

# Use of Affective Visual Information for Summarization of Human-Centric Videos

Berkay Köprü, Engin Erzin

**Abstract**—Increasing volume of user-generated human-centric video content and their applications, such as video retrieval and browsing, require compact representations that are addressed by the video summarization literature. Current supervised studies formulate video summarization as a sequence-to-sequence learning problem and the existing solutions often neglect the surge of human-centric view, which inherently contains affective content. In this study, we investigate the affective-information enriched supervised video summarization task for human-centric videos. First, we train a visual input-driven state-of-the-art continuous emotion recognition model (CER-NET) on the RECOLA dataset to estimate emotional attributes. Then, we integrate the estimated emotional attributes and the high-level representations from the CER-NET with the visual information to define the proposed affective video summarization architectures (AVSUM). In addition, we investigate the use of attention to improve the AVSUM architectures and propose two new architectures based on temporal attention (TA-AVSUM) and spatial attention (SA-AVSUM). We conduct video summarization experiments on the TvSum database. The proposed AVSUM-GRU architecture with an early fusion of high level GRU embeddings and the temporal attention based TA-AVSUM architecture attain competitive video summarization performances by bringing strong performance improvements for the human-centric videos compared to the state-of-the-art in terms of F-score and self-defined face recall metrics.

**Index Terms**—Affective computing, video summarization, continuous emotion recognition, neural networks.



## 1 INTRODUCTION

Multimedia applications and services have seen a surge in recent years. There is more than 500 hours upload per minute to video streaming platforms such as Youtube and Twitch. Especially user-generated human-centric videos allocate major part of the available videos in these video streaming platforms. Such a massive video resource creates two crucial needs: (i) personalized human-computer interactions (HCI), and (ii) efficient representations and retrieval of video contents. Emotion recognition addresses the former by providing affect aware applications [1], while video summarization addresses the latter [2].

Emotion recognition is a critical task that enables personalized HCI applications and understanding of personal choices. Humans are emotional creatures, even for many cognitive tasks decision making processes are driven primarily by emotions [3]. Emotions are represented both in discrete and continuous domains. Discrete categorical emotions, such as happiness and sadness, can also be represented in the 3-dimensional continuous affect space of activation, valence, and dominance attributes, which are the indicators of activeness-passiveness, positiveness-negativeness, and dominance-submissiveness, respectively [4], [5].

In the literature, estimation of activation, valence and dominance attributes is referred as Continuous Emotion Recognition (CER). Although CER is widely studied by the speech processing community [6], [7], [8], it has been

studied over the visual channels [9], [10], [11]. In [9], a deep attention based convolutional network is proposed to detect facial expressions and emotions. Aspandi et al. [10] proposed an adversarial training approach to jointly estimate whether the image is fake or not while estimating the activation and valence (AV) attributes. In [11], VGG-16 driven visual features are used as the input of a stacked convolutional recurrent neural network (CRNN) for affect recognition in the wild.

Compact and efficient representations of videos can be extracted by selecting subsets of frames from the videos through video summarization techniques. Video summarization is categorized regarding the output type which could be video frames or video fragments [2], [12], [13]. The former, finding key frames, is known as video storyboard generation. The latter, chronologically stitched segment selection, is referred as video skimming. Due to the discrete nature of storyboard generation, it lacks in smoothness and naturalness. However, since it does not require synchronization, it offers more degrees of freedom in data organization than video skimming [2].

Early video summarization approaches are unsupervised and made use of low level similarity measures between the frames [14], [15], [16], [17], [18]; while recent unsupervised studies apply Generative Adversarial Networks (GANs) [19] and attention [20], [21] to solve the video summarization problem. They use GANs to reconstruct the input video from the selected key frames for unsupervised video summarization. Also low-level measures, such as color histogram intersections [16], [17] and mutual information [18], are investigated for unsupervised video summarization.

Supervised approaches for video summarization have

B. Köprü is with the KUIS AI Lab and Electrical & Electronics Engineering Department, Koç University, Istanbul, Turkey.

E. Erzin is with the KUIS AI Lab, Computer Engineering Department and Electrical & Electronics Engineering Department, Koç University, Istanbul, Turkey.

E-mail: {bkopru17, eerzin}@ku.edu.tr

recently become popular that tend to perform better in other computer vision tasks such as object detection and segmentation. In [12], [22], [23], [24], the authors model the video summarization task as a sequence-to-sequence mapping problem. In [12], Long Short-Term Memory (LSTM) is adopted to tackle with variable range dependencies within the encoder-decoder type architectures. In contrast to recurrent models, [22] proposes a fully convolutional solution, where all frames are processed together. Recent approaches [23], [24], successfully adapted attention mechanisms into the video-summarization task.

Coupling of affective computing and video summarization has not been extensively studied in the literature. In two related video summarization studies, key frames are selected based on the physiological responses of the viewers that are expected to be highly correlated with the viewers' emotional states [25], [26]. In [25], the facial activity of the viewer is tracked, and then the frames are ranked to extract personal highlights from the videos by using heuristics. Money and Agius track physiological responses of the viewers, such as heart rate and electrodermal response, to analyze sub-segments of the videos [26]. These studies leverage physiological information that are highly related to emotional states of the humans to generate personal highlights/summaries. On the other hand, these approaches require at least one subject as a video viewer and do not provide feasible solutions for the automatic video summarization.

In this study, we address the use of affective information, which is extracted from visual data, for the summarization of human-centric videos. Unlike studies evaluating viewers' perspective [25], [26], we model affective states from humans in the video. That makes automatic affective video summarization feasible and specific for human-centric videos. For affective video summarization, we adopt a two-step approach. First, affective information is extracted using the convolutional recurrent neural networks. Then, we explore the affective information and attention mechanisms to enrich a fully convolutional network based video summarization. To summarize, the main contributions of this study are as follows:

- We formulate a novel end-to-end learning problem for affective-information enriched video summarization targeting human-centric videos.
- We model affective information in terms of emotional attributes and learned embeddings extracted from the CER module.
- We investigate the use of attention mechanisms in video summarization for human-centric videos and develop temporal and spatial attention-based frameworks.
- We carry out affective video summarization evaluations on the state-of-the-art using both standard and self-defined evaluation metrics.

The rest of this paper is organized as follows. Section 2 reviews related work, and Section 3 describes main building blocks of the proposed framework. Section 4 presents the experiments conducted together with the performance evaluations. Finally, conclusion is presented in Section 5.

## 2 RELATED WORK

Our proposed affective video summarization framework integrates emotion recognition into video summarization. In this regard, we first point out several emotion recognition studies that relates video summarization. Then, we briefly discuss unsupervised and supervised video summarization literature.

Emotion recognition studies formulate the problem as discrete emotion recognition (DER) [27], [28], [29] or continuous emotional attribute regression [1], [30]. In [27], stacked CNN-RNN architecture is proposed to extract local and global features to classify emotions. For the CER problems, Schmitt et al. investigates RNNs with CCC loss function after extracting low-level descriptors (LLDs) such as mel-frequency cepstral coefficients (MFCCs) and zero crossing rate from speech signal [30]. Recently, multi-modal approaches are explored for CER and DER [1], [31]. Tzirakis et al. design an end-to-end network utilizing raw video, audio and text, where visual features are extracted using 3 stage High-Resolution Network and audio features are extracted via multiple 1-D convolutional layers [31]. Contextual features are extracted from text by first generating point-wise n-grams using convolutional layers, then linearly projecting these sub-features with multiple heads to increase the diversity. Finally, extracted features are fused using attention. In [1], visual information is expressed using facial attributes and audio information is represented by the MFCCs. Audio-visual information is fused at the feature level, and a CRNN model is trained with multi-task learning for the CER problem.

Two recent emotion recognition studies are interesting in the context of video summarization [32], [33]. Xu et al. investigates information transfer from image and textual data of videos for emotion recognition, emotion attribution and emotion-oriented summarization [32]. First, they learn video representations using Image Transfer Encoding and textual representations using zero-shot learning from auxiliary datasets. Then, they perform a categorical emotion recognition using the Support Vector Machine (SVM) classifier. Later, emotion attribution sets the contribution of each frame to the video's overall emotion. Finally, video summarization is formulated as a selection of key frames by maximizing an emotion attribute based score function. In the second related study, Tu et al. train a joint model to capture emotion attribution and recognition using multitask learning [33]. Later, similar to the first study [32], video summarization task is formulated as a post-processing optimization problem and solved using MINMAX dynamic programming. Note that both of these studies formulate summarization task as a optimization problem which is executed in post-processing. Hence their video summarization frameworks are not learning based and they do not explore how affective information alter behavior of the proposed summarization architecture. Furthermore, the proposed solutions are not evaluated on the state-of-the-art video summarization datasets.

