

# Video Summarization through Reinforcement Learning with a 3D Spatio-Temporal U-Net

Tianrui Liu, Qingjie Meng, Jun-Jie Huang, Athanasios Vlontzos, Daniel Rueckert, *Fellow, IEEE*,  
and Bernhard Kainz *Senior member, IEEE*,

**Abstract**—Intelligent video summarization algorithms allow to quickly convey the most relevant information in videos through the identification of the most essential and explanatory content while removing redundant video frames. In this paper, we introduce the 3DST-UNet-RL framework for video summarization. A 3D spatio-temporal U-Net is used to efficiently encode spatio-temporal information of the input videos for downstream reinforcement learning (RL). An RL agent learns from spatio-temporal latent scores and predicts actions for keeping or rejecting a video frame in a video summary. We investigate if real/inflated 3D spatio-temporal CNN features are better suited to learn representations from videos than commonly used 2D image features. Our framework can operate in both, a fully unsupervised mode and a supervised training mode. We analyse the impact of prescribed summary lengths and show experimental evidence for the effectiveness of 3DST-UNet-RL on two commonly used general video summarization benchmarks. We also applied our method on a medical video summarization task. The proposed video summarization method has the potential to save storage costs of ultrasound screening videos as well as to increase efficiency when browsing patient video data during retrospective analysis or audit without losing essential information.

**Index Terms**—Video summarization, Reinforcement learning, 3D convolutions, 3D U-Net, Medical video processing, Ultrasound.

## I. INTRODUCTION

**T**he amount of video data, including videos shared across social media, has been growing at an exponential rate. It is reported that for example approximately 300 hours of videos are uploaded to YouTube every single minute. This amount of information is impossible to process for humans. Currently, searching for desired content is mainly relied on textual key words and video thumbnails, thus search times often exceed genuine video utilization. This inhibits consumers and content creators alike.

There has been a steadily growing interest in machine learning techniques for video summarization. Early works adopt low-level or mid-level visual features to locate important segments of a video with a particular strategy such as clustering [1], [2] and sparse dictionary learning [3], [4]. Other methods rely on heuristic frame sampling mechanisms like static video storyboard summaries [5], [6] or frame-similarity-based dynamic video skimming [7]. However, most of these methods are not capable to model a temporal component that can provide additional contextual information for enhanced discriminative power [8], [9].

Incorporating recurrent neural networks (RNNs) for sequence modeling has led to significant progresses in this

field [10]–[13]. Long short term memory (LSTM)-based methods using bi-directional LSTMs [10] as well as hierarchical LSTMs [12], [13] have been proposed for sequence modeling in video summarization. Recently, Rochan *et al.* [14] have demonstrated that it is possible to model the video summarization task as an element-wise segmentation problem using fully convolutional networks (FCNs). One of the most common FCN architectures are U-Nets [15], which have originally been proposed for image segmentation tasks where they exhibit good semantic image feature extraction abilities in the spatial domain [16]. For video summarization, Rochan *et al.* proposed to utilize fully convolutional sequential network (FCSN) [14] in the temporal domain for sequential modeling. As a segmentation tool, FCSN [14] works for the video summarization problem by determining which parts of the video should be “segmentated” as important. Compared to methods using recurrent models [10], [11], [13], FCSNs have the advantage of operating on the whole video, which provides the maximum amount of context and allows improved GPU utilization. Encouraged by the success of FCSN [14], we propose a new 3D U-Net-based architecture to model both the spatial and the temporal dependencies among video frames. While the FCSN in [14] takes 2D CNN features for each single frame and applies 1-dimensional (1D) convolutions along the temporal dimension for sequential modeling, we hypothesize that spatio-temporal features as well as full 3D convolutions are better suited for video data.

In [17]–[19], spatio-temporal CNN features have been used to improve the video action recognition task. However, to the best of our knowledge, a 3D spatio-temporal U-Net architecture using spatio-temporal features exploited for video summarization remains a rarely addressed problem.

In this paper, we propose a 3D spatio-temporal U-Net (3DST-UNet) using spatio-temporal features for video summarization. We accommodate the paradigm of Reinforcement Learning (RL) and exploit the combination of RL and 3DST-UNet. Our 3DST-UNet directly links with spatio-temporal feature extraction networks by taking as input 4-dimensional (4D) video features, encoding both spatial and temporal video information. A feature extraction CNN outputs spatio-temporal features of frame sequences and feeds them into the 3DST-UNet. The proposed 3DST-UNet exploits contextual information and maps the spatio-temporal video features effectively into a continuous latent space. Both the spatio-temporal features and the 3DST-UNet have the potential to encode spatio-temporal dependencies between video frames. The role of the RL agent is to learn policies which can maximize the

accumulated reward of taking actions, where the actions are defined as whether to select or discard the current frame as a key frame. The proposed 3DST-UNet model can be trained either in a supervised way or a fully unsupervised way, that is, the training process of our summarization network accommodates cases where little to no human annotations are available.

To evaluate the effectiveness of the proposed method, we conduct comprehensive experiments using different training settings on public video summary datasets, *i.e.*, SumME [2], TVSum [20] as well as OVP and YouTube [21], [22]. Furthermore, we extend our video summarization method for medical videos and take fetal screening with ultrasound imaging as an example. We demonstrate that our method achieves better performance than previous approaches proposed for video summarization on both general videos and medical videos.

Overall, the contribution of this paper is as follows:

(1) We propose 3DST-UNet-RL, a deep RL-based framework using a 3D spatio-temporal U-Net (3DST-UNet) for video summarization. To the best of our knowledge, this is the first approach to use a 3D U-Net with spatial-temporal features for video summarization.

(2) We introduce a 3D spatio-temporal U-Net (3DST-UNet) for sequential spatio-temporal video feature modeling. Compared RNN-based models, our 3DST-UNet directly links with spatio-temporal CNN features to exploit both spatial and temporal information.

(3) We provide comprehensive experiments and demonstrate that our method achieves better performance than previous approaches on two commonly used video summarization benchmarks. We also discussed the limitations of the current evaluation criteria for general video summarization and evaluate the video summary under various length constrains.

(4) We extend our video summarization method for ultrasound scanning videos, where automated report generation with short but relevant video clips is desirable. We show experimental evidence for the effectiveness of our approach in medical video data from the clinical practice.

The rest of the paper is organized as follows: Section II gives a review on the related works, Section III introduces the proposed 3DST-UNet-RL framework for video summarization, Section IV shows the ablation studies of the proposed method and compares with other video summarization methods, and finally Section V concludes the paper.

