pair data of consecutive frames, it may not be accurate enough to reflect the motion relationship between two frames with long time interval. Intuitively, the closer to the reference frame, the more accurate the estimated motion flow. We can gradually employ flow maps of closer frames to refine that of larger time interval. Based on the above consideration, we propose to combine a ConvLSTM [39] with CNN based FlowNet to jointly learn the flow map and refine in reverse order.

ConvLSTM is an extension of traditional fully connected LSTM which has convolutional structures in both input-to-state and state-to-state connections. All of the data transferred in ConvLSTM can be regarded as 3D tensors with the last two dimensions being spatial dimensions. Let  $\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_t$  denote the input to ConvLSTM and  $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_t$  stand for its hidden states. At each time step, the output hidden state of ConvLSTM is updated based on its own input as well as the encoded past states from its previous input, which is formulated as

$$\mathcal{H}_t = \text{ConvLSTM}(\mathcal{H}_{t-1}, \mathcal{C}_{t-1}, \mathcal{X}_t), \quad (1)$$

where  $\mathcal{C}$  is the memorized cell state of the ConvLSTM at its previous time-step. Following [39], the ConvLSTM module consists of the input gate  $i_t$ , forget gate  $f_t$  and output gate  $o_t$ , the overall updating equations can be listed in (2), where ‘ $*$ ’ denotes the convolution operator, ‘ $\circ$ ’ denotes the Hadamard product, and  $\sigma(\cdot)$  stands for the sigmoid function:

$$\begin{aligned} i_t &= \sigma(W_{xi} * \mathcal{X}_t + W_{hi} * \mathcal{H}_{t-1} + W_{ci} \circ \mathcal{C}_{t-1} + b_i) \\ f_t &= \sigma(W_{xf} * \mathcal{X}_t + W_{hf} * \mathcal{H}_{t-1} + W_{cf} \circ \mathcal{C}_{t-1} + b_f) \\ C_t &= f_t \circ \mathcal{C}_{t-1} + i_t \circ \tanh(W_{xc} * \mathcal{X}_t + W_{hc} * \mathcal{H}_{t-1} + b_c) \\ o_t &= \sigma(W_{xo} * \mathcal{X}_t + W_{ho} * \mathcal{H}_{t-1} + W_{co} \circ \mathcal{C}_{t-1} + b_o) \\ \mathcal{H}_t &= o_t \tanh(C_t) \end{aligned} \quad (2)$$

To update the optical flow field with ConvLSTM, the LSTM layer is unrolled for a window of  $k$  flow fields and the size of the hidden state is set to be the same as the input flow map. We sequentially feed the  $k$  initial motion flow to the ConvLSTM cells in reverse order, i.e.,  $\mathcal{X}_{1:k} = O_{i \rightarrow (i-1)}, O_{i \rightarrow (i-2)}, \dots, O_{i \rightarrow (i-k)}$ . The hidden states are the encoding of the updated flow field, which is further fed to a convolutional layer with kernel size of  $1 \times 1$  to produce the refined flow map  $RO_{i \rightarrow j}$ , formulated as:

$$\begin{aligned} j &= i - t \\ \mathcal{H}_t &= \text{ConvLSTM}(\mathcal{H}_{t-1}, \mathcal{C}_{t-1}, O_{i \rightarrow j}) \\ RO_{i \rightarrow j} &= \text{Conv}_{1 \times 1}(\mathcal{H}_t) \end{aligned} \quad (3)$$

### 3.2. Motion Guided Feature Warping

As motivated by [45], given a refined flow map  $RO_{i \rightarrow j}$ , the feature maps  $\mathcal{N}_{\text{fea}}(I_j)$  on the  $j^{\text{th}}$  frame are warped to the reference frame by applying the following warping function,

$$\text{WarpF}_{i \rightarrow j} = \mathcal{W}(\mathcal{N}_{\text{fea}}(I_j), RO_{i \rightarrow j}) \quad (4)$$

where  $\text{WarpF}_{i \rightarrow j}$  refers to the feature maps warped from frame  $j$  to frame  $i$ .  $\mathcal{W}(\cdot)$  is the bilinear warping function, which is applied on all the spatial locations for each channel of the feature maps. It is implemented as a bilinear interpolation of  $\mathcal{N}_{\text{fea}}(I_j)$  at the desired positions w.r.t the optical flow  $RO_{i \rightarrow j}$ .