Scarcity of the labeled data leads unsupervised learning studies for the video summarization problem. In a recent study, Jung et al. address unsupervised video summarization problem by first learning discriminative features over

a Variational Auto Encoder (VAE) and GAN-based architecture using variance loss to alleviate ineffective feature learning [34]. Then, they define a chunk and stride network (CSNet) to overcome the difficulty of learning for long-length videos. In another study, Zhao et al. presents a dual learning framework for the unsupervised video summarization [35]. They integrate the summary generation and video reconstruction tasks using multi-task learning so to reward the summary generator under the assistance of the video reconstructor. Zhou et al. formulates summarization task as sequential decision-making using an end-to-end reinforcement learning based framework [36]. They utilize a reward function that jointly accounts for diversity and representativeness of the generated summaries in an unsupervised setting.

Supervised studies on video summarization adapts encoder-decoder architectures [22], [23], [24], [37]. In [37], video is modeled as a 3-dimensional tensor, and 3D convolutional networks are used in the encoder to extract shallow and deep spatio-temporal features. Then, extracted multi-level features are fed into an LSTM based decoder. The proposed architecture is trained with the Sobolev loss, which constrains the derivative of the sequential data. With the great success of attention modules [38], [39], recent works in video summarization adapt attention into the decoder part. Ji et al. extract visual features using the GoogleNet and later encode with a Bidirectional LSTM network [23]. Encoded vectors combine Bahdanau attention [38], and then feed into an LSTM based decoder. Later the Bidirectional LSTM approach is extended to prevent semantic information loss [40]. In the extended architecture, an additional network analyzing the semantic information loss is added and used as a feedback mechanism from the decoder to the encoder using the Huber loss.

Although we apply supervised learning for video summarization task, our study differs from these studies in the exploitation of affective information for the video summarization task.

### 3 METHODOLOGY

In this paper, we investigate the use of affective information for enriching video summarization by capturing emotionally salient regions of the human-centric videos. First, we state and formulate the video summarization problem and define the visual feature extraction for both the CER and video summarization tasks. Then, an end-to-end framework for the CER is presented. Emotional attributes and high-level embeddings from the CER framework are later used as affective representations by the video summarization framework. Finally, we introduce the proposed affective video summarization architectures by first defining a video summarization baseline and then enriching this baseline with the fusion of affective information.

#### 3.1 Problem Statement

Video summarization is widely formulated as either a binary classification or a frame-level regression task. In the binary classification task, summarization outputs are either key-frames [12], [22] or key-shots [12], [41] from the video.

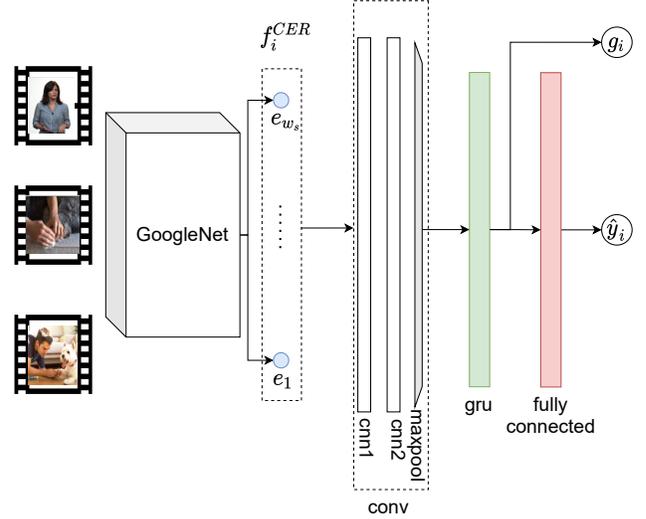


Fig. 1. CER-NET: The continuous emotion recognition network

On the other hand, frame-level importance scores are extracted in the regression task [12], [23].

In this study, we formulate the video summarization as a binary classification task where the positive labels correspond to the selected key-frames. The summarization network receives a feature matrix,  $\mathbf{v} \in \mathbb{R}^{N \times D}$ , and emits an output matrix,  $\mathbf{s} \in \mathbb{R}^{N \times C}$ , where  $N$  is the number of frames in the video,  $D$  is the dimensionality of the frame-level visual feature, and  $C$  is the number of classes. We take  $C = 2$  representing the positive and negative classes for the key-frame selection and these two nodes output class probability values at the output of the network. Then the positive class for key-frame selection or the negative class for frame skip are set by picking the node with higher probability. Eventually, the video summary is constructed from the key-frames that are labeled as positive.

In this study, we extract the visual information using the *GoogleNet* [42]. Output of the *pool5* layer of the pre-trained *GoogleNet* is used as the visual feature and represented as  $\mathbf{e}_i \in \mathbb{R}^D$  at frame  $i$  with dimension  $D = 1024$ .

#### 3.2 CER Network (CER-NET)

The continuous emotion recognition problem is set as the continuous regression of an emotional attribute from the temporal visual features. For this purpose, we construct a CER network (CER-NET), which consists of two back-to-back convolutional layers, a max pool layer in the temporal domain, a Gated Recurrent Unit (GRU) layer, and a fully connected layer as shown in Figure 1. We use the CER-NET to train two separate networks to estimate the activation and valence (AV) attributes separately. A temporal visual feature matrix is defined to be the input of the CER-NET. For this purpose, visual features around the  $i$ -th frame are cascaded to define the temporal visual feature as

$$\mathbf{f}_i^{\text{CER}} = [\mathbf{e}_{i-\Delta+1}, \dots, \mathbf{e}_{i-1}, \mathbf{e}_i, \mathbf{e}_{i+1}, \dots, \mathbf{e}_{i+\Delta}] \in \mathbb{R}^{D \times T} \quad (1)$$

at frame  $i$ , where  $T$  is the temporal window size and it is set as  $T = 2\Delta = 20$  frames.

We refer the group of back-to-back two convolutional and a max-pool layers as the *conv* layer for the sake of

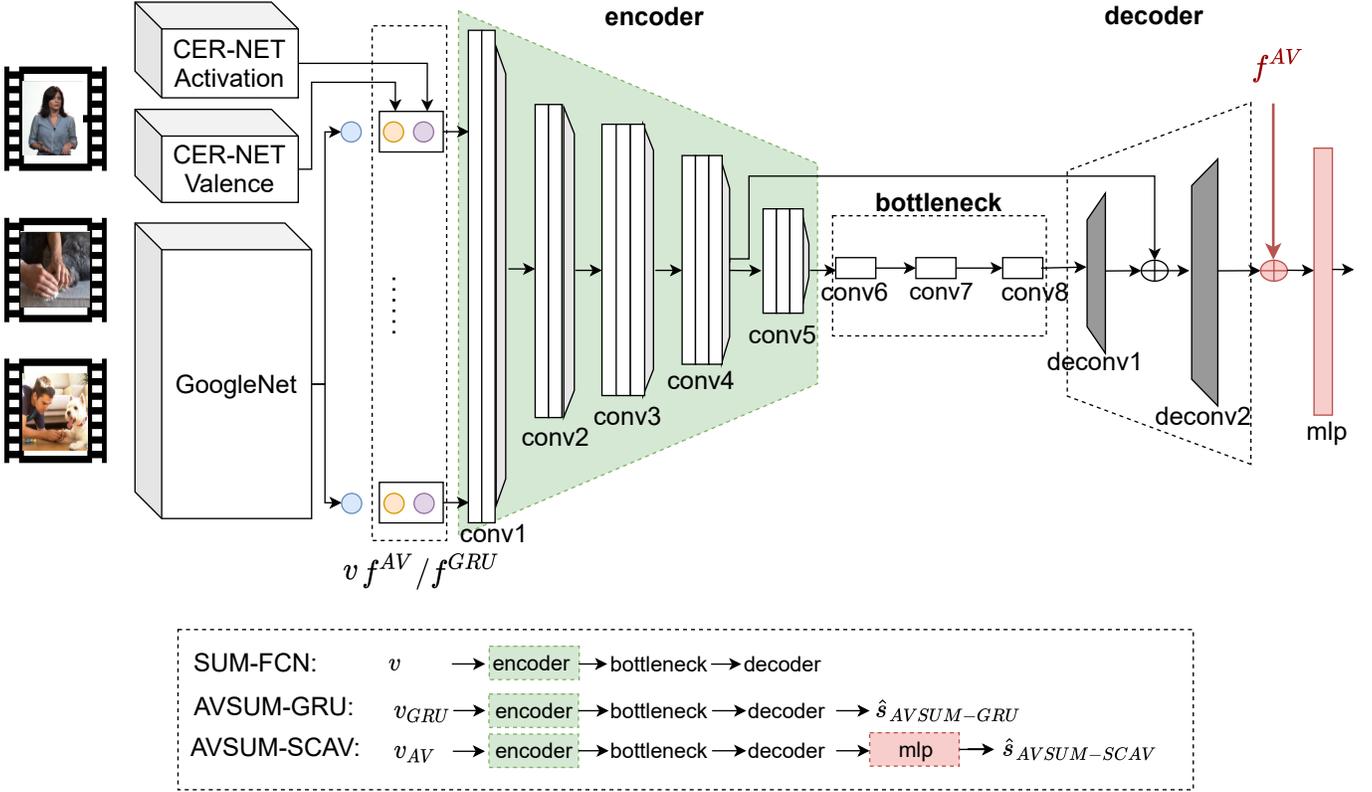


Fig. 2. An overview of the proposed AVSUM-GRU and AVSUM-SCAV architectures. AVSUM-GRU combines high-level affective information  $f_{GRU}$  from the CER-NET estimators with the GoogleNet driven visual features. AVSUM-SCAV defines a skip-connection for the emotional attributes  $f_{AV}$  as well as combining them with the GoogleNet driven visual features for the video summarization.

simplicity. In CER-NET, the *conv* layer models the spatial relations and provides a compact representation of the temporal window. In this manner, both temporally related features are highlighted and dimensionality of the representation is reduced. Dimensionality reduction is the key to complexity reduction and it also prevents overfitting. The compact representation from the *conv* layer is fed into GRU to model long-term temporal relations.

