

ビデオから人目を惹く領域を検出する

# 時空間深層特徴による映像中の顕著物体検出


Video Salient Object Detection Using Spatiotemporal Deep Features

Trung-Nghia Le (SOKENDAI)

Akihiro Sugimoto (NII)

## Introduction

### Background



映像中の顕著物体検出はコンピュータービジョンおよび  
マルチメディアにおいて重要  
SOD from video is important for computer vision and multimedia


### Contribution

深層特徴と視覚的顕著性モデルの両方に時間的情報を利用  
Exploiting temporal information in both deep feature and saliency model


時空間深層特徴を提案  
Proposing spatiotemporal deep feature

時空間深層特徴を用いた STCRF モデルを開発  
Developing SpatioTemporal Conditional Random Field (STCRF)  
model using spatiotemporal deep features

## Proposed method



### Spatiotemporal Conditional Random Field




$$\text{Energy Optimization: } \hat{\ell} = \underset{\ell}{\operatorname{argmin}} E(\ell, x; \theta) = \underset{\ell}{\operatorname{argmin}} \left[ \sum_{i \in V} \psi_u(\ell_i, x; \theta_u) + \sum_{(i,j) \in E} \psi_b(\ell_i, \ell_j, x; \theta_b) \right]$$

$$\text{Unary Potential: } \psi_u(\ell_i, x; \theta_u) = \theta_u \ell_i(x)$$

$$\text{Binary Potential: } \psi_b(\ell_i, \ell_j, x; \theta_b) = \begin{cases} \theta_{bs}[\ell_i \neq \ell_j] \exp\left(-\frac{\|F_i(x) - F_j(x)\|^2}{2\sigma^2}\right) & (i, j) \in \text{Edge}_s \\ \theta_{bt}[\ell_i \neq \ell_j] & (i, j) \in \text{Edge}_t \end{cases}$$

## Experiments


### 10-Clips



SegTrack2



DAVIS



### Comparision with State-of-the-art Methods

Dataset Metric	10-Clips	SegTrack2	DAVIS
	F-Adap	F-Adap	F-Adap
STCRF (Our method)	<b>0.927</b>	<b>0.817</b>	<b>0.794</b>
LD [T.Liu, PAMI 2011]	0.637	0.286	0.252
LGFOGR [W.Wang, TIP 2015]	0.629	0.500	0.537
RST [T.N.Le, PSIVT 2015]	0.827	0.510	0.627
SAG [W.Wang, CVPR 2015]	0.755	0.504	0.494
STS [F.Zhou, CVPR 2014]	0.591	0.471	0.379
DCL [G.Li, CVPR 2016]	<b>0.935</b>	<b>0.734</b>	0.664
DHS [N.Liu, CVPR 2016]	0.923	0.733	<b>0.715</b>
ELD [G.Lee, CVPR 2016]	0.893	0.611	0.572

### SegTrack2



Input



Ours LD (PAMI 2011)




Ours LD (PAMI 2011)




Ours LD (PAMI 2011)


### DAVIS




Input



Ours LD (PAMI 2011)



Ours LD (PAMI 2011)



Ours LD (PAMI 2011)

### Comparison with State-of-the-art Methods

Setting Description	10-Clips	SegTrack2	DAVIS
	F-Adap	F-Adap	F-Adap
Local	-	0.868	0.590
STF	-	0.887	0.658
STF	CRF	0.916	0.789
STF	STCRF	<b>0.927</b>	<b>0.817</b>
STCRF			<b>0.794</b>

### Detail Analysis

### DAVIS



Input



RST (PSIVT 2015)



RST (PSIVT 2015)



RST (PSIVT 2015)

### Comparison with State-of-the-art Methods

Setting Description	10-Clips	SegTrack2	DAVIS
	F-Adap	F-Adap	F-Adap
Local	-	0.868	0.600
STF	-	0.887	0.658
STF	CRF	0.916	0.789
STF	STCRF	<b>0.927</b>	<b>0.817</b>
STCRF			<b>0.794</b>