### 3.3. Temporal Coherence Feature Encoding

Although feature warping operation can compensate for the misalignment of features caused by object or camera motion. It is still not enough to characterize the continuous dynamic evolution of video frames as well as the evolution of appearance contrast in temporal domain. Base on the above considerations, we proposed to exploit another ConvLSTM for sequential feature encoding. Specifically, this ConvLSTM takes a sequence of warped features (including the feature of the reference frame) as input, i.e.,  $\mathcal{X}_{1:k}$  denoted in Equation (1) is set to  $\mathcal{X}_{1:k} = \text{WarpF}_{i \rightarrow (i-k)}, \text{WarpF}_{i \rightarrow (i-k+1)}, \dots, \text{WarpF}_{i \rightarrow (i-1)}, \mathcal{N}_{\text{fea}}(I_i)$ , and works by computing the forward hidden sequence for temporal feature encoding from  $t = 1$  to  $t = k + 1$ , and then updating the output layer. The state updating function in (1) can be rewritten as follows:

$$\begin{aligned} j &= i - k + (t - 1) \\ \mathcal{H}_t &= \text{ConvLSTM}(\mathcal{H}_{t-1}, \mathcal{C}_{t-1}, \text{WarpF}_{i \rightarrow j}), t \leq k \\ \mathcal{H}_{k+1} &= \text{ConvLSTM}(\mathcal{H}_k, \mathcal{C}_k, \mathcal{N}_{\text{fea}}(I_i)) \end{aligned} \quad (5)$$

The hidden states are the encoding of the memorized future till now. And the hidden state of the last time-step  $k + 1$  is our final feature encoding.

## 4. Experimental Results

### 4.1. Experimental Setup

#### 4.1.1 Datasets

We evaluate the performance of our method on two public datasets: Freiburg-Berkeley Motion Segmentation (FBMS) dataset [2, 25], and DAVIS [27] dataset. The FBMS dataset contains 59 videos with 720 annotated sparsely annotated frames. DAVIS is a newly developed dataset for video object segmentation, which contains 50 high quality and full HD video sequences with 3,455 densely annotated pixel-level and per-frame ground-truth. It is one of the most challenging benchmark which covers various video object segmentation challenges such as occlusions, motion blur and appearance changes.

There exists another dataset SegTrack V2, which is an extended dataset from the original SegTrack dataset proposed in [30] and contains 14 videos about bird, animal, car and human with 1,066 densely annotated frame images. As referred to [36], we combine the whole SegTrackV2, the training sets of FBMS and DAVIS as our training set, and

evaluate our trained model on the testing sets of DAVIS and FBMS.

### 4.1.2 Evaluation Criteria

Similar to image-based salient object detection, we adopt precision-recall curves (PR), maximum F-measure and mean absolute error (MAE) as the evaluation metrics. The continuous saliency map is rescaled to  $[0, 255]$  and is binarized using all integer thresholds in the interval. At each threshold value, a pair of precision and recall value can be obtained by comparing the binary saliency map against the groundtruth. The PR curve is obtained from the average precision and recall over saliency maps of all images in the dataset. The F-measure is defined as

$$F_\beta = \frac{(1 + \beta^2) \cdot \text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}, \quad (6)$$

where  $\beta^2$  is set to 0.3 as suggested in [1]. We report the maximum F-measure (maxF) computed from the PR curve. MAE is defined as the average pixelwise absolute difference between the binary ground truth  $G$  and the saliency map  $S$  [26],