At the inference phase, CER-NET receives  $\mathbf{f}_i^{\text{CER}}$  and provides the GRU layer outputs as  $\mathbf{g}_i \in \mathbb{R}^G$  and the fully connected layer outputs as  $\hat{y}_i \in \mathbb{R}$ . We define the GRU layer output,  $\mathbf{g}_i$ , as an affective embedding, which can carry affective information to the video summarization model. On the other hand, the fully connected layer output,  $\hat{y}_i$ , represents the estimated emotional attribute (A or V) at frame  $i$  and delivers another source of affective information.

### 3.3 Affective Video Summarization

The proposed affective-information enriched video summarization (AVSUM) is based on a fully convolutional neural network (FCN) for semantic segmentation [43], which is adapted by [22] into video summarization (SUM-FCN). In the following, we briefly describe SUM-FCN and then present two fusion architectures to combine affective information with video summarization, which are later extended by temporal and spatial attention mechanisms.

#### 3.3.1 SUM-FCN

SUM-FCN adopts an encoder-decoder architecture where the encoder is based on fully convolutional layers and

the decoder is based on deconvolutions. Figure 2 includes the architecture of the SUM-FCN comprising encoder-bottleneck-decoder layers driven by the visual input  $\mathbf{v}$ . The encoder compresses the temporal information while increasing spatially, and the bottleneck passes it to the decoder to reconstruct necessary information from this compact representation. SUM-FCN receives the visual features for the whole video at once and provides the summarization outputs at once. Video frame rate is typically down-sampled before the summarization. Let  $\mathbf{v}$  be the input of SUM-FCN, then  $\mathbf{v} = [\mathbf{e}_1, \dots, \mathbf{e}_N]' \in \mathbb{R}^{N \times D}$  where  $N$  is number of frames in the down-sampled stream. Given  $\mathbf{v}$ , SUM-FCN emits  $\mathbf{s} \in \mathbb{R}^{N \times C}$ , where  $C$  is the dimension of the summarization annotations and set as  $C = 2$  as defined in section 3.1.

#### 3.3.2 Affective Feature Fusion for AVSUM

Affective information is extracted from the two CER-NET models, which are trained to estimate the activation and valence (AV) attributes separately. Two types of affective information cues are extracted from the CER-NET models: (i) the estimated AV attributes ( $\mathbf{f}^{\text{AV}}$ ), and (ii) the learned high-level CER-NET representations based on the GRU embeddings ( $\mathbf{g}^{\text{AV}}$ ). The estimated AV attribute vector  $\mathbf{f}_j^{\text{AV}}$  is constructed as a column vector from the estimated activation and valence attributes as  $\mathbf{f}_j^{\text{AV}} = [\hat{y}_j^A, \hat{y}_j^V]' \in \mathbb{R}^2$  at frame  $j$  in the down-sampled stream. Similarly,  $\mathbf{g}_j^{\text{AV}}$  is constructed by concatenating the outputs of the GRU layers from the two CER-NET models as  $\mathbf{g}_j^{\text{AV}} = [\mathbf{g}_j^A, \mathbf{g}_j^V]' \in \mathbb{R}^{2G}$ . Then,

the emotional attribute  $\mathbf{f}^{\text{AV}}$  and the affect embedding  $\mathbf{f}^{\text{GRU}}$  representations of the video are defined as

$$\mathbf{f}^{\text{AV}} = [\mathbf{f}_1^{\text{AV}}, \dots, \mathbf{f}_N^{\text{AV}}]' \in \mathbb{R}^{N \times 2} \quad (2)$$

$$\mathbf{f}^{\text{GRU}} = [\mathbf{g}_1^{\text{AV}}, \dots, \mathbf{g}_N^{\text{AV}}]' \in \mathbb{R}^{N \times 2G}. \quad (3)$$

The first proposed AVSUM architecture, referred as AVSUM-GRU, combines the affect embedding  $\mathbf{f}^{\text{GRU}}$  with the visual input  $\mathbf{v}$ . Figure 2 presents the AVSUM-GRU architecture receiving the  $\mathbf{v}_{\text{GRU}}$  input as

$$\mathbf{v}_{\text{GRU}} = \mathbf{v} \oplus \mathbf{f}^{\text{GRU}} \in \mathbb{R}^{N \times (D+2G)}, \quad (4)$$

where  $\oplus$  operator is representing the feature vector combining over the whole video.

Alternatively, we define the AVSUM-SCAV architecture, as illustrated in Figure 2, by combining the emotional attribute  $\mathbf{f}^{\text{AV}}$  with the visual input  $\mathbf{v}$  to receive  $\mathbf{v}_{\text{AV}}$  as

$$\mathbf{v}_{\text{AV}} = \mathbf{v} \oplus \mathbf{f}^{\text{AV}} \in \mathbb{R}^{N \times (D+2)} \quad (5)$$

and incorporating a long skip-connection by concatenating the emotional attribute  $\mathbf{f}^{\text{AV}}$  to the final layer of the summarization network. AVSUM-SCAV inserts a skip-connection to the output of *deconv2* layer and has a final fully connected *mlp* layer to reduce the dimension from  $C + 2$  to  $C$ .

### 3.3.3 Temporal Attention for AVSUM

We adapt multi-headed attention (MHA) mechanism into the AVSUM architecture to efficiently model temporal dependencies across the video frames. Figure 3 depicts the MHA based temporal attention structure, referred as TA-AVSUM, where MHA is placed to the output of *conv4* layer which emits  $\mathbf{X} \in \mathbb{R}^{M \times S}$ . Here,  $M$  and  $S$  are respectively the temporal and spatial dimensions of  $\mathbf{X}$ .

MHA receives three inputs as Query ( $\mathbf{Q}$ ), Key ( $\mathbf{K}$ ) and Value ( $\mathbf{V}$ ), then outputs a weighted summation of the rows of  $\mathbf{V}$ . The weights are calculated from the similarity between the  $\mathbf{Q}$  and  $\mathbf{K}$ . In this context, the MHA is defined as

$$\text{head}_h = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{W}_h^{\text{Q}}(\mathbf{K}\mathbf{W}_h^{\text{K}})^T}{\sqrt{M}}\right)\mathbf{V}\mathbf{W}_h^{\text{V}} \quad (6)$$

$$\text{MHA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \dots, \text{head}_H)\mathbf{W}^{\text{O}}, \quad (7)$$

where  $\mathbf{W}_h^{\text{Q}}$ ,  $\mathbf{W}_h^{\text{K}}$ ,  $\mathbf{W}_h^{\text{V}}$  and  $\mathbf{W}^{\text{O}}$  are the learned linear projections and  $H$  is the number of heads. We employ multi-headed self-attention, by setting  $\mathbf{Q}$ ,  $\mathbf{K}$  and  $\mathbf{V}$  to  $\mathbf{X}$ . Hence, the learned projections matrices are formed as  $\mathbf{W}_h^{\text{Q}}$ ,  $\mathbf{W}_h^{\text{K}}$ ,  $\mathbf{W}_h^{\text{V}} \in \mathbb{R}^{S \times \frac{S}{H}}$ , and  $\mathbf{W}^{\text{O}} \in \mathbb{R}^{M \times M}$ .

### 3.3.4 Spatial Attention for AVSUM

Motivated by the Squeeze and Excitation Networks [44], which attend to different channels of an image, we propose the fourth AVSUM network with spatial domain attention and refer it as SA-AVSUM. Figure 3 depicts the SA-AVSUM structure where we employ attention to the output of the *deconv2* layer. Unlike TA-AVSUM, we adopt single-headed attention, which is formulated in (6). Then,  $\mathbf{K}$  and  $\mathbf{V}$  are set to transpose of  $\hat{\mathbf{s}}_{\text{AVSUM}}$  and  $\mathbf{Q}$  is set from the affective embedding  $\mathbf{f}^{\text{GRU}}$ .

