

Two-Stream Temporal Convolutional Network for Dynamic Facial Attractiveness Prediction

Nina Weng¹, Jiahao Wang², Annan Li^{2*}, Yunhong Wang² Shenyuan Honors College of Beihang University, Beijing, China¹ School of Computer Science and Engineering, Beihang University, Beijing, China² wengnn@buaa.edu.cn



Introduction

In the field of facial attractiveness pre diction, while deep models using stati c pictures have shown promising resul ts, little attention is paid to dynamic f acial information. Meanwhile, the incr easing popularity of short video apps creates an enormous demand for facia l attractiveness prediction from short video clips. In this work, we target on the dynamic facial attractiveness pred iction problem. A large-scale video-ba sed facial attractiveness prediction dat aset (VFAP) with more than one thous and clips from TikTok is collected. W e propose a two-stream temporal conv olutional network (2S-TCN) to captur e dynamic attractiveness features fro m both facial appearance and landmar ks. We employ attentive feature enhan cement along with specially designed modality and temporal fusion strategi es to better explore the temporal dyna mics. Extensive experiments on the pr oposed VFAP dataset demonstrate the superiority of 2S-TCN.

Video-based Facial Attractiveness Prediction Dataset



All videos in VFAP are collected from TikTok. After carefully manual selecti on, we build the VFAP dataset with 1,430 short videos. The dataset is divided into four subsets according to the different TikTok channels. The average fra me number of each video is about 314, and the whole dataset contains over 4 49k facial frames. Our attractiveness score is generated according to the inter active behaviors in social media, including pressing the like button, leaving a comment and forwarding the video to others. The ranking of facial attractive ness can be observed in the figure above, where four examples sampled from the top 1%, 25% and last 1% are presented.

Two-Stream Temporal Convolutional Network



As shown above, we first extract frame-level facial appearance and landmark features in the spatial feature extraction module. The proposed 2S-TCN mode l then performs temporal modeling on these spatial features. Within 2S-TCN, the appearance features are augmented by the attentive feature enhancement module. Modality and temporal score fusions are successively performed to g et the predicted attractiveness score.

Experimental Results

Method	SO	S1	S2	S3	All
AlexNet	0.08920	0.01559	0.03243	0.13180	0.01388
Resnet-18	0.11078	0.00787	0.03016	0.12721	-0.00240
ResNeXt-50	0.12934	0.06878	0.02485	0.12898	-0.00388
2S-TCN	0.38621	0.26273	0.32138	0.38699	0.18965

In order to validate the significance of te mporal information in this task, we com pare our proposed dynamic framework

with three state-of-the-art static FAP methods (i.e. AlexNet, Resnet-18, ResN eXt-50), which are selected from the top performers of the SCUT-FBP5500. I t can be observed from the results that the static methods do not predict well on the dynamic facial data. In comparison, our 2S-TCN model has much high er performance in both single-subset and all-subset experiments, demonstrati ng the importance of temporal modeling. We also conduct extensive ablation studies in order to validate the design of our 2S-TCN.

Motivation



The continuous face sequence or vide o is naturally more competent in expr essing facial attractiveness compare w ith static pictures. On one hand, psych ological and neuroscience studies hav e shown that temporal cues play an im portant role in the perception of huma n face. Dynamic face sequences are b etter at conveying emotions and expre ssions. On the other hand, the increasi ng popularity of short video apps crea tes an enormous demand for FAP fro m short video clips.