### 4.1.3 Implementation Details

Our proposed FRGNE has been implemented on the Mxnet [4], a flexible open source deep learning framework. FRGNE is compatible with any FCN based still-image salient object detectors. In this paper, we choose the state-of-the-art deeply supervised salient object detection (DSS) [10] method with public trained model as a baseline and take the updated DSS with FRGNE embedded as our final model for video salient object detection when performing ablation study and compared with other benchmarks. In Section 4.3, we will list more results of our proposed FRGNE on other host networks to demonstrate the effectiveness of our proposed algorithm. During training, the frame images are resized to  $256 \times 256$  before feeding into the network. While inference, we resize the image to a shorter side of 256 pixels. We train all the components incorporated in our framework in an end-to-end mode using SGD with a momentum of 0.9. The learning rate is initially set to  $2.5 \times 10^{-4}$  and decayed by 0.9 at every 8k training round. The loss function is set as same as the host network (e.g. DSS [10] employs an image-level class-balanced cross-entropy loss). The window size  $k$  is limited by the memory, with its default value set to 5 in our experiment. We have also explored the impact of different settings in Section 4.3. Experiments are performed on a workstation with an NVIDIA Titan X GPU and a 3.4GHz Intel processor.

DATASET	Metric	MST	MB+	RFCN	DHSNet	DCL	DSS	SAG	GF	DLVSD	FGRNE
DAVIS	maxF	0.455	0.520	0.732	<b>0.778</b>	0.740	<b>0.775</b>	0.528	0.628	0.699	<b>0.798</b>
	MAE	0.165	0.183	<b>0.047</b>	<b>0.035</b>	0.061	<b>0.047</b>	0.080	0.067	0.064	<b>0.032</b>
FBMS	maxF	0.540	0.525	0.741	<b>0.744</b>	0.740	<b>0.760</b>	0.572	0.607	0.696	<b>0.783</b>
	MAE	0.179	0.204	0.089	<b>0.076</b>	0.133	<b>0.077</b>	0.145	0.101	<b>0.077</b>	<b>0.063</b>

Table 1. Comparison of quantitative results including maximum F-measure (larger is better) and MAE (smaller is better). The best three results are shown in red, blue, and green, respectively.

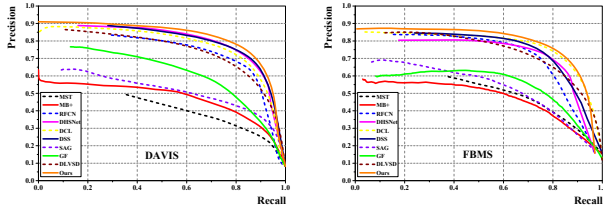


Figure 3. Comparison of precision-recall curves of 10 saliency detection methods on DAVIS and FBMS. Our proposed FGRNE consistently outperform other methods across the two testing dataset.

## 4.2. Comparison with the State of the Art

We compare our method (FGRNE) against 9 recent state-of-the-art methods, including MST [31], MB+ [41], RFCN [33], DHSNet [23], DCL [20], DSS [10], SAG [34], GF [35] and DLVSD [36]. The first six are the latest state-of-the-art salient object detection methods for static images while the last three are video-based saliency models. For fair comparison, we use either the implementations or the saliency maps provided by the authors. We also fine-tune all the public static saliency models using the training set as same as we train our FGRNE, and use the refined model for comparison.

A visual comparison is illustrated in Fig. 4. As can be seen, deep learning based static saliency models can generate seemingly promising saliency maps when watched independently, they are unsurprisingly inconsistent when putting in a whole sequence. Though existing video-based models can produce consistent results on videos with relatively slight object motions, they still can not handle videos with dramatic changes in appearance (object or camera motion). It is particularly noteworthy that our proposed method incorporates the off-the-shelf DSS [10] model as our baseline, it can learn to improve the original feature with temporal coherence and eventually produce optimized results far better than the original ones. In general, our method generates much more accurate and consistent saliency maps in various challenging cases.