## 3.4 Model Training

Training of the CER-NET and video summarization models are executed in two phases. First the CER-NET models are trained, later they are fixed and integrated for the affective video summarization to train the AVSUM-GRU, AVSUM-SCAV, TA-AVSUM, and SA-AVSUM models.

### 3.4.1 CER-NET

The CER-NET models are trained separately to estimate the activation and valence attributes using the concordance correlation coefficient (CCC) based loss function. The loss function of the CER-NET is defined as the negated CCC value,

$$L_{\text{CER}} = -\frac{2\sigma_{y\hat{y}}^2}{\sigma_y^2 + \sigma_{\hat{y}}^2 + (\mu_y - \mu_{\hat{y}})^2}, \quad (8)$$

where  $y$  is the ground truth and  $\hat{y}$  is the estimated attribute.

### 3.4.2 Video Summarization Networks

The key frame selection problem has an imbalanced nature, since only a small number of frames are selected for the summary [22]. In order to overcome the imbalance problem, a weighted binary cross entropy loss is defined as

$$L_{\text{SUM}} = -\frac{1}{N} \sum_{j=1}^N w_{z_j} (z_j \log \hat{z}_j + (1 - z_j) \log (1 - \hat{z}_j)), \quad (9)$$

where  $z_j \in 0, 1$  is the binary ground truth,  $\hat{z}$  is the predicted score and  $w_{z_j}$  is the weight of the  $j^{\text{th}}$  frame. The weights for the binary target are defined as

$$w_0 = \frac{1}{N} \sum_{j=1}^N z_j \quad \text{and} \quad w_1 = 1 - w_0. \quad (10)$$

## 4 EXPERIMENTS

Experimental evaluations of the proposed models and comparisons with the state-of-the-art are performed using two datasets. In this section, we first introduce the datasets and evaluation metrics, then explain implementation details. Finally, experimental results are presented and discussed.

### 4.1 Datasets

We train and evaluate the CER-NET on the RECOLA dataset, which is a popular multi-modal dataset for emotion recognition [45]. The RECOLA dataset is composed of multi-modal recordings of dyadic conversations from 27 French speakers. From these 27 recordings, 18 of them are annotated, and the rest of the records are used for testing. The annotations are at the rate of 40 msec and from 6 different annotators. In total, we use 90 minutes of recordings from the RECOLA dataset in this study.

Experimental evaluations on the proposed AVSUM architectures are executed on the frequently used TvSum dataset [41]. Note that there is no available video summarization dataset containing only human-centric videos in the literature. The TvSum dataset contains 50 user generated videos from 10 different categories, such as vehicle tire changing, sandwich making, grooming an animal, etc. The demographics of the dataset in terms of face including frames in the summary and in total video clip is presented

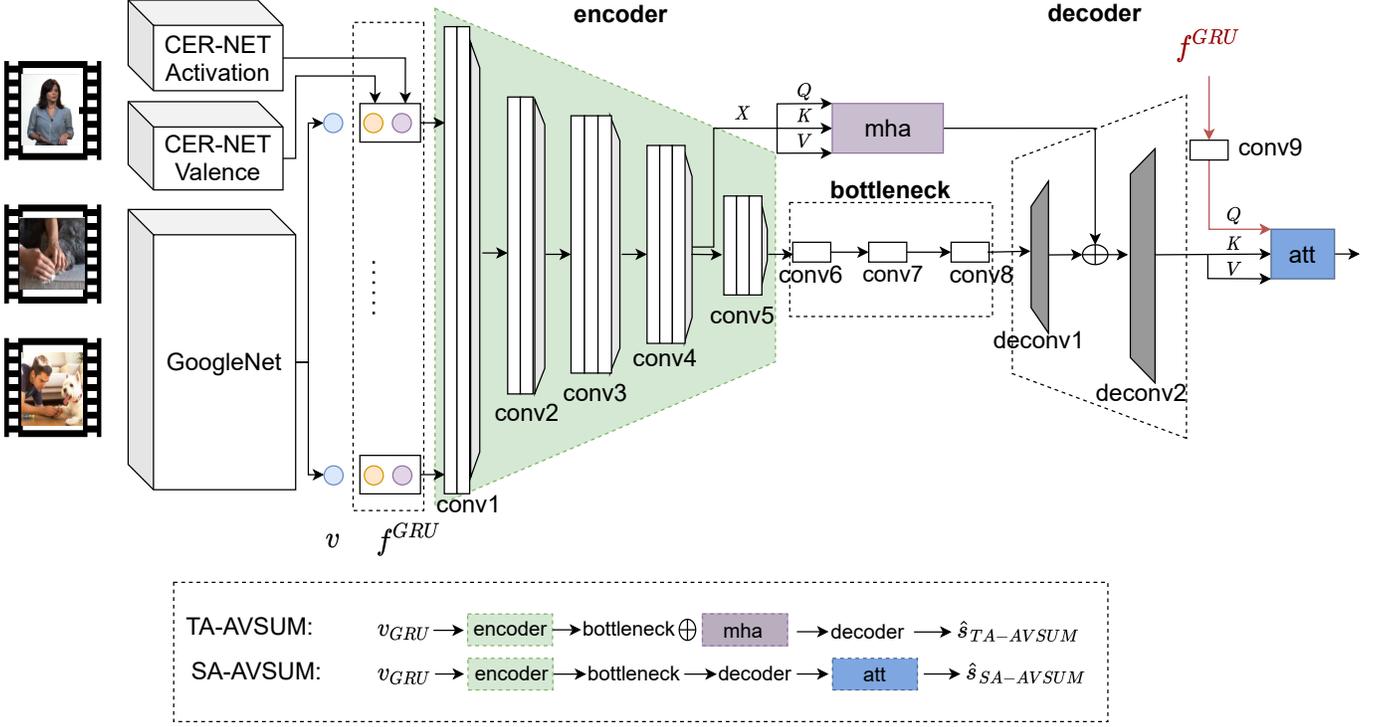


Fig. 3. An overview of the proposed TA-AVSUM and SA-AVSUM architectures. TA-AVSUM applies multi-head-self attention to the output of  $\text{conv4}$  layer ( $X$ ), on top of AVSUM architecture. On the other hand, SA-AVSUM modifies the long skip connection of affective features by applying a spatial attention.

in Figure 4. We categorize videos into human-centric and rest regarding the number of faces in the summary where we set the threshold to 80 frames labeling 15 videos as human-centric. Ground truths are provided by the frame-level importance score for each video from 20 raters. We follow the approach in [12], [22] to convert the frame-level importance scores into the keyshot-based summaries.

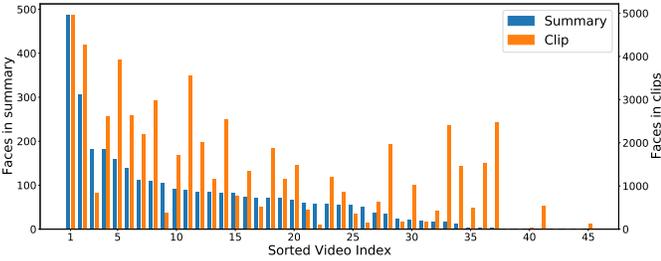


Fig. 4. Number of face including frames in the summary and total video for TvSum dataset where the videos are sorted regarding the number of faces in the summary.

## 4.2 Evaluation Metrics

Emotional attribute estimation with the CER-NET is evaluated using the CCC metric following [1], which is the negative of the  $L_{\text{CER}}$  loss in (8).

For the video summarization task, following [12], [22], [23], F-score is used as the evaluation metric. For the  $k$ -th video, let the ground truth binary summarization vector be  $s^k$  and estimated binary summarization vector to be  $\hat{s}^k$

where  $s^k, \hat{s}^k \in \mathbf{R}^N$ . Then precision  $P^k$  and recall  $R^k$  for the  $k$ -th video are calculated as

$$P^k = \frac{s^k \cdot \hat{s}^k}{\sum_j^N \hat{s}_j^k} \quad \text{and} \quad R^k = \frac{s^k \cdot \hat{s}^k}{\sum_j^N s_j^k}, \quad (11)$$

where  $j$  runs over the frames. Video level F-score is defined as the harmonic mean of the precision and recall,

$$F1^k = \frac{2P^k R^k}{P^k + R^k}. \quad (12)$$

Then, the final F-score metric  $F1$  is defined as the un-weighted average of the video level F-score values as

$$F1 = \frac{1}{K} \sum_{k=1}^K F1^k, \quad (13)$$

where  $K$  is the number of videos in the dataset.