## II. RELATED WORKS

### A. Video Summarization

Video summarization methods have been explored in literature using either supervised learning or unsupervised learning. Supervised methods learn video summarization from manually labeled data consisting of videos and their corresponding user annotations as ground-truth, usually key frames that are subjectively perceived as important. In [10], a LSTM-based key frame selection model is trained by minimizing the cross-entropy loss between the estimated key frames and the user annotated ground-truth key frames. In order to ensure diversity of the selected frames, an additional objective based on a

determinantal point process (DPP) is used. Zhang *et al.* [23] combine sequential models for the summary creation with a retrospective encoder, which maps the summaries to an abstract latent space. The work of [24] learns a mapping function from a set of web videos to a set of summary videos via an adversarial process. There are also video summarization method using Generative Adversarial Network (GAN). In [11], a subset of representative key frames is selected by training a summarizer to minimize the distances between videos and a distribution of their summaries using generative adversarial networks as critics. Similarly, [25] aims to maximize the mutual information between a summary and video using an information-preserving metric, two trainable discriminators and a cycle consistent adversarial learning objective.

While deep neural networks (DNNs) have achieved significant improvements for video summarization performance, they impose a heavy burden on algorithm designers to collect a huge amount of labeled data for fully supervised DNN models. Therefore, unsupervised DNN methods are attractive and there have been pioneering works showing promising results. Some summarization methods provide weak supervision through additional cues such as images and videos from the web [26]–[28] and their accompanying category information [29], [30] to improve performance. In [28], a variational auto-encoder is applied for learning the latent semantics from web videos. In addition, an attention network for saliency estimation is used to improve the performance.

### B. Reinforcement Learning for Video Analysis

The goal of reinforcement learning (RL) is to learn a good policy for the agent from experimental trials by maximizing expected future rewards. Reinforcement Learning (RL) has been succeeded in solving various vision tasks, such as visual tracking [31], video face recognition [32], video captioning [33], and video object segmentation [34].

Recently, deep reinforcement learning (RL) has been explored for video summarization and videos fast-forwarding since these two tasks fit the narrative of RL well. Zhou and Qiao [35] propose a deep summarization network, which formulates the video summarization task as a sequential decision making process. The summaries are generated by predicting the probabilities of a given frame being a key-frame. The summary frames are then sampled based on this probability. In [30], a Q-learning-based summarization network is explored to guide an artificial RL agent to use a recognizability reward. The reward is derived from a category classification network and the summarization network relies on video-level category labels to address the summarization problem in a weakly supervised manner [30]. In [36], [37] RL has been used to fast-forward lengthy videos. A FFNet (*i.e.* FastForward Network) is proposed in [37] based on a Markov decision process to automatically decide the number of frames required for efficient fast-forwarding.

RL has also been applied in the field of medical image and video analysis. In [38], [39], RL strategies are successfully combined with spatio-temporal frame histories for video landmark detection. In [40], an RL based video analyse

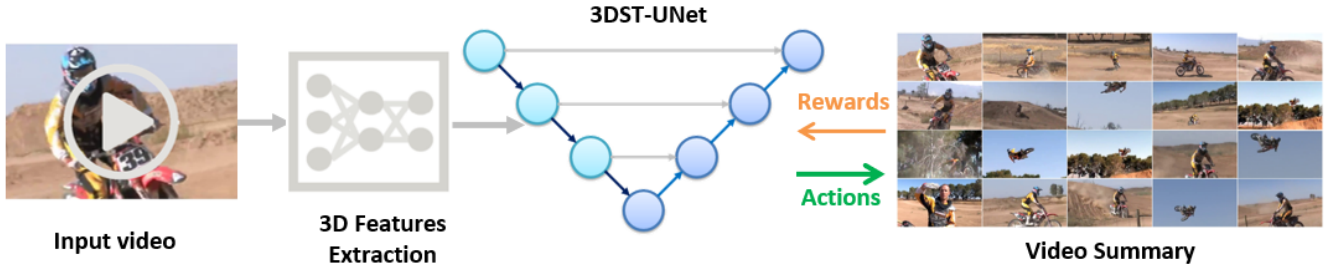


Fig. 1. Overview of the proposed 3DST-UNet-RL method for video summarization. The framework consists of three main parts: a spatio-temporal CNN with 3D convolution for video features extraction, a 3D spatio-temporal U-Net (3DST-UNet) for sequence modeling, and an RL network which takes actions to select key frames for the summary set. Given an input video, the 3D feature extraction network first computes spatio-temporal feature and feeds them into a spatio-temporal U-Net for sequence modeling. Following, we make use of reinforcement learning to predict actions keeping/rejecting a frame for inclusion in a video summary.

method is used to preserve essential information in ultrasound diagnostic videos. A diagnostic plane detection reward has been proposed which guides the agent to learn according to the clinical diagnostic standards. The features representing the video frames in [40] are extracted from 2D CNN networks and are further modelled using LSTM. Differently, we use the proposed 3DST-UNet to model the 3D CNN features for video summarization.

### C. 3D U-Net

U-Net [15] was originally proposed based on FCN for the biomedical image segmentation. Thereafter, the U-Net-like structures have been widely used in the field of image segmentation [15]. U-Net resembles a fully convolutional network (FCN) structure and uses deconvolution to restore image size and feature. Different from FCN which the fusion operation during upsampling is direct feature addition, the U-Net upsampling process uses the concatenate operation to splicing the feature maps. Following concatenation is the feature map deconvolution. The skip connection strategy directly utilizing shallow features.

Several studies [41], [42] have demonstrated that a 3D version of U-Net can produce better results than the 2D architectures. Çiçek et al. [41] proposed a 3D U-Net model that generates dense volumetric segmentations. It realizes 3D image segmentation by inputting a continuous 2D slice sequence of 3D images. The network structure is similar to U-Net, with one encoding path and one decoding path, each has four resolution levels. The encoder gradually reduces the spatial dimension by continuously merging the layers to extract feature information, and the decoder portion gradually restores the target detail and the spatial dimension.

Our proposed architecture is based on UNet and is influenced by the fully convolutional sequential network in [14] which uses 1-dimensional (1d) convolutional along the temporal dimension for sequential modeling. We extend the idea of (FCSN) [14] and design a 3D temporal U-Net (3DST-UNet) that takes spatio-temporal video features as input from a preceding spatio-temporal 3D CNN.

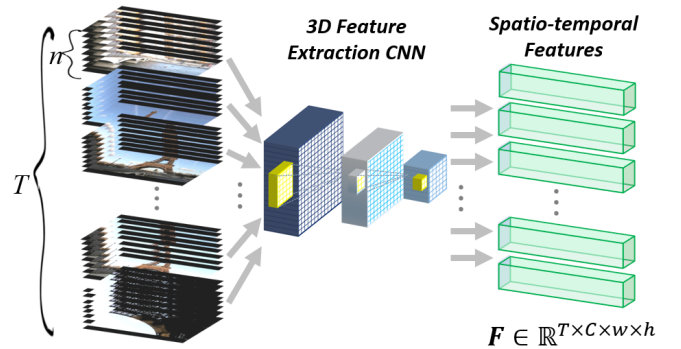


Fig. 2. Sequential spatio-temporal video 3D feature extraction network. The input video of length  $L = T \times n$  has been divided into  $T$  segments of length  $n$ . The spatio-temporal video 3D feature extraction network takes  $n$  frames at a time to compute video features of size  $C \times w \times h$  and results in video feature  $F \in \mathbb{R}^{T \times C \times w \times h}$ .