As a part of quantitative evaluation, we show a comparison of PR curves in Fig. 3. As shown in the figures, our method (FGRNE) significantly outperforms all state-of-the-art static and dynamic salient object detection algorithms on both DAVIS and FBMS. Moreover, a quantitative comparison of maximum F-measure and MAE is listed in Table. 1,

our proposed method improves the maximum F-measure achieved by the best-performing static algorithm by 5.24% and 2.57% respectively on FBMS and DAVIS, and lowers the MAE by 17.10% and 8.57% accordingly. When compared with the best-performing video-based model, our FGRNE improves the maximum F-measure by 12.50% and 14.16% respectively on the FBMS and DAVIS dataset, and lowers the MAE by 18.18% and 50% accordingly. An interesting phenomenon is that the current best static saliency model actually outperforms the state-to-state video-based salient object detection methods because of the outstanding fully convolutional network.

Methods	$S_a$	$S_b$	$S_c$	$S_d$	$S_e$	$S_f$
feature aggregation?		✓		✓		
flow guided feature warping?				✓	✓	✓
flow update with LSTM?						✓
feature encoding with LSTM?			✓		✓	✓
maxF	0.775	0.768	0.777	0.780	0.793	<b>0.798</b>
MAE	0.047	0.052	0.036	0.036	0.035	<b>0.032</b>
runtime(ms)	97	112	137	162	184	191

Table 2. Effectiveness of flow guided recurrent neural encoder.

## 4.3. Ablation Studies

### 4.3.1 Effectiveness of Flow Guided Recurrent Neural Encoder

As discussed in Section 3, our proposed FGRNE involves three major modules, including motion flow updating, motion guided feature warping and temporal coherence feature encoding. To validate the effectiveness and necessity of each of these three modules, we compare FGRNE with its five variants in Table. 2.

$S_a$  refers to the saliency maps generated from the single-frame baseline model. To facilitate comparison, we also finetune the model using the individual frames of our used training set. It reaches the max  $F_\beta = 0.775$  and MAE = 0.047 in the test set of DAVIS, which already outperforms most of the state-of-the-art methods. This indicates that the fine-tuned baseline model is competitive and serves as a valid reference for evaluation. Compared to our entire framework, it is shown that embedding FGRNE to the baseline model totally leads to a 2.97% F-measure increase while reducing the MAE by 31.91%.

$S_b$  refers to a naive feature aggregation algorithm on the baseline model. The feature of the reference frame is simply updated as the weighted sum of the feature maps in

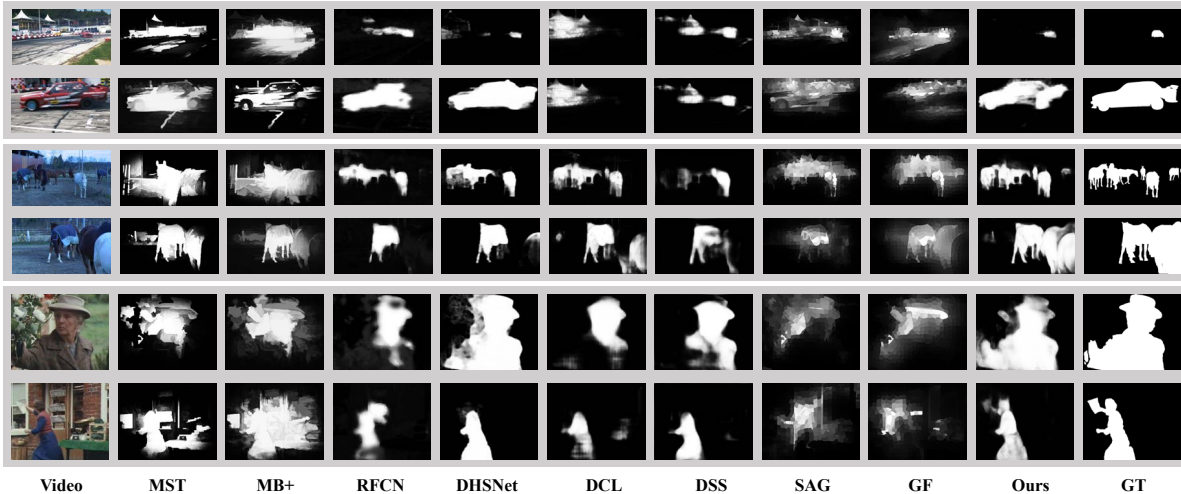


Figure 4. Visual comparison of saliency maps generated from state-of-the-art methods, including our FGRNE. The ground truth (GT) is shown in the last column. Our model consistently produces saliency maps closest to the ground truth.