Since the affective information is learned from a dataset of human-centric videos, capability of capturing affective frames of the video during the summarization is expected to be better with human-centric videos. By highlighting this fact, we define two new metrics for the evaluation of the affective video summarization. The first metric computes normalized F1 score differences with the baseline over the videos with the highest number of face appearances. To define this metric, let us first define the normalized F1 score difference for the  $k$ -th video with respect to the SUM-FCN baseline model as

$$\Delta F1^k = \frac{F1^k - F1_{SUM-FCN}^k}{F1_{SUM-FCN}^k}, \quad (14)$$

where  $F1^k$  refers to the F1 score of the model in evaluation. Also assume that all the videos in the dataset are sorted with the highest number of face appearances in descending order and are indexed with  $k_l$  for  $l = 1, \dots, K$ . Then, the cumulative F1 score difference metric for the *Top-L* human-centric videos is defined as

$$\Delta F1_L = \sum_{l=1}^L \Delta F1^{k_l}. \quad (15)$$

Associated with the  $\Delta F1_L$ , we also compute the F1 score of the *Top-L* human-centric videos and refer it as  $F1_L$ .

Human-centric nature of the videos can be associated with the face appearances in the video frames. Motivated with this fact, we set a second metric for evaluation of the affective video summarization as the recall rate of face appearing frames in the extracted summary. Let  $\mathbf{d}_F^k$  be the binary vector representing whether a frame includes a face appearance or not for the  $k$ -th video. Binary face appearance vectors are extracted by the histogram of oriented gradients (HOG) based face detector [46]. Then, the recall rate of face appearing frames, let's refer it as face recall,  $R$  is defined as

$$R = \frac{1}{K} \sum_{k=1}^K \frac{\hat{\mathbf{s}}^k \cdot (\mathbf{d}_F^k \odot \mathbf{s}^k)}{\sum_j^N s_j^k}, \quad (16)$$

where  $(\mathbf{d}_F^k \odot \mathbf{s}^k)$  is the Hadamard product of  $\mathbf{d}_F^k$  and  $\mathbf{s}^k$ .

We also study statistical differences of a given affective feature dimension,  $f$ , across face appearing and not-appearing frames. For this purpose, Kullback-Leibler (KL) divergence of  $f$  in these two classes is defined as

$$D(P_f || Q_f) = \sum_{x \in \mathcal{X}} P_f(x) \log \left( \frac{P_f(x)}{Q_f(x)} \right), \quad (17)$$

where  $P_f$  and  $Q_f$  are respectively probability distributions of the affective feature dimension  $f$  across face appearing and face not-appearing frames. Note that the affective feature dimension  $f$  is driven from the affective feature set as  $f \in \{f^A, f^V, g_1^{AV}, g_2^{AV}, \dots, g_{2G}^{AV}\}$ .

## 4.3 Implementation Details

### 4.3.1 CER-NET

The window size  $T$  is selected as 20 frames. The *cnm1* and *cnm2* layers have 10 filters, max-pooling layer reduces the temporal dimension from 20 to 5, *gru* layer has 10 cells and *FCN* layer has 1 node. Hence, the dimensionality of  $\mathbf{f}^{AV}$  is 2 and  $\mathbf{g}_{AV}$  is 20 ( $G = 10$ ).

Annotated recordings of the RECOLA dataset are used during the training of the CER-NET. RECOLA recordings are divided as 10% for the test, 10% for the validation, and 80% for the training. We applied Adam optimizer with a learning rate of  $10^{-4}$  and with the batch size of 256.

TABLE 1  
CCC performance of the CER-NET emotion recognition model on the RECOLA dataset

Model	Activation	Valence
CER-NET	0.40	0.17
CER-MTL Facial [1]	0.15	0.06
End-to-end Visual Network [47]	0.36	0.48

### 4.3.2 Video Summarization Networks

Following [12], [22], TvSum videos are downsampled to 2 fps, and frames are fed into GoogleNet. For the AVSUM-SCAV, the input feature dimension is set as 1026. The input dimension of the AVSUM-GRU, TA-AVSUM, and SA-AVSUM models became 1044 with the GRU embeddings. For the TA-AVSUM, we set the number of heads  $H = 4$ .

We mimic fixed size cropping in semantic segmentation by uniformly sampling the video frames and using  $N = 320$  [22]. We adopted the leave one-group out cross-validation technique to compare the performances. At each fold, 9 videos are selected, and the rest of the videos are used for training. The training is held for 50 epochs with a batch size of 5 videos. During the training phase, Adam optimizer with the learning rate of  $10^{-3}$  is used. For each training fold, a model achieving the highest F-score ( $F1$ ), and a model achieving the highest face recall ( $R$ ) are selected for the performance evaluations.

## 4.4 CER-NET Performance

Table 1 presents the CCC performance of the CER-NET emotion recognition model on the RECOLA dataset. The proposed architecture performs better at estimating the activation than the valence. The end-to-end CER-NET architecture outperforms CER-MTL Facial model [1] in both activation and valence by achieving CCC of 0.40 and 0.17 respectively. The end-to-end visual network in [47] performs strongly for the valence estimation and fairly close with the CER-NET for the activation estimation. CER-NET outperforms [47] by 4% at estimating the activation, while [47] achieves CCC of 0.48 for valence and CER-NET achieves 0.17. Different than CER-NET, CER-MTL Facial receives visual information as facial activation units and optical flow vectors. Due to the input dimension, CER-NET has more trainable coefficients leading to a more complex structure than the CER-MTL Facial.

In comparison with these two baseline visual models [1], [47], CER-NET performs competitively in representing affective information on the visual channel.

## 4.5 Cumulative AVSUM Performance

In this section, we first present performance evaluations of the proposed affective video summarization models in terms of cumulative F-score and face recall metrics. Then, the video level performances are investigated to better highlight characteristics of videos that have improved summarization performance with the affective cues.

Table 2 presents cumulative F-score ( $F1$ ) and face recall ( $R$ ) together with the *Top-15* F-score ( $F1_{15}$ ) and face recall ( $R_{15}$ ) performances for the proposed AVSUM and the baseline SUM-FCN models. In each column, top-two scoring performances are highlighted in bold. Recall that we apply two model selection criteria based on F-score and face recall maximization. Each model selection criterion is observed to favor its related performance metric in the evaluations. That is, maximization of  $F1$  (Max  $F1$ ) yields higher F-score while maximization of  $R$  (Max  $R$ ) yields higher face recall  $R$ .

Observing the cumulative  $F1$  performances, AVSUM-GRU model is competitive for both model selection criteria and observed as the best performing model with the

TABLE 2

Cumulative F-score ( $F1$ ) and face recall ( $R$ ) together with the  $Top-15$  F-score ( $F1_{15}$ ) and face recall ( $R_{15}$ ) performances (top two models are in bold) of the AVSUM and the SUM-FCN models with the maximization of  $F1$  and  $R$  model selection criteria

Model	Max $F1$				Max $R$			
	$F1$ (%)	$F1_{15}$ (%)	$R$ (%)	$R_{15}$ (%)	$F1$ (%)	$F1_{15}$ (%)	$R$ (%)	$R_{15}$ (%)
SUM-FCN [22]	57.46	60.00	53.20	<b>65.52</b>	<b>54.10</b>	57.40	60.62	60.28
AVSUM-GRU	<b>57.50</b>	<b>60.25</b>	<b>54.04</b>	59.14	<b>54.02</b>	57.40	60.77	<b>66.20</b>
AVSUM-SCAV	56.64	<b>60.60</b>	52.11	63.80	53.80	57.42	<b>62.06</b>	59.64
TA-AVSUM	<b>57.47</b>	59.95	<b>53.22</b>	<b>65.55</b>	53.79	<b>59.23</b>	<b>65.12</b>	<b>70.31</b>
SA-AVSUM	55.92	59.98	49.25	62.38	52.76	<b>59.90</b>	58.85	62.55

Max  $F1$  criterion. On the other hand, the cumulative  $R$  score highlights AVSUM-GRU and the temporal attention based TA-AVSUM models with the Max  $F1$  criterion and AVSUM-SCAV and TA-AVSUM models with the Max  $R$  criterion. The temporal attention based TA-AVSUM model especially performs significantly better with the Max  $R$  criterion achieving 65.12% face recall rate.

Affective information is modeled with the continuous emotion recognition task trained on the RECOLA dataset. Since RECOLA is a human-centric dataset, which includes facial videos, we also choose to evaluate the video summarization performance for the  $Top-L$  human-centric videos, where  $L$  is set as 15 with the discussion in section 4.1. Table 2 presents the  $Top-15$  F-score  $F1_{15}$  and face recall  $R_{15}$  performances. While attention based SA-AVSUM and TA-AVSUM models perform best for the  $F1_{15}$  score with the Max  $R$  criterion, AVSUM-SCAV and AVSUM-GRU models perform best with the Max  $F1$  criterion. Overall, the best performance is 60.60%  $F1_{15}$  score with the Max  $F1$  criterion for the AVSUM-SCAV model. Observing the  $Top-15$  face recall  $R_{15}$  performances, while TA-AVSUM model is competitive with the baseline SUM-FCN model with the Max  $F1$  criterion, it performs significantly superior with the Max  $R$  criterion achieving 70.31% face recall  $R_{15}$  rate.