## III. METHOD

### A. Overview

As illustrated in Fig. 1, our video summarization framework consists of three main parts: a spatio-temporal CNN with 3D convolution for video feature extraction, a 3D Spatio-Temporal U-Net (3DST-UNet), and a RL agent's network. The video feature extraction network outputs spatio-temporal features from a sequence of video frames. These features are then combined and fed into the 3DST-UNet, which can further model both the spatial and temporal relationships within the video sequence. The 3DST-UNet is followed by a Sigmoid layer to produce scores for each frame of the input video. Based on these scores, the RL agent's network takes actions to decide whether or not to select a frame for the summary.

### B. Video Feature Extraction Network

CNNs utilizing spatio-temporal convolutions are expected to suit video data better than 2D features because of their potential to directly encode spatio-temporal dependencies between frames. We investigate two options for video feature representation, *i.e.*, spatio-temporal 3D features (ST3D) [18] and inflated 3D (I3D) features [19]. Specifically, we remove the last two fully connect layers of the backbone networks

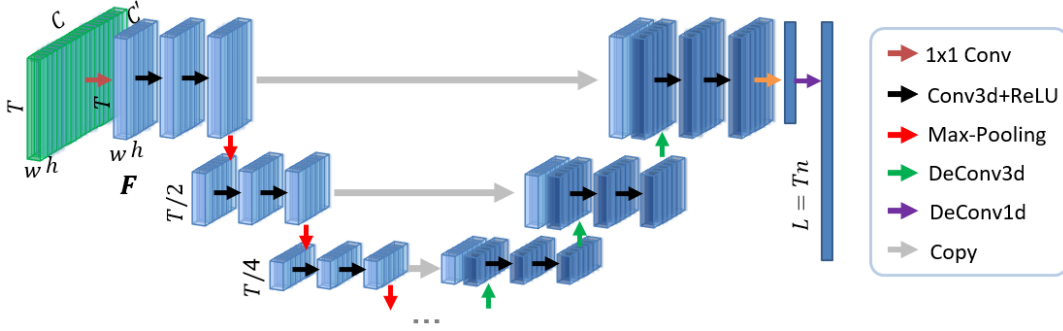


Fig. 3. The proposed 3D Spatio-Temporal UNet (3DST-UNet). It takes 4-dimensional video feature  $\mathbf{F} \in \mathbb{R}^{T \times C \times w \times h}$  as input where  $C$ ,  $w$ ,  $h$  denotes the temporal, channel, spatial width, and spatial height dimensions of the feature, respectively. 3DST-UNet gradually extracts features at different temporal resolutions in the encoder part which are then successively recombined in the decoder part using upsampling operations to propagate context information to higher temporal resolution layers.

[43]. In this way, the extracted video feature representation is kept with the spatial dimensions  $w$  and  $h$ .

1) *Spatio-temporal 3D Features*: As a direct solution, Spatio-temporal 3D (ST3D) features can generate spatio-temporal feature representations by allocating the third dimension of the input tensor as the temporal component. ST3Ds use 3D convolutional layers and have been proposed for video action recognition tasks [17]–[19]. The preceding spatio-temporal CNN consists of three 3D convolutional layers which takes  $n$  frames at a time to compute  $T = L/n$  spatio-temporal feature vectors from the total  $L$  frames for every  $n$  neighboring frames.

2) *Inflated 3D Features*: The performance of high-level computer vision tasks using real 3D convolutions can be limited when the pre-trained dataset is not large enough. In [19], Inflated 3D features (I3D) are proposed to make use of 2D CNNs with ImageNet [44] pre-trained weights. I3D duplicates the 2D kernels along the axial direction to produce 3D kernels. We accommodate the inflated 3D convolution network to efficiently encoding the 3D video sequences. This allows our summarization framework to fully exploit 3D context information while re-purposing off-the-shelf deep 2D network structures and inheriting their large capacities to cope with image variances. Similar to the ST3D features, the convolutional layers compute  $T = L/n$  feature vectors for every  $n$  neighboring frames.

Hence, for both the ST3D and I3D features, the video feature representation is represented as  $\mathbf{F} \in \mathbb{R}^{T \times C \times w \times h}$  where  $C$ ,  $w$ ,  $h$  denotes the feature temporal, channel, spatial width, and spatial height dimensions, respectively. The block diagram of the 3D feature extraction network can be seen in Fig. 2.

### C. 3DST-UNet for Sequential Spatio-temporal CNN Feature Modeling

We propose to use a 3D Spatio-Temporal U-Net (3DST-UNet) to model the spatio-temporal video feature representation  $\mathbf{F}$  in both, the temporal and the spatial domain for video summarization. The 3DST-UNet aims to generate latent scores for each video frame and takes features which encode both the spatial and the temporal video information as input. As such,

the proposed 3DST-UNet is directly linked with the spatio-temporal feature extraction network for sequential modeling.

Different from the ordinary U-Net [15] for semantic segmentation in the spatial domain [16], the proposed 3D spatio-temporal U-Net (3DST-UNet) is proposed to be used for sequential modeling in both spatial and temporal domain. Also different from the fully convolutional sequential network (FCSN) which uses 1-dimensional convolutional along the temporal dimension, 3DST-UNet takes spatio-temporal video features from a preceding spatio-temporal 3D CNN as input. The 3DST-UNet for sequential modeling has the advantage of processing the entire video sequence at once. It does not employ recurrent or sequential processing, therefore can be driven in a single forward/backward pass during inference/training for sequences with variable length. Another merit of 3DST-UNet is that it is able to extract context information from different temporal resolutions and exploit the spatio-temporal information which are very helpful for video summarization.

The proposed 3DST-UNet is an end-to-end volumetric architecture consisting of an encoder, decoder stage, and upsampling stage for video frame latent score generation. As illustrated in Fig. 3, the 3DST-UNet encoder generates features at different temporal resolutions and a decoder which successively fuses multi-resolution features and produces latent scores for each frame. Given the 4-dimensional (4D) video feature sequence of dimension  $T \times C \times w \times h$ , the encoder part of 3DST-UNet takes the 4D video feature  $\mathbf{F} \in \mathbb{R}^{T \times C \times w \times h}$  as input; the decoder part outputs  $L = T \times n$  latent scores where  $L$  equals to the total number of frames in the input video,  $T$  is the total number of spatio-temporal features fed into the 3DST-UNet.

