

# Learning Dynamic GMM for Attention Distribution on Single-face Videos

Yun Ren, Zulin Wang, Mai Xu, Haoyu Dong, Shengxi Li

The School of Electronic and Information Engineering, Beihang University  
Beijing, 100191 China

MaiXu@buaa.edu.cn

## Abstract

The past decade has witnessed the popularity of video conferencing, such as FaceTime and Skype. In video conferencing, almost every frame has a human face. Hence, it is necessary to predict attention on face videos by saliency detection, as saliency can be used as a guidance of region-of-interest (ROI) for the content-based applications. To this end, this paper proposes a novel approach for saliency detection in single-face videos. From the data-driven perspective, we first establish an eye tracking database which contains fixations of 70 single-face videos viewed by 40 subjects. Through analysis on our database, we investigate that most attention is attracted by face in videos, and that attention distribution within a face varies with regard to face size and mouth movement. Inspired by the previous work which applies Gaussian mixture model (GMM) for face saliency detection in still images, we propose to model visual attention on face region for videos by dynamic GMM (DGMM), the variation of which relies on face size, mouth movement and facial landmarks. Then, we develop a long short-term memory (LSTM) neural network in estimating DGMM for saliency detection of single-face videos, so called LSTM-DGMM. Finally, the experimental results show that our approach outperforms other state-of-the-art approaches in saliency detection of single-face videos.

## 1. Introduction

Visual saliency [5] aims at predicting how much each pixel or region of an image/video attracts human's attention, which has been widely used in areas of object detection [16], video quality assessment [39] and perceptual video coding [35]. The studies on visual saliency can be traced back to 1998, when Itti and Koch [19] explored that intensity, color and orientation information in an image can be employed to predict image's saliency map. Afterwards, they extended their work to video saliency detection [18]. During the past two decades, extensive approaches, such as [17, 4, 38, 9, 30, 24, 13, 8, 12], have been proposed for

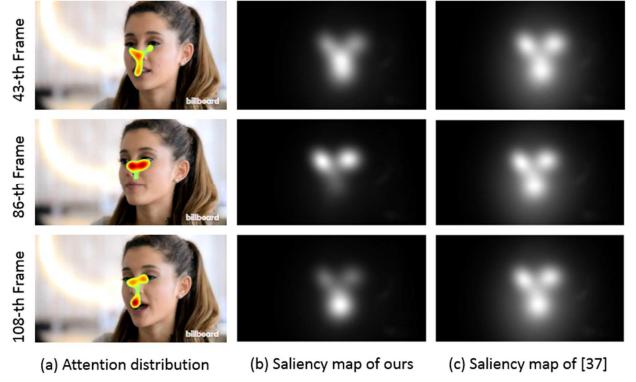


Figure 1: An example of video saliency maps generated by our approach and [37]. Note that [37] is a saliency detection approach for images, while ours works on videos. Here, the saliency maps of [37] are generated by regarding each video frame as a still image. The visual attention distribution by 40 subjects is also shown in this figure.

detecting saliency in videos. All these saliency detection approaches are heuristic ones, as they are generally driven by incorporating biologically-inspired features. However, the biologically-inspired features of these approaches rely heavily on the unmatured study of the human visual system (HVS), leading to inferior performance in saliency detection.

Recently, the top-down approaches ([21, 14, 15, 27, 31, 7, 22, 28, 36, 6, 40, 20, 37]) have become more prevalent in both image and video saliency detection, which learn saliency model from human fixations on training images/videos. These top-down approaches found out that some high-level features are indeed attractive to visual attention. In particular, face is an obvious high-level feature to attract visual attention, and thus many top-down approaches have incorporated face as a channel for saliency detection of face images [6, 40, 20, 37]. To be more specific, Cerf *et al.* [6] investigated from eye tracking data that face is highly correlated with attention, and they therefore proposed to integrate face channel with the channels of Itti's model [19] using equal weights, for detecting saliency of face images. Later, Zhao *et al.* [40] found that the face channel is more important than other channels. Accordingly, they proposed to learn

















- [35] M. Xu, X. Deng, S. Li, and Z. Wang. Region-of-interest based conversational hevc coding with hierarchical perception model of face. *IEEE Journal of Selected Topics on Signal Processing*, 8(3), 2014.
- [36] M. Xu, L. Jiang, Z. Ye, and Z. Wang. Bottom-up saliency detection with sparse representation of learnt texture atoms. *Pattern Recognition*, 60:348–360, 2016.
- [37] M. Xu, Y. Ren, and Z. Wang. Learning to predict saliency on face images. In *ICCV*, pages 3907–3915, 2015.
- [38] L. Zhang, M. H. Tong, and G. W. Cottrell. Sunday: Saliency using natural statistics for dynamic analysis of scenes. In *Proceedings of the 31st Annual Cognitive Science Conference*, pages 2944–2949. AAAI Press Cambridge, MA, 2009.
- [39] L. Zhang, M. Wang, L. Nie, R. Hong, Y. Xia, and R. Zimermann. Biologically inspired media quality modeling. In *ACM international conference on Multimedia (ACM MM)*, pages 491–500, 2015.
- [40] Q. Zhao and C. Koch. Learning a saliency map using fixated locations in natural scenes. *Journal of vision*, 11(3):9, 2011.