Table 2 highlights two runner up models, AVSUM-GRU and TA-AVSUM. AVSUM-GRU model sustains strong F-score rates with the Max  $F1$  criterion, especially for human-centric videos targeted with  $F1_{15}$  performance. Alternatively, while temporal attention based TA-AVSUM model performs strongly for the face recall with the Max  $R$  criterion and attains 70.31% face recall  $R_{15}$  rate, it also sustains a competitive performance for the F-score and face recall rates with the Max  $F1$  criterion. Hence temporal attention is observed to better integrate the affective information for the affective video summarization.

## 4.6 Explainability

We conduct explainability evaluations to better understand contributions of the affective feature dimensions and the proposed model architectures for the affective video summarization. First, we present KL divergence analysis for the affective feature dimensions. Then, we investigate video level summarization performances of the AVSUM models in terms of the cumulative F-score difference metric  $\Delta F1_L$  for the  $Top-L$  human-centric videos.

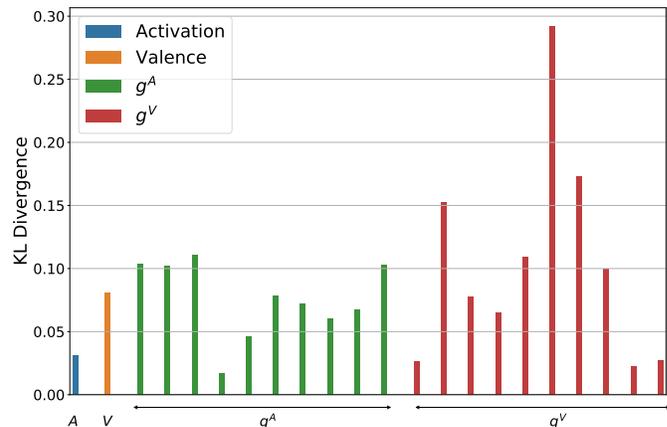


Fig. 5. KL divergence  $D(P_f||Q_f)$  of each affective feature dimension

### 4.6.1 Affective Feature Evaluations

Figure 5 depicts the KL divergence (KLD) of the affective features across distributions gathered on the frames with faces and without faces. A higher KL divergence indicates a bigger discrimination for the distributions of the feature dimension across the with/without face classes. In Figure 5, the first two KLD values are for the activation and valence attributes, and the later values are color coded for the dimensions of  $g^A$  and  $g^V$ . Note that all the dimensions of the  $g^A$ , except the fourth, exhibit higher KLD values than the activation attribute, and at certain dimensions KLD is almost 4 times higher than the KLD of the activation. A similar trend can be observed for the valence embedding vector  $g^V$ , where five dimensions exhibit higher KLD values than the valence attribute, and the largest KLD is extracted for the 6th dimension of the  $g^V$ . These higher KLD values for the GRU based feature dimensions can be observed as the discriminative cues for the human-centric videos that also contribute to the performance of the proposed AVSUM architectures.

### 4.6.2 Video Level Evaluations

Figure 6 depicts the performance comparison of the AVSUM models with the  $\Delta F1_L$  metric for the  $Top-30$  human-centric videos and yields valuable insights. Note that positive accumulation of  $\Delta F1_L$  metric indicates a better performance than the baseline SUM-FCN model. In Figure 6, the AVSUM-SCAV, TA-AVSUM and SA-AVSUM models perform better than the baseline till the  $Top-15$  human-centric videos, whereas AVSUM-GRU model sustains a higher performance

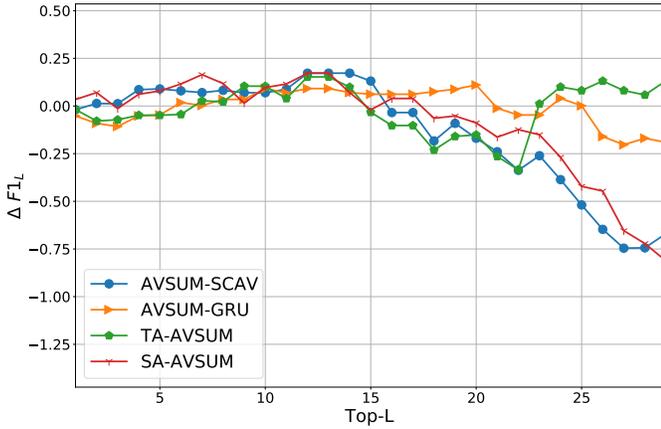


Fig. 6. Performance comparison of the AVSUM models with the  $\Delta F1_L$  metric for the *Top-30* human-centric videos

till the *Top-20* human-centric videos. Similar trends of the AVSUM-SCAV and SA-AVSUM  $\Delta F1_L$  performances can be due to the common late fusion mechanism in these architectures. As  $L$  increases, we can assume human-centric characteristic of the videos is getting weaker. Hence both AVSUM-SCAV and SA-AVSUM and as well as AVSUM-GRU models are observed to performed better for the top human-centric videos and their performances start degrading as human-centric characteristic is getting weaker.

We also investigate video level F-score performances for the proposed AVSUM models. Figure 7 presents scatter plot of the video level  $F1$  scores on the left column and video level  $R$  scores on the right column, where the diagonal in each figure represents similar performances of the compared models. Furthermore, the *Top-15* human-centric videos are color coded in blue to observe their comparative performances.

Video level F-score performances of the AVSUM-GRU vs SUM-FCN tend to cluster around the diagonal that makes these two models to be most similar in terms of F-score performance. Furthermore, majority of the *Top-15* human-centric videos are on or above the diagonal, which indicates a stronger performance for the human-centric videos. Unlike, F1-score, face recall performance comparison depicted by Figure 7(b) has a scattered behavior providing improvements to some videos while degrading other. However, it is seen that majority of the *Top-15* human-centric videos are affected positively. This observation is in line with the  $F1_{15}$  performance of AVSUM-GRU.

Video level F-score performances of the AVSUM-SCAV vs SUM-FCN have a higher deviation from the diagonal. However, almost all the *Top-15* human-centric videos are on or above the diagonal. This indicates a strong performance improvement for the human-centric videos. In terms of face recall on Figure 7(d), majority of the videos are accumulated around the diagonal, indicating that AVSUM-SCAV and SUM-FCN are most similar in terms of face recall.

Video level F-score performances of the TA-AVSUM and SA-AVSUM models have also high deviation from the diagonal. However, like AVSUM-SCAV majority of the *Top-15* human-centric videos are on or above the diagonal for both. Similar to F-score, face recall comparisons of TA-AVSUM

and SA-AVSUM have a scattered behavior depicted by Figure 7(f), and 7(h). However, different than SA-AVSUM, scattered points accumulated on the positive side for TA-AVSUM, stating a major performance improvement which is inline with its best achieving  $R_{15}$  performance for the Max  $R$  criterion.

## 5 CONCLUSION

In this study, we proposed a new affective information enriched end-to-end video summarization framework for human-centric videos. As a first step, we modeled affective information in terms of AV attributes and GRU embeddings, which were extracted from the CER models. The CER-NET, a CER model achieving state-of-the-art CCC performance, was introduced. We explored the use of affective information with the proposed AVSUM-SCAV and AVSUM-GRU fusion models and attention mechanisms based TA-AVSUM and SA-AVSUM models. Experimental investigations of the proposed models were conducted on the RECOLA and TvSum datasets.

We observed that with the fusion of affective information, F-score performance of the video summarization on the human-centric videos can be improved. To further analyze the effect of injected features we defined a face recall ( $R$ ) metric and showed that AVSUM-GRU and AVSUM-SCAV models outperform SUM-FCN with more than 1% increase in face recall  $R$ . The AVSUM-GRU model has strong performance improvements for the human-centric videos on  $F1$  score and face recall  $R$ , and as well it is the most competitive model with the baseline SUM-FCN. On the other hand, we observed that attention enhanced mechanisms exhibits strong performance gains with the Max  $R$  criterion for the human-centric videos. We should also note that affective GRU embedding features exhibit higher KLD across with-face and without-face frames.

Comparing the proposed AVSUM models, AVSUM-GRU has a consistent and competitive performance regardless of the model selection criterion at  $F1$  and  $R$  metrics and has a good balance between the performance improvement at *Top-15* and the performance degradation at the remaining videos. On the other hand, temporal attention based TA-AVSUM performs competitive with the Max  $F1$  criterion and attains strong improvement with the Max  $R$  criterion. The proposed AVSUM models integrate affective information to the summarization architectures and attain important video summarization improvements for the human-centric videos. Compilation of affective human-centric video datasets for video summarization tasks stays as a critical and valuable future study. As a future work, we would like to collect a dataset which is labeled for both video summarization and CER. To extend current study, we would like to investigate multi-modal architectures for both CER and human-centric video summarization components.