the watching window, with the weight of  $j^{th}$  frame  $w_{i \rightarrow j}$  set to  $\frac{1}{i-j+1}$ . It is denoted as  $F_i = \sum_{l=0}^k \frac{1}{l+1} \mathcal{N}_{fea}(F_{i-l})$ . It is also trained end-to-end in the same way as we train our FGRNE. As shown in the table, the F-measure of this variant drops to 0.768 while the MAE increases to 0.052, which is even shy of the baseline model. It suggests that this naive feature aggregation is not suitable for sequential feature modeling. We speculate that the reason lies in the misalignment of features caused by changes in scene structure and appearance.

$S_c$  refers to a simple feature encoding algorithm on the baseline model, and a degenerate variant of FGRNE. The motion updating module is turned off and no flow motion is used, i.e., the motion flow  $O_{i \rightarrow j}$  is set to all zeros during training. The variant is also trained end-to-end in the same way as FGRNE. As shown in the table, the F-measure obtains a very slight increase to 0.777, while MAE greatly decreases by 23.40% to 0.036. However, the performance is still much inferior to the proposed FGRNE. This indicates that recurrent neural encoder can learn to exploit feature of previous frames to improve the temporal coherence of the reference frame. However, LSTM based feature encoding alone is not enough.

$S_d$  adds motion guided feature warping to the model of  $S_b$ , without the motion evolution update module turned on. It is actually a flow guided feature aggregation program. It increases the F-measure by 1.56% to 0.780 while lowers the MAE by 30.77% to 0.036 w.r.t the performance of  $S_b$ . It implies that feature alignment is an important operation before feature aggregation. The evident performance gain towards that of  $S_a$  also reveals the importance of motion modeling for video salient object detection.

$S_e$  adds motion guided feature warping to the model of

$S_c$ . It is a degenerate version of FGRNE without motion flow updating. All the other factors remain the same. It increases the max F-measure by 2.06% to 0.793 and lowers the MAE by 2.78% to 0.035 w.r.t the performance of  $S_c$ , which implies that the performance gain of motion guided feature warping is complementary to the LSTM based temporal coherence modeling. In fact, both object motion and the change of its appearance contrast are two core influencing factors to video saliency, which correspond exactly to the design of two complementary modules in our proposed FGRNE.

$S_f$  refers to the proposed FGRNE method, which turns on the motion flow evolution update module in  $S_e$ . It further brings a 0.63% boost to the F-measure to 0.798 while reducing the MAE by 8.57% to 0.032. This demonstrates the reverse LSTM can help to refine the motion flow, which makes up for the lack of FlowNet in estimating optical flow for frame pairs with large time interval.

Moreover, we have also listed the comparison on runtime cost of each variant of our proposed FGRNE. As shown in the figure, incorporating FGRNE to a static model cost an extra 94ms per-frame. Noted that the feature extraction are shared during the saliency inference of all the frames in a given window and our algorithm runs in a sliding window mode. Therefore, enlarging the window size does not contribute to severe increase of time computation cost.

#### 4.3.2 Sensitivities to Feature Extractor Selection

As described in Section 3, our FGRNE is dependent on a pre-trained static saliency detector as our host network. The host network is split into a feature extractor and a pixel-wise classification module. In principle, it can be split at any

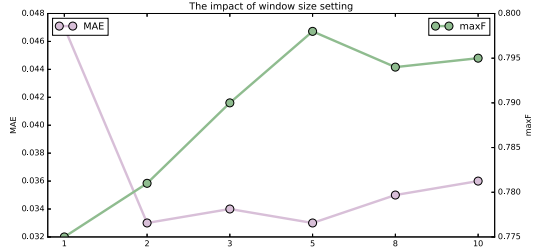


Figure 5. Sensitivity analysis on different window size settings.