1) *Encoder*: We introduced squeeze-and-excitation blocks [45], [46] to the 3D U-Net architecture. The 4D features maps (denoted as green box in Fig. 3) are first compressed by a squeeze layer to compress the features through  $1 \times 1$  convolutions. The channel-wise dimension  $C$  is thereafter squeezed from  $C$  to  $C'$ , i.e.,

$$\hat{\mathbf{f}}_i = \mathbf{v}_i * \mathbf{F}, \quad (1)$$

where  $\hat{\mathbf{f}}_i$  is the  $i$ -th feature element of the output feature map  $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_{C_{in}}] \in \mathbb{R}^{T \times H \times W \times C_{in}}$ ,  $\mathbf{v}_i$  is the  $i$ -th learned

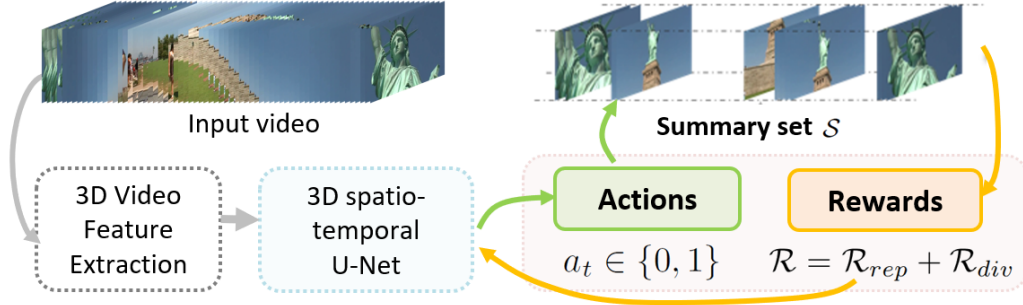


Fig. 4. Reinforcement learning on the latent scores produced by the 3DST-UNet. With the latent scores from the 3DST-UNet, the reinforcement learning network selects actions  $a_t$  on whether a frame should be selected into the summary set  $\mathcal{S}$  or not by maximizing the rewards including the representativeness reward  $\mathcal{R}_{rep}$  and diversity reward  $\mathcal{R}_{div}$ .

filter in the squeeze layers for  $i = 1, \dots, C'$ , and ‘\*’ denotes convolution.

The encoder of 3DST-UNet consists of repeated application of 3D convolution, each followed by a rectified linear unit (ReLU) and a max pooling operation with stride 2 for downsampling. The dimension of feature channels are doubled at each downsampling stage. The 3D convolution operations of kernel size  $3 \times 3 \times 3$  sliding along the channel, width and height dimensions to exploit both the spatial and the temporal context information. The pooling operations are performed along the temporal dimension to enable context information extraction from different temporal resolutions.

2) *Decoder*: Every step in the decoding path consists of an upsampling of the feature map followed by a  $2 \times 2 \times 2$  deconvolution that halves the number of feature channels. 3DST-UNet performs 3D deconvolution along the temporal dimension to upsample the features from lower temporal resolutions. There are skip connections between the encoder and decoder features at the same temporal resolutions to pass information from the encoder part to the decoder part. The concatenated features are then processed using two  $3 \times 3 \times 3$  convolutions (each followed by a ReLU). At the final stage, 3DST-UNet performs global average pooling and 1-dimensional deconvolution to transform the 4D feature into latent scores of length  $L = T \times n$  for the downstream reinforcement learning.

#### D. Reinforcement Learning of Latent Scores for Video Summarization

An RL value network takes as an input the representation of the environment state and the action choice of the agent. In our case, where the summarization task is interpreted as an RL decision making process, the RL takes input the frame latent scores from the 3DST-UNet and an action is taken at every frame on whether to include it in the summary set or not.

Our RL network receives spatio-temporal features and uses fully convolution layers with a Sigmoid activation function to output a probability score of each action. By denoting the spatio-temporal feature sequence as  $\{\mathbf{x}_t\}_{t=1}^L$ , the frame-level probability scores can be defined as  $p_t = \text{Sigmoid}(\mathbf{W} * \mathbf{x}_t + b)$ , where  $\mathbf{W}$  and  $b$  are the parameters of the fully convolution layer.

As illustrated in Fig. 4, the RL agent can take *actions* on whether a frame should be selected in the summary set  $\mathcal{S}$  according to the frame-level probability scores  $p_t$ . For the video summarization problem, the *actions* are binary values, *i.e.*,  $a_t \in \{0, 1\}$  where  $a_t = 1$  indicates the frame  $t$  should be selected in the summary set  $\mathcal{S} = \{s_i | a_{s_i} = 1, i = 1, \dots, |\mathcal{S}|\}$  and vice versa. The frame probability scores  $p_t$  are sampled from a Bernoulli distribution, *i.e.*,  $a_t \sim B(p_t)$ . The objective of the RL summarization network is to maximize the expected rewards:

$$\mathcal{R} = \mathcal{R}_{rep} + \mathcal{R}_{div}, \quad (2)$$

where the representativeness reward  $\mathcal{R}_{rep}$  and the diversity reward  $\mathcal{R}_{div}$  evaluate the quality of the selected summary  $\mathcal{S}$ .

1) *The representativeness reward*:  $\mathcal{R}_{rep}$  measures how well the generated summary can represent the original video. By maximizing  $\mathcal{R}_{rep}$  for the original video, the temporal information across the entire video can be maximally preserved. The degree of representativeness of a video summary is formulated as a  $k$ -medoids problem [47]. The agent is encouraged to select a set of medoids such that the MSE between video frames and their nearest medoids is minimal, *i.e.*,

$$\mathcal{R}_{rep} = \exp\left(-\frac{1}{L} \sum_{t=1}^L \min_{t' \in \mathcal{S}} \|\mathbf{x}_t - \mathbf{x}_{t'}\|_2\right), \quad (3)$$

where  $\mathbf{x}_i$  is the spatio-temporal feature representation of a selected frame in the summary set  $\mathcal{S}$ .

2) *The diversity reward*:  $\mathcal{R}_{div}$  measures the dissimilarity between the selected frames for the summary video. It enforces the agent to select frames with different visual representations into the summary set  $\mathcal{S}$ . Thus, redundancy is kept small in  $\mathcal{S}$ .  $\mathcal{R}_{div}$  is computed as the pairwise frame dissimilarity, *i.e.*,

$$\mathcal{R}_{div} = \frac{1}{|\mathcal{S}|(|\mathcal{S}| - 1)} \sum_{t \in \mathcal{S}} \sum_{\substack{i \in \mathcal{S} \\ i \neq t}} d(\mathbf{x}_t, \mathbf{x}_i), \quad (4)$$

where  $d(\cdot, \cdot)$  is the pair-wise dissimilarity of images in the feature space. In this paper, we adopt the cosine distance for measuring pair-wise dissimilarity.