## REFERENCES

- [1] B. Köprü and E. Erzin, "Multimodal continuous emotion recognition using deep multi-task learning with correlation loss," *arXiv:2011.00876*, 2020.
- [2] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris, and I. Patras, "Video summarization using deep neural networks: A survey," *arXiv:2101.06072v1*, 2021.

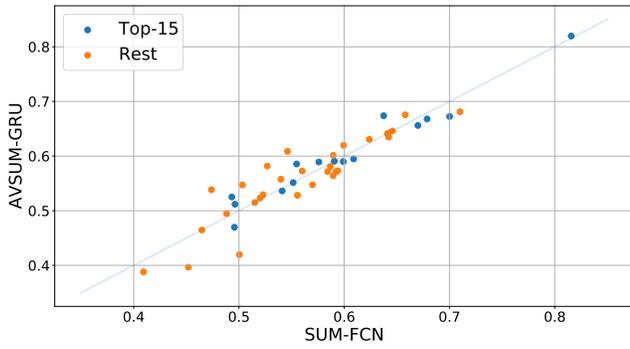
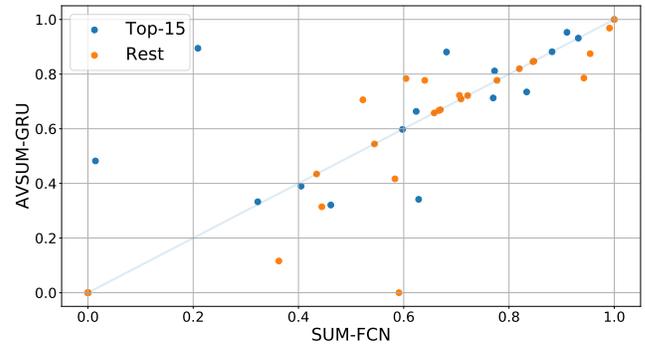
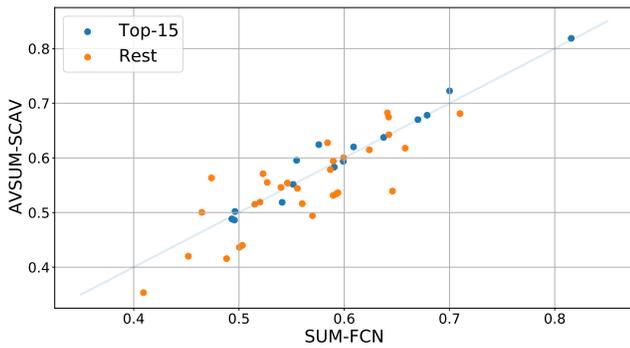
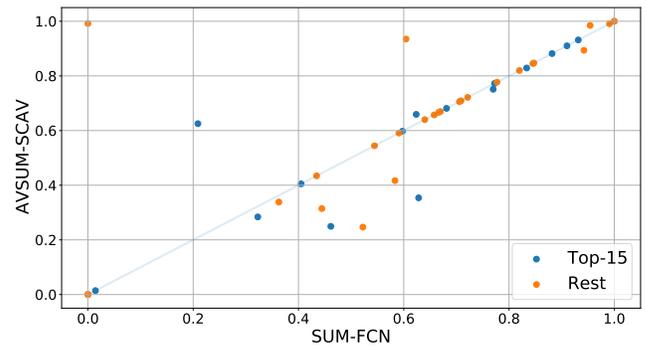
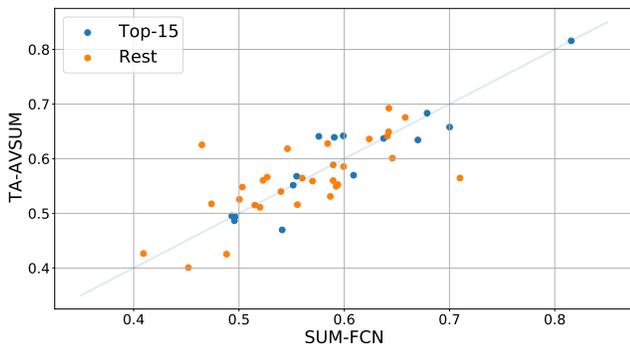
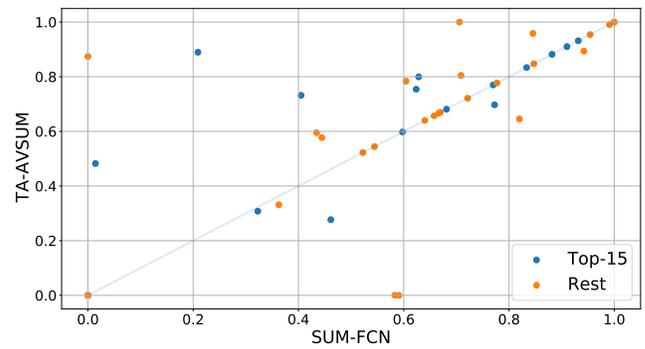
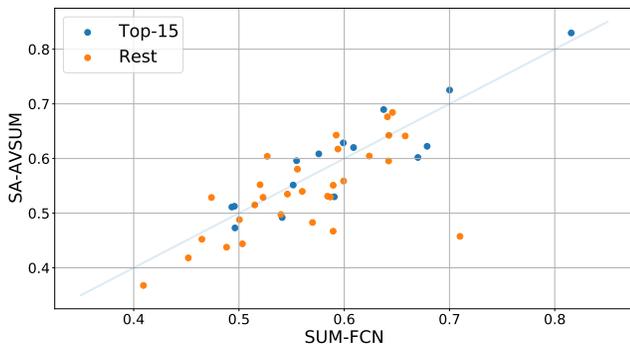
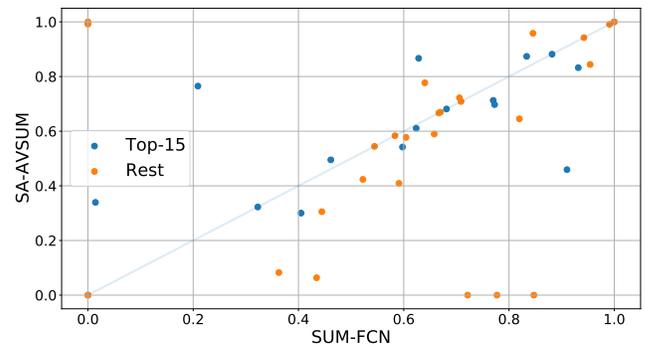
(a)  $F1$  scores for AVSUM-GRU vs SUM-FCN(b)  $R$  scores for AVSUM-GRU vs SUM-FCN(c)  $F1$  scores for AVSUM-SCAV vs SUM-FCN(d)  $R$  scores for AVSUM-SCAV vs SUM-FCN(e)  $F1$  scores for TA-AVSUM vs SUM-FCN(f)  $R$  scores for TA-AVSUM vs SUM-FCN(g)  $F1$  scores for SA-AVSUM vs SUM-FCN(h)  $R$  scores for SA-AVSUM vs SUM-FCN

Fig. 7. Video level  $F1$  score and  $R$  score comparisons of the AVSUM models against the baseline SUM-FCN:  $F1$  scores are with Max  $F1$  criterion,  $R$  scores are with the Max  $R$  criterion, and the *Top-15* face including videos are color coded in blue.