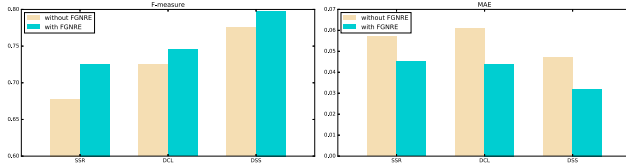


Figure 6. Sensitivity to host model selection.

layer as the host network is fully convolutional. We explore the effect of adding FGRNE to different levels of feature extraction on the performance of final results. We respectively experiment on adding feature encoding to the output feature map of Conv3\_3, Conv4\_3 and Conv5\_3 of the host DSS model. Experimental results shown that FGRNE is able to improve the temporal coherence on all scales of feature maps, which results in maxF value of 0.777, 0.789 and 0.798 respectively when choosing feature maps of Conv3\_3, Conv4\_3 and Conv5\_3. Among them, incorporating FGRNE with the feature extracted from Conv5\_3 results in the maximum performance gain, which increases the F-measure by 2.97% and decreases the MAE by 8.57% w.r.t to its single-frame static version.

### 4.3.3 Sensitivity to Window Size Setting

Our proposed FGRNE learns to facilitate the temporal coherence of the encoded feature by exploiting a window  $k$  former frames. Limited by the memory of our workstation,  $k$  can be set to a maximum of 10. We have explored the impact of different settings of  $k = \{1, 2, 3, 5, 8, 10\}$  on the salient object detection performance. Results in Fig. 5 show that training with 5 and 8 former frames achieves very close accuracy, with  $k = 5$  performing slightly better. By default, we set  $k = 5$  during training and inference in our experiments.

### 4.3.4 Sensitivity to Host Model Selection

As described in Section 3, we adopt a FCN based static saliency detector as the host model for our FGRNE. To demonstrate that our proposed method is widely applicable to any FCN based host network model, we apply to incorporate our FGRNE in two other recently published FCN based salient object detection methods, including DCL [20]

and MSRNet [18]. For the latter, due to the limitation of machine memory, we only experiment on its single scale version, i.e. SSRNet. As shown in Fig. 6, experimental evaluation on both F-measure and MAE have shown that our FGRNE can be trained to effectively enhance the spatio-temporal coherence of the feature representation, which greatly boost the performance of video salient object detection.

DATASET	LVO	LVO+CRF	FSEG	LMP	SFL	OUS	OUS+CRF
DAVIS	70.9	75.9	70.7	70.0	67.4	73.0	77.1
FBMS	63.5	65.1	68.4	35.7	55.0	72.4	76.2

Table 3. Performance comparison on unsupervised video object segmentation in terms of mean IoU

## 5. Comparison with Unsupervised Video Object Segmentation Methods

The problem setting of video salient object detection is very similar to that of unsupervised video object segmentation, except that its goal is to calculate a saliency probability value for each pixel instead of a binary classification. To make a fair comparison with the state-of-the-art unsupervised video object segmentation methods, we incorporate our FGRNE with a static ResNet-101 based pixel-wise binary classification model with feature extracted from the final output feature map of Conv5\_x. We evaluate our proposed method on both the DAVIS and FBMS datasets in terms of mean IoU and make a comparison with some state-of-the-art methods. As shown in Table 3, our proposed method outperforms LVO [29], the previous state of the art, by 2.96% and 14.0% on the IoU measure respectively on DAVIS and FBMS. Noted that as described in [29], the mIoU value of 75.9% reported on the leaderboard of DAVIS includes CRF as post-processing, the result of LVO without CRF is 70.9 as reported in their paper. For fair comparison, we also report our mIoU results with and without CRF in the table. As can be seen, our proposed method with CRF also greatly outperforms LVO by 1.6% and 16.90% respectively on DAVIS and FBMS.

## 6. Conclusion

In this paper, we have presented an accurate and end-to-end framework for video salient object detection. Our proposed flow guided recurrent encoder aims at improving the temporal coherence of the deep feature representation. It can be considered as a universal framework to extend any FCN based static saliency detector to video salient object detection, and can easily benefit from the future improvement of image based salient object detection methods. Moreover, as we focus on the learning an enhanced feature encoding, it can be easily extended to other applications of video analysis and it is worth exploring in the future.