#### E. Network Training

In our summarization network, the loss terms and reward terms are jointly optimized in an end-to-end manner. In addi-

tion to the introduced award terms, we use two regularization loss terms  $\mathcal{L}_{\text{reg}}^p$  and  $\mathcal{L}_{\text{reg}}^b$  to regularize the properties of the selected summary set.

The proportion regularization loss  $\mathcal{L}_{\text{reg}}^p$  is used to penalize the selection of a large number of frames in the summary set  $\mathcal{S}$ :

$$\mathcal{L}_{\text{reg}}^p = \left\| \frac{1}{L} \sum_{t=1}^L p_t - \epsilon \right\|^2, \quad (5)$$

where  $\epsilon$  is a scalar controlling the proportion of the selected frames.

The binary regularization loss  $\mathcal{L}_{\text{reg}}^b$  is to encourage the learned frame probability to be binary since the annotations in video summarization datasets are given with binary labels:

$$\mathcal{L}_{\text{reg}}^b = \left( \frac{1}{T} \sum_{t=1}^L |p_t - 0.5| \right)^{-1}. \quad (6)$$

In this way, the total regularization loss is:

$$\mathcal{L}_{\text{reg}} = \mathcal{L}_{\text{reg}}^p + \lambda \mathcal{L}_{\text{reg}}^b, \quad (7)$$

where  $\lambda$  controls the relative importance of the two loss terms.

1) *Unsupervised learning*: In the unsupervised case, the ground-truth user annotated scores are not available during training. The regularization loss and the reward terms can be used to regularize the properties of the selected summary and guide the learning of selecting key frames from the video sequence. The combined cost function for the video summarization network is formulated as

$$\mathcal{L}_{\text{uns}} = \mathcal{L}_{\text{reg}} - \mathcal{R}, \quad (8)$$

where the subscript in  $\mathcal{L}_{\text{uns}}$  indicates that the loss terms and reward terms are optimized in a fully unsupervised manner in the 3DST-UNet<sub>unsup</sub> variant of our method.

2) *Supervised learning*: In case user annotations of key frames are available, we can extend our 3DST-UNet-RL method to a supervised learning model with an objective that minimizes the distance between the predicted frames-wise importance scores and user annotated scores. We denote this supervised model as 3DST-UNet<sub>sup</sub>. Thus, the cost function for the supervised summarization model can be defined as

$$\mathcal{L}_{\text{sup}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{reg}} - \mathcal{R}, \quad (9)$$

where  $\mathcal{L}_{\text{pred}}$  is the mean squared errors (MSE) between the predicted frame-wise importance scores  $p_t$  and ground-truth user annotated scores  $p_t^*$ , *i.e.*,

$$\mathcal{L}_{\text{pred}} = \frac{1}{L} \sum_{t=1}^L \|p_t - p_t^*\|^2. \quad (10)$$

The objective of the RL network is to train a video summarization agent for the optimal policy  $\pi$  which indicates the actions to take to maximize the overall reward. The expected reward  $\mathcal{J}(\theta)$  is  $\mathcal{J}(\theta) = \mathbb{E}_{p_\theta(a|\pi)}[\mathcal{R}_{\text{rep}} + \mathcal{R}_{\text{div}}]$ , where  $p_\theta(a|\pi)$  denotes the probability distribution over the actions of sequences.

We follow [35] to compute the derivative of the objective function and approximate the gradient by taking the average of

$k$  repeated episodes for each video while subtracting a constant value  $c$ . This is used to avoid high variance which would make it difficult for the network to converge. The constant  $c$  is computed as the moving average of the previous rewards. The derivative of the expected reward  $\mathcal{J}(\theta)$  can be expressed as:

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{p_\theta(a_{1:T})} \left[ \sum_{t=1}^T (\mathcal{R} - c) \nabla_\theta \log \pi_\theta(a_t | h_t) \right]. \quad (11)$$

## F. Key Frames to Key Shots Conversion

Once we obtained frame-level importance scores via the deep RL network, we generate summary videos following existing protocols [10], [14]. To transfer a key frame-based summary into a keyshot-based summary, the testing video is first segmented into shot intervals using kernel temporal segmentation (KTS) [29]. The shot intervals containing at least one keyframe are marked as key shots. Finally, the knapsack algorithm [20] is applied to select key shots with the highest averaging keyframe scores while keeping the duration of summaries below a threshold. The default length of the summary videos are restricted to be 15% duration of the original video for fair comparison with other video summarization methods. We also perform experiments on summaries generated with different summary length constraints to analyse the impact of prescribed summary lengths.

## IV. EXPERIMENTS

### A. Datasets

In order to test the performance of our proposed 3DST-UNet-RL method on both general video sets and medical videos, we perform experiments on widely used video summarization benchmarks, *i.e.*, SumMe [2], TVSum [20], OVP and YouTube [21], [22]. We also demonstrate the effectiveness of the video summarization method on a fetal screening ultrasound video dataset [40].

1) *The TVSum Dataset*: The Title-based Video Summarization dataset [20] contains 50 videos of various genres (e.g., news, documentary, egocentric) and 1,000 annotations of shot-level importance scores (20 user annotations per video). The duration varies from 2-10 minutes.

2) *The SumMe Dataset*: The SumMe dataset [2] consists of 25 videos, each annotated with at least 15 human-annotated summaries. The duration of videos varying from 1.5–6.5 minutes. The data consists of videos, annotations and source code for standardized evaluation, which we also use in our experiments.

3) *OVP and YouTube*: We utilize additional videos as augmentation datasets to alleviate overfitting as in [10], [14]. The OVP [21], [22] and YouTube datasets [22] are constructed for keyframe-based video summarization. OVP contains 50 videos downloaded from the Open Video Project. The duration of each video is 1–4 min. YouTube contains 50 videos collected from the YouTube website of duration 1–10 minutes. Both of OVP and YouTube are provided with five key-frame based summaries annotated by human.

4) *The Ultrasound Dataset:* The ultrasound dataset are screen capture video recordings from fetal screening ultrasound examinations. There are 50 videos of 13-65 minutes length in our dataset from 50 different patients acquired between 24-30 weeks of gestation. The videos have been acquired and labelled during routine screenings according to the guidelines in the UK National Health Service (NHS) FASP handbook [48]. The feature extraction network is trained on annotations indicating the type of standard ultrasound diagnostic plane. From all available FASP planes we have selected Brain (Cerebellum), Brain (Ventricle), Profile, Lips, Abdominal, Kidneys, Femur, Spine (Coronal), Spine (Sagittal), 4-Chamber (4CH) cardiac view, 3 Vessel View (3VV) cardiac view, Right Ventricular Outflow Tract (RVOT), Left Ventricular Outflow Tract (LVOT) as the most frequent exemplars. For ultrasound video summarization, we take the freeze-frame images which are saved by the sonographers during the scan as the ground-truth key frames.