- [3] J. S. Lerner, Y. Li, P. Valdesolo, and K. S. Kassam, "Emotion and decision making," *Annual Review of Psychology*, vol. 66, pp. 799–823, jan 2015.
- [4] J. A. Russell and A. Mehrabian, "Evidence for a three-factor theory of emotions," *Journal of Research in Personality*, vol. 11, no. 3, pp. 273–294, 1977.
- [5] H. Schlosberg, "Three dimensions of emotion," *Psychological Review*, vol. 61, no. 2, pp. 81–88, 1954.
- [6] G. Trigeorgis, F. Ringeval, R. Brueckner, E. Marchi, M. A. Nicolaou, B. Schuller, and S. Zafeiriou, "Adieu features? end-to-end speech emotion recognition using a deep convolutional recurrent network," in *ICASSP 2016, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016, pp. 5200–5204.
- [7] M. Neumann and N. T. Vu, "Attentive convolutional neural network based speech emotion recognition: A study on the impact of input features, signal length, and acted speech," in *INTERSPEECH 2017, Annual Conference of the International Speech Communication Association*, 2017, pp. 1263–1267.
- [8] Y. Li, T. Zhao, and T. Kawahara, "Improved end-to-end speech emotion recognition using self attention mechanism and multitask learning," in *INTERSPEECH 2019, Annual Conference of the International Speech Communication Association*, 2019, pp. 2803–2807.
- [9] S. Minaee and A. Abdolrashidi, "Deep-Emotion: Facial expression recognition using attentional convolutional network," *arXiv:1902.01019*, 2019.
- [10] D. Aspandj, A. Mallol-Ragolta, B. Schuller, and X. Binefa, "Adversarial-based neural networks for affect estimations in the wild," *arXiv:2002.00883*, 2020.
- [11] D. Kollias, M. A. Nicolaou, I. Kotsia, G. Zhao, and S. Zafeiriou, "Recognition of affect in the wild using deep neural networks," in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 1972–1979.
- [12] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, "Video summarization with long short-term memory," in *ECCV 2016, European Conference on Computer Vision*. Springer International Publishing, 2016, pp. 766–782.
- [13] Y.-F. Ma, L. Lu, H.-J. Zhang, and M. Li, "A user attention model for video summarization," in *MULTIMEDIA '02, Tenth ACM International Conference on Multimedia*, 2002, pp. 533–542.
- [14] G. Kim and E. P. Xing, "Reconstructing storyline graphs for image recommendation from web community photos," in *CVPR 2014, IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 3882–3889.
- [15] P. Mundur, Y. Rao, and Y. Yesha, "Keyframe-based video summarization using Delaunay clustering," *International Journal on Digital Libraries*, vol. 6, no. 2, pp. 219–232, apr 2006.
- [16] C. W. Ngo, T. C. Pong, H. J. Zhang, and R. T. Chin, "Motion-based video representation for scene change detection," in *ICPR 2000, International Conference on Pattern Recognition*, vol. 15, no. 1, 2000, pp. 827–830.
- [17] C. W. Ngo, Y. F. Ma, and H. J. Zhang, "Automatic video summarization by graph modeling," in *ICCV 2003, IEEE International Conference on Computer Vision*, vol. 1, 2003, pp. 104–109.
- [18] Z. Cernekova, I. Pitas, and C. Nikou, "Information theory-based shot cut/fade detection and video summarization," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 16, no. 1, pp. 82–91, jan 2006.
- [19] B. Mahasseni, M. Lam, and S. Todorovic, "Unsupervised video summarization with adversarial LSTM networks," in *CVPR 2017, 30th IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2982–2991.
- [20] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris, and I. Patras, "Unsupervised video summarization via attention-driven adversarial learning," in *MultiMedia Modeling*, Y. M. Ro, W.-H. Cheng, J. Kim, W.-T. Chu, P. Cui, J.-W. Choi, M.-C. Hu, and W. De Neve, Eds. Springer International Publishing, 2020, pp. 492–504.
- [21] X. He, Y. Hua, T. Song, Z. Zhang, Z. Xue, R. Ma, N. Robertson, and H. Guan, "Unsupervised video summarization with attentive conditional generative adversarial networks," in *Multimedia 2019, 27th ACM International Conference on Multimedia*, 2019, p. 2296–2304.
- [22] M. Rochan, L. Ye, and Y. Wang, "Video summarization using fully convolutional sequence networks," in *ECCV 2018, European Conference on Computer Vision*, 2018, pp. 358–374.
- [23] Z. Ji, K. Xiong, Y. Pang, and X. Li, "Video summarization with attention-based encoder-decoder networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 6, pp. 1709–1717, 2020.
- [24] P. Li, Q. Ye, L. Zhang, L. Yuan, X. Xu, and L. Shao, "Exploring global diverse attention via pairwise temporal relation for video summarization," *Pattern Recognition*, vol. 111, p. 107677, 2021.
- [25] H. Joho, J. Staiano, N. Sebe, and J. M. Jose, "Looking at the viewer: Analysing facial activity to detect personal highlights of multimedia contents," *Multimedia Tools and Applications*, vol. 51, no. 2, pp. 505–523, jan 2011.
- [26] A. G. Money and H. Agius, "Analysing user physiological responses for affective video summarisation," *Displays*, vol. 30, no. 2, pp. 59–70, apr 2009.
- [27] J. Zhao, X. Mao, and L. Chen, "Speech emotion recognition using deep 1D & 2D CNN LSTM networks," *Biomedical Signal Processing and Control*, vol. 47, pp. 312–323, jan 2019.
- [28] S. Mirsamadi, E. Barsoum, and C. Zhang, "Automatic speech emotion recognition using recurrent neural networks with local attention," in *ICASSP 2017, IEEE International Conference on Acoustics, Speech and Signal Processing*, 2017, pp. 2227–2231.
- [29] H. Meng, T. Yan, F. Yuan, and H. Wei, "Speech emotion recognition from 3D log-mel spectrograms with deep learning network," *IEEE Access*, vol. 7, pp. 125 868–125 881, 2019.
- [30] M. Schmitt, N. Cummins, and B. W. Schuller, "Continuous emotion recognition in speech — do we need recurrence?" in *INTERSPEECH 2019, Annual Conference of the International Speech Communication Association*. ISCA, 2019, pp. 2808–2812.
- [31] P. Tzirakis, J. Chen, S. Zafeiriou, and B. Schuller, "End-to-end multimodal affect recognition in real-world environments," *Information Fusion*, vol. 68, pp. 46–53, apr 2021.
- [32] B. Xu, Y. Fu, Y. G. Jiang, B. Li, and L. Sigal, "Heterogeneous knowledge transfer in video emotion recognition, attribution and summarization," *IEEE Transactions on Affective Computing*, vol. 9, no. 2, pp. 255–270, apr 2018.
- [33] G. Tu, Y. Fu, B. Li, J. Gao, Y. G. Jiang, and X. Xue, "A Multi-Task Neural Approach for Emotion Attribution, Classification, and Summarization," *IEEE Transactions on Multimedia*, vol. 22, no. 1, pp. 148–159, 2020.
- [34] Y. Jung, D. Cho, D. Kim, S. Woo, and I. S. Kweon, "Discriminative feature learning for unsupervised video summarization," in *33rd AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 8537–8544.
- [35] B. Zhao, X. Li, and X. Lu, "Property-constrained dual learning for video summarization," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 10, pp. 3989–4000, oct 2020.
- [36] K. Zhou, Y. Qiao, and T. Xiang, "Deep reinforcement learning for unsupervised video summarization with diversity-representativeness reward," in *AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [37] Y. Yuan, H. Li, and Q. Wang, "Spatiotemporal modeling for video summarization using convolutional recurrent neural network," *IEEE Access*, vol. 7, pp. 64 676–64 685, 2019.
- [38] D. Bahdanau, K. H. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *ICLR 2015, 3rd International Conference on Learning Representations*, 2015.
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv:1706.03762*, 2017.
- [40] Z. Ji, F. Jiao, Y. Pang, and L. Shao, "Deep attentive and semantic preserving video summarization," *Neurocomputing*, vol. 405, pp. 200–207, sep 2020.
- [41] Y. Song, J. Vallmitjana, A. Stent, and A. Jaimes, "TVSum: Summarizing web videos using titles," in *CVPR 2015, IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5179–5187.
- [42] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR 2015, IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [43] J. Long, E. Shelhamer, and T. Darrell, "Fully Convolutional Networks for Semantic Segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 640–651, nov 2014.
- [44] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR 2018, IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.
- [45] F. Ringeval, A. Sonderegger, J. Sauer, and D. Lalanne, "Introducing the RECOLA multimodal corpus of remote collaborative and af-

fective interactions,” in *FG 2013, 10th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition*, 2013.

- [46] D. E. King, “Dlib-ml: A machine learning toolkit,” *Journal of Machine Learning Research*, vol. 10, no. 60, pp. 1755–1758, 2009.
- [47] P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. W. Schuller, and S. Zafeiriou, “End-to-end multimodal emotion recognition using deep neural networks,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, 2017.



**Berkay Köprü** is a Ph.D. student at Koc University, Istanbul, Turkey. He received his M.Sc. degree and B.Sc. degree from Technical University of Munich, Munich, Germany in 2017, and Bilkent University, Ankara, Turkey in 2014 respectively. His research interest include human-computer interaction, computer vision and affective computing.



**Engin Erzin** (S'88-M'96-SM'06) received his Ph.D. degree, M.Sc. degree, and B.Sc. degree from the Bilkent University, Ankara, Turkey, in 1995, 1992 and 1990, respectively, all in Electrical Engineering. During 1995-1996, he was a postdoctoral fellow in Signal Compression Laboratory, University of California, Santa Barbara. He joined Lucent Technologies in September 1996, and he was with the Consumer Products for one year as a Member of Technical Staff of the Global Wireless Products Group. From 1997 to 2001, he was with the Speech and Audio Technology Group of the Network Wireless Systems. Since January 2001, he has been with the Electrical & Electronics Engineering and Computer Engineering Departments of Koc University, Istanbul, Turkey. Engin Erzin is currently a member of the IEEE Speech and Language Processing Technical Committee and Associate Editor for the IEEE Transactions on Multimedia, having previously served as Associate Editor of the IEEE Transactions on Audio, Speech & Language Processing (2010-2014). His research interests include speech-audio-visual signal processing, affective computing, human-computer interaction and machine learning.