## References

- [1] R. Achanta, S. Hemami, F. Estrada, and S. Susstrunk. Frequency-tuned salient region detection. In *CVPR*, pages 1597–1604. IEEE, 2009. 5
- [2] T. Brox and J. Malik. Object segmentation by long term analysis of point trajectories. *ECCV*, pages 282–295, 2010. 5
- [3] C. Chen, S. Li, Y. Wang, H. Qin, and A. Hao. Video saliency detection via spatial-temporal fusion and low-rank coherency diffusion. *TIP*, 26(7):3156–3170, 2017. 1, 2, 3
- [4] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. *arXiv preprint arXiv:1512.01274*, 2015. 5
- [5] M.-M. Cheng, N. J. Mitra, X. Huang, P. H. Torr, and S.-M. Hu. Global contrast based salient region detection. *TPAMI*, 37(3):569–582, 2015. 2
- [6] Y. Fang, Z. Wang, W. Lin, and Z. Fang. Video saliency incorporating spatiotemporal cues and uncertainty weighting. *TIP*, 23(9):3910–3921, 2014. 3
- [7] P. Fischer, A. Dosovitskiy, E. Ilg, P. Häusser, C. Hazırbaş, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. *arXiv preprint arXiv:1504.06852*, 2015. 2, 3, 4
- [8] D. Gao and N. Vasconcelos. Bottom-up saliency is a discriminant process. In *ICCV*, pages 1–6. IEEE, 2007. 2
- [9] C. Guo, Q. Ma, and L. Zhang. Spatio-temporal saliency detection using phase spectrum of quaternion fourier transform. In *CVPR*, pages 1–8. IEEE, 2008. 3
- [10] Q. Hou, M.-M. Cheng, X.-W. Hu, A. Borji, Z. Tu, and P. Torr. Deeply supervised salient object detection with short connections. *arXiv preprint arXiv:1611.04849*, 2016. 1, 2, 3, 5, 6
- [11] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. *arXiv preprint arXiv:1612.01925*, 2016. 3
- [12] L. Itti. Automatic foveation for video compression using a neurobiological model of visual attention. *TIP*, 13(10):1304–1318, 2004. 1
- [13] L. Itti, C. Koch, and E. Niebur. A model of saliency-based visual attention for rapid scene analysis. *TPAMI*, 20(11):1254–1259, 1998. 1
- [14] Y. Jia and M. Han. Category-independent object-level saliency detection. In *ICCV*, pages 1761–1768, 2013. 2
- [15] H. Jiang, J. Wang, Z. Yuan, Y. Wu, N. Zheng, and S. Li. Salient object detection: A discriminative regional feature integration approach. In *CVPR*, pages 2083–2090, 2013. 2
- [16] T.-N. Le and A. Sugimoto. Video salient object detection using spatiotemporal deep features. *arXiv preprint arXiv:1708.01447*, 2017. 2, 3
- [17] G. Lee, Y.-W. Tai, and J. Kim. Deep saliency with encoded low level distance map and high level features. In *CVPR*, pages 660–668, 2016. 2
- [18] G. Li, Y. Xie, L. Lin, and Y. Yu. Instance-level salient object segmentation. *arXiv preprint arXiv:1704.03604*, 2017. 1, 2, 3, 8
- [19] G. Li and Y. Yu. Visual saliency based on multiscale deep features. In *CVPR*, pages 5455–5463, 2015. 2
- [20] G. Li and Y. Yu. Deep contrast learning for salient object detection. In *CVPR*, pages 478–487, 2016. 1, 2, 3, 6, 8
- [21] G. Li and Y. Yu. Visual saliency detection based on multi-scale deep cnn features. *TIP*, 25(11):5012–5024, 2016. 1
- [22] Y. Li, X. Hou, C. Koch, J. M. Rehg, and A. L. Yuille. The secrets of salient object segmentation. In *CVPR*, pages 280–287, 2014. 2
- [23] N. Liu and J. Han. Dhsnet: Deep hierarchical saliency network for salient object detection. In *CVPR*, pages 678–686, 2016. 3, 6