### B. Experimental Settings

We implement the 3DST-UNet-RL method in PyTorch based on the architecture of the DSN [35] network. We take 3D ResNet50 architecture [49] as the backbone network, which is a good compromise between performance and computational complexity. ResNet50 has shown competitive performance against other deeper architectures, for example for the task of action recognition.

The experiments were run on a single TITAN RTX GPU. Stochastic Gradient Descent with momentum 0.9 and weight decay of  $10^{-6}$  is used to train the 3DST-UNet-RL model. The initial learning rate was set to  $10^{-5}$  for SumMe and  $10^{-6}$  for TVSum, and was subsequently reduced by a factor of 0.5 for every 30 epochs. We set  $\lambda = 0.01$  and  $\epsilon = 0.5$  for Eq. (7). For ST3D and I3D, the squeeze layer of in 3DST-UNet compress the features  $F$  through  $1 \times 1$  convolutions from  $C = 2048$  to  $C' = 32$ .

The supervised learning model relies on the use of a single ground-truth summary to compute the prediction loss. For ultrasound dataset, the ground-truth key frame are freeze-frame images which are saved by a single sonographer. For SumMe and TVSum dataset, however, there are multiple human-annotated summaries. Therefore, we follow prior works [10], [14], [50] to generate an ‘‘oracle’’ summary [50] for SumMe and TVSum dataset that maximally agrees with all annotators for each video. During training, the ‘‘oracle’’ summary is served as the ground-truth frame importance scores.

### C. Experimental Results

We performed experiments on five different splits of training and testing subsets on SumMe and TVSum datasets. The percentages for training and testing subsets are 80% and 20% for both TVSum and SumMe datasets. The averaged F1 scores of the five splits are computed for both unsupervised and supervised learning paradigms.

TABLE I  
COMPARISONS OF F1 SCORES USING UNSUPERVISED APPROACHES ON THE SUMME AND TVSUM DATASETS. (THE BEST SCORES ARE IN BOLD.)

Method	TVSum	SumMe
TVSum [20]	50.0	36.0
SUM-GAN <sub>unsup</sub> [11]	51.7	39.1
Backprop-Grad [26]	52.7	-
SUM-FCN <sub>unsup</sub> [14]	52.7	41.5
UnpairedVSN <sub>unsup</sub> [24]	53.6	44.8
SASUM [51]	53.9	40.6
DR-DSN <sub>unsup</sub> [35]	57.6	41.4
DR-DSN <sub>sup</sub> (baseline)	56.7	43.2
3DST-UNet <sub>unsup</sub> (ours)	<b>58.1</b>	<b>44.6</b>

Following the protocols in [2], [20], [47], we compute the precision ( $P$ ) and recall ( $R$ ) as well as their harmonic mean  $F1$ -score against the user summary for evaluation, *i.e.*,

$$F = \frac{(2P \times R)}{(P + R)}, \quad (12)$$

where  $P$  and  $R$  are computed according to the temporal overlap, *i.e.*,  $|\mathcal{A} \cup \mathcal{B}|$  between a user annotated summary  $\mathcal{A}$  and a network predicted summary  $\mathcal{B}$ :

$$P = \frac{|\mathcal{A} \cap \mathcal{B}|}{|\mathcal{A}|}, \quad R = \frac{|\mathcal{A} \cap \mathcal{B}|}{|\mathcal{B}|}, \quad (13)$$

where  $|\mathcal{A}|$  and  $|\mathcal{B}|$  denotes the duration of summary  $\mathcal{A}$  and  $\mathcal{B}$ , respectively.

1) *Comparison to State-of-the-Arts:* We compare our proposed 3DST-UNet<sub>unsup</sub> with state-of-the-art unsupervised video summarization methods including TVSum [20], SUM-GAN [11], Backprop-Grad [26], SASUM [51], SUM-FCN [14] and DR-DSN [35]. The results are given in Table I.

The supervised training result of the proposed method 3DST-UNet<sub>sup</sub> are compared in Table II. SUM-GAN<sub>sup</sub> [11], DR-DSN<sub>sup</sub> [35], and SASUM<sub>sup</sub> [51] are extended from the unsupervised version by adding the discriminative loss of generated summaries and human-annotated summaries.

From the experimental results in Table I and II, our method outperforms the recurrent sequence modeling approach, *i.e.*, VSUMM [22], dppLSTM [10] and SUM-GAN [11]. To be noted, the DR-DSN<sub>sup</sub> (baseline) in Table I and II indicated the results that we obtained by using the same training/testing data splits as in our methods. Under the strictly fair comparison setting, our performance surpass the comparison method not only for supervised model but also for unsupervised model on both two testing datasets.

### D. Ablation Study

We first conduct ablation studies on the public video dataset [20] to investigate (i) the effectiveness of spatio-temporal 3D CNN features, (ii) the performance of 3DST-UNet for sequential modeling compared to bi-directional LSTM (Bi-LSTM), and (iii) the effectiveness of the reward terms.

For all the experiments in this section, we use a single split of training and testing videos on the TVSum dataset. We compared the proposed 3DST-UNet with the Bi-LSTM network applied in our baseline method [35] in terms of

TABLE II  
COMPARISONS OF F1 SCORES USING SUPERVISED APPROACHES ON THE SUMME AND TVSUM DATASETS. (THE BEST SCORES ARE IN BOLD.)

Method	TVSum	SumMe
LSTM [10]	54.2	37.6
dppLSTM [10]	54.7	38.6
SUM-GAN <sub>sup</sub> [11]	56.3	41.7
SASUM <sub>sup</sub> [51]	<b>58.2</b>	45.3
SUM-FCN <sub>sup</sub> [14]	56.8	<b>47.5</b>
UnpairedVSN <sub>sup</sub> [24]	55.6	<b>47.5</b>
DR-DSN <sub>sup</sub> [35]	<b>58.1</b>	42.1
DR-DSN <sub>sup</sub> (baseline)	57.2	45.7
3DST-UNet <sub>sup</sub> (ours)	<b>58.3</b>	<b>47.4</b>

the sequential model performance. The performance has been evaluated using 2D, Inflated 3D (I3D) and spatio-temporal 3D (ST3D) features and has been given in Table III.

1) *Bi-LSTM vs. 3DST-UNet*: The 3DST-UNet and Bi-LSTM in Table III indicates the summarization networks using our 3DST-UNet for sequential modeling and the baseline method [35] uses a bi-directional LSTM network topped with a fully connected layer, respectively. For both 3DST-UNet and the Bi-LSTM model, a Sigmoid function is used to predict frame probability scores.