- [24] V. Mahadevan and N. Vasconcelos. Spatiotemporal saliency in dynamic scenes. *TPAMI*, 32(1):171–177, 2010. 3
- [25] P. Ochs, J. Malik, and T. Brox. Segmentation of moving objects by long term video analysis. *TPAMI*, 36(6):1187–1200, 2014. 5
- [26] F. Perazzi, P. Krahenbuhl, Y. Pritch, and A. Hornung. Saliency filters: Contrast based filtering for salient region detection. In *CVPR*, 2012. 5
- [27] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, pages 724–732, 2016. 5
- [28] A. Ranjan and M. J. Black. Optical flow estimation using a spatial pyramid network. *arXiv preprint arXiv:1611.00850*, 2016. 3
- [29] P. Tokmakov, K. Alahari, and C. Schmid. Learning video object segmentation with visual memory. *arXiv preprint arXiv:1704.05737*, 2017. 3, 8
- [30] D. Tsai, M. Flagg, A. Nakazawa, and J. M. Rehg. Motion coherent tracking using multi-label mrf optimization. *IJCV*, 100(2):190–202, 2012. 5
- [31] W.-C. Tu, S. He, Q. Yang, and S.-Y. Chien. Real-time salient object detection with a minimum spanning tree. In *CVPR*, pages 2334–2342, 2016. 6
- [32] L. Wang, H. Lu, X. Ruan, and M.-H. Yang. Deep networks for saliency detection via local estimation and global search. In *CVPR*, pages 3183–3192, 2015. 2
- [33] L. Wang, L. Wang, H. Lu, P. Zhang, and X. Ruan. Saliency detection with recurrent fully convolutional networks. In *ECCV*, pages 825–841. Springer, 2016. 2, 3, 6
- [34] W. Wang, J. Shen, and F. Porikli. Saliency-aware geodesic video object segmentation. In *CVPR*, pages 3395–3402, 2015. 6
- [35] W. Wang, J. Shen, and L. Shao. Consistent video saliency using local gradient flow optimization and global refinement. *TIP*, 24(11):4185–4196, 2015. 1, 2, 3, 6
- [36] W. Wang, J. Shen, and L. Shao. Video salient object detection via fully convolutional networks. *TIP*, 27(1):38–49, 2018. 1, 2, 5, 6
- [37] Y. Wei, X. Liang, Y. Chen, X. Shen, M.-M. Cheng, J. Feng, Y. Zhao, and S. Yan. Stc: A simple to complex framework for weakly-supervised semantic segmentation. *TPAMI*, 39(11):2314–2320, 2017. 1
- [38] H. Wu, G. Li, and X. Luo. Weighted attentional blocks for probabilistic object tracking. *The Visual Computer*, 30(2):229–243, 2014. 1

- [39] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *NIPS*, pages 802–810, 2015. 4
- [40] J. Yang and M.-H. Yang. Top-down visual saliency via joint crf and dictionary learning. In *CVPR*, pages 2296–2303. IEEE, 2012. 2
- [41] J. Zhang, S. Sclaroff, Z. Lin, X. Shen, B. Price, and R. Mech. Minimum barrier salient object detection at 80 fps. In *ICCV*, pages 1404–1412, 2015. 6
- [42] R. Zhao, W. Ouyang, H. Li, and X. Wang. Saliency detection by multi-context deep learning. In *CVPR*, pages 1265–1274, 2015. 2
- [43] R. Zhao, W. Ouyang, and X. Wang. Unsupervised saliency learning for person re-identification. In *CVPR*, pages 3586–3593, 2013. 1
- [44] X. Zhu, Y. Wang, J. Dai, L. Yuan, and Y. Wei. Flow-guided feature aggregation for video object detection. *arXiv preprint arXiv:1703.10025*, 2017. 3
- [45] X. Zhu, Y. Xiong, J. Dai, L. Yuan, and Y. Wei. Deep feature flow for video recognition. *arXiv preprint arXiv:1611.07715*, 2016. 3, 4