As shown in Table III, the 3DST-UNet outperforms the Bi-LSTM model on sequential estimation of importance scores using both 2D and 3D video features. These comparison results conform to the conclusion that the proposed 3DST-UNet can encode spatio-temporal information of the input videos more efficiently for the video summarization task.

2) *2D CNN features vs Spatio-temporal 3D features*: The 2D CNN features are extracted from the *pool5* layer of a GoogleNet [52] pre-trained with ImageNet [44] and are of dimension 1024 which is identical to that in [10], [14], [35]. The backbone networks of 2D features are pre-trained with ImageNet, while the backbone networks of I3D are pre-trained on the Kinetics 400 database [53]. The ST3D features and I3D features takes  $n = 16$  frames at a time to compute representation features for the input video. In order to fit the I3D and ST3D features into the Bi-LSTM network for comparison, the spatial dimension of the features are pooled into 1 by 1, that is  $w = h = 1$  for  $\mathbb{R}^{T \times C \times w \times h}$ .

As compared in Table III, the performance of using the I3D and the ST3D features are better than that of using 2D features, both with the Bi-LSTMs model and the 3DST-UNet model. These results show evidence that spatio-temporal 3D features have stronger representation ability for video sequential modeling.

3) *The effectiveness of reward(s)*: We further conduct ablation experiments by dropping the rewards terms to analyze the effect of the rewards using for RL learning. According to Eq. (9) and Eq. (2),  $\mathcal{L}_{\text{sup}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{reg}} - (\mathcal{R}_{\text{rep}} + \mathcal{R}_{\text{div}})$ . By dropping both of the two reward terms, the model is trained using the regression loss  $\mathcal{L}_{\text{reg}}$  as well as the MSE loss  $\mathcal{L}_{\text{pred}}$  between the ground-truths and the predictions in a supervised manner. The results of ablation experiments on the effectiveness of reward terms are compared in Table IV. For the supervised learning settings, the F1 score increases by 0.6 by adding  $\mathcal{R}_{\text{rep}}$ .

TABLE III  
ABLATION STUDY COMPARING THE PROPOSED 3DST-UNET WITH THE BI-DIRECTIONAL LSTM NETWORK APPLIED IN OUR BASELINE METHOD [35] IN TERMS OF THE SEQUENTIAL MODEL PERFORMANCE. THE PERFORMANCE HAS BEEN EVALUATED USING 2D, INFLATED 3D (I3D) AND SPATIO-TEMPORAL 3D (ST3D) FEATURES. (THE BEST SCORES ARE IN BOLD.)

Method	Supervised			Unsupervised		
	2D	I3D	ST3D	2D	I3D	ST3D
Bi-LSTM	50.6	51.1	52.9	51.8	51.6	53.0
Ours	52.1	52.8	<b>54.1</b>	52.3	52.7	<b>53.6</b>

TABLE IV  
F1 SCORES OF THE VIDEO SUMMARIZATION RESULTS ON THE TVSUM DATASET USING AN UNSUPERVISED AND A SUPERVISED LEARNING PARADIGM WITH ( $\checkmark$ ) OR WITHOUT ( $\times$ )  $\mathcal{R}_{\text{rep}}$  AND  $\mathcal{R}_{\text{div}}$  REWARD TERMS. (THE BEST SCORES ARE IN BOLD.)

Rewards		Learning paradigm	
$\mathcal{R}_{\text{rep}}$	$\mathcal{R}_{\text{div}}$	Unsupervised	Supervised
$\times$	$\times$	-	56.3
$\checkmark$	$\times$	55.1	<b>56.9</b>
$\times$	$\checkmark$	54.5	56.1
$\checkmark$	$\checkmark$	<b>55.8</b>	56.6

This indicates that the usage of RL can further enable the model to predict higher quality summaries than using only the supervised loss. It is worth noting that the supervised model implicitly achieves diversity to some extent by directly learning from the human ground-truth annotations. Hence, it is reasonable that the performance of using  $\mathcal{R}_{\text{div}}$  is not better than the pure supervised model without using rewards.

### E. Augmented and Transfer Data Settings

To further analyze the results for our model from the main paper, we follow the prior works [10], [14], [24] to utilize augmentation videos from OVP [21], [22] and YouTube [22] as supplementary of the main SumMe and TVSum datasets.

1) *Canonical*: The canonical setting is the standard form where the training and testing sets are from the same dataset, as in our paper.

2) *Augmented*: In the augmented setting, for a given dataset, we randomly leave 20% of for testing, and augment the remaining 80% with the other three datasets to form an augmented training dataset.

3) *Transfer*: In the transfer setting, for a given dataset, we use the other three datasets for training and testing the learned models on the dataset.

Table V and Table VI show the performances of different methods in the Canonical (C), Augmented (A) and Transfer (T) settings on both TVSum and SumMe datasets.

As we can see from the results, under the *Augmented* setting, our proposed method has been further improved for both the supervised model and the unsupervised model. It is also noticed that, the supervised model does not improve as much as the unsupervised ones. This may due to the usage of ‘oracle’ as the ground-truth frame importance score for supervised training. Although the ‘oracle’ summary maximally agrees with all annotators for each video, however, as has been argued



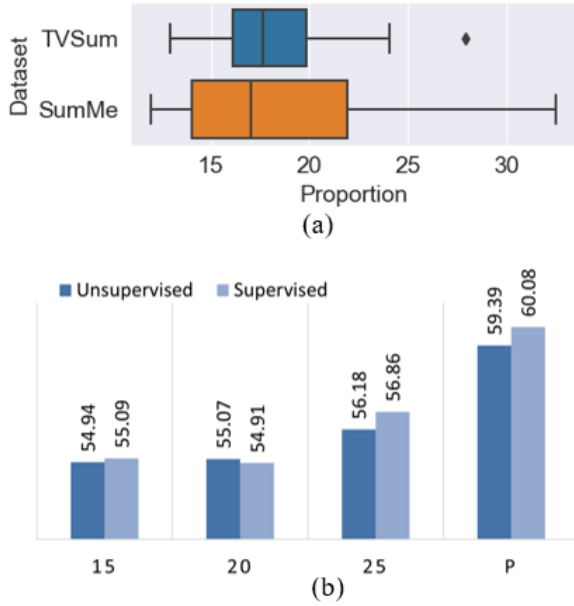


Fig. 5. (a) Box plot of the proportion of important frames given by human annotators in TVSum and SumMe datasets. (b) Comparison of F1 scores with four different summarization length constrains:  $l = 15\%$ ,  $l = 20\%$ , and  $l = 25\%$  and  $l = P$ , where  $P$  equals to proportion of important frames in the ground-truth for each video.

in [54] and [55], due to the highly-subjective of general video summarization task, it is not guaranteed that the inconsistent annotations from multiple users can fully explore the learning potential of a supervised model.

For the *Transfer* setting, the the surprisingly good performance, as has also been reported in [10], suggests a high correlation of the domains for the training datasets and the test dataset. These results are encouraged for the cases that little to no human annotations are available.

### F. Results with Different Summary Lengths

The above experimental results are conducted with a constant summary length constraint of 15%, *i.e.*, the length of the summary videos are restricted to be shorter than 15% of the input video length. Nevertheless, as we investigated in our experiments, the proportion of the important frames labeled by human annotators are not restricted with this threshold. We provide the distribution of summary length in the SumMe and TVSum datasets in Fig. 5(a). As evident, the proportion of the important frames labeled by human annotators can be as large as 30%.

We believe that the video summary performance will be impacted if we restrict the summary length to be 15% for input videos which have a ground-truth summary length of being longer than 15%. The video summarization network is forced to select only short video intervals of relatively high importance, instead of the intervals with highest scores. This will degrade the summary performance especially for the supervised models.

To mitigate this problem, we perform experiments on summaries generated with four different summary length constrains  $l$ : 15%, 20%, and 25% as well as  $P$ , where  $P$  is

TABLE V  
PERFORMANCES OF DIFFERENT METHODS IN THE CANONICAL (C), AUGMENTED (A) AND TRANSFER (T) SETTINGS ON THE TVSUM DATASETS.

Method	C	A	T
LSTM <sub>sup</sub> [10]	54.2	57.9	56.9
dppLSTM <sub>sup</sub> [10]	54.7	59.6	58.7
SUM-FCN <sub>unsup</sub> [14]	52.7	-	52.5
SUM-FCN <sub>sup</sub> [14]	56.8	59.2	58.2
UnpairedVSN <sub>unsup</sub> [24]	55.6	-	55.7
3DST-UNet <sub>unsup</sub> (ours)	58.3	58.4	58.0
3DST-UNet <sub>sup</sub> (ours)	58.3	58.9	56.1

TABLE VI  
PERFORMANCES OF DIFFERENT METHODS IN THE CANONICAL (C), AUGMENTED (A) AND TRANSFER (T) SETTING ON THE SUMME DATASETS.

Method	C	A	T
LSTM <sub>sup</sub> [10]	37.6	41.6	40.7
dppLSTM <sub>sup</sub> [10]	38.6	42.9	41.8
SUM-FCN <sub>unsup</sub> [14]	41.5	-	39.5
SUM-FCN <sub>sup</sub> [14]	47.5	51.5	44.1
UnpairedVSN <sub>sup</sub> [24]	47.5	-	41.6
3DST-UNet <sub>unsup</sub> (ours)	44.6	49.5	45.7
3DST-UNet <sub>sup</sub> (ours)	47.4	49.9	47.9

the averaged proportion of the important frames labeled by human annotators. The results are shown in Fig. 8 (b). When the video summarization network is allowed to select as many key frames as the human annotators, *i.e.*  $l = P$ , the F1 scores increase significantly. Compared to our unsupervised model 3DST-UNet<sub>unsup</sub>, the supervised model 3DST-UNet<sub>sup</sub> benefits even more from the relaxed summary length constrains.

### G. Application to Medical Video Summarization

For general videos, different individuals can have very different and subjective views regarding to the importance of video segments. The subjectiveness of human-annotators may lead to difficulties in the evaluation of general video summarization tasks. The common solution is to take the average [20] or maximum F1-scores [2] over the number of human created summaries. This would not become a problem for medical video summarization where the importance of video frames explicitly follow the clinical diagnostic standards.

We take a fetal ultrasound screening videos dataset to evaluate our method on the medical video summarization task. We show that our method is superior to alternative video summarization methods and that it preserves essential information required by clinical diagnostic standards. Some example summary frames for ultrasound examination recordings are given in Fig. 6.

We have compared our ultrasound video summarization results with the baseline method [35] with different summary lengths. Fig. 7 gives the F1 scores for three different summarization length constrains: again  $l = 15\%$ ,  $l = 25\%$ , and  $l = P$  where  $P$  is the proportion of important frames in the ground-truth for each video in the ultrasound dataset. As we can see, we outperform [35] for ultrasound video summarization. Our results are even on-par with that of [40] (*i.e.*, 63.29 with

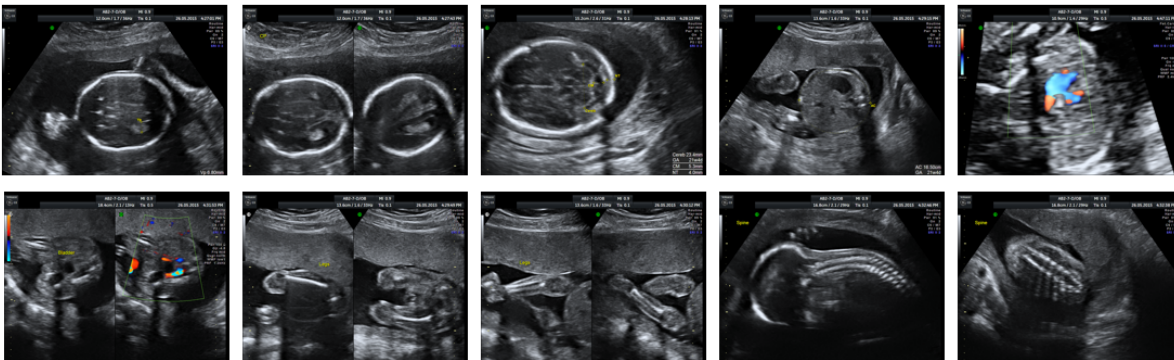


Fig. 6. Example key frames for an ultrasound examination recording videos, including the ultrasound plane of brain, abdominal, kidneys, femur, cardiac view and profile and spine. The key frames are acquired during routine ultrasound screenings according to the guidelines in the UK NHS FASP handbook.

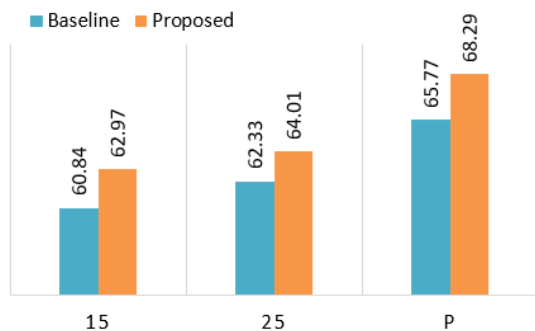


Fig. 7. Comparison of performance of the baseline method [35] and the proposed method on medical video summarization. The bar plot gives F1 scores with three different summarization length constrains:  $l = 15\%$ ,  $25\%$  and  $P$ , where  $P$  equals to proportion of the ground-truth important frames for each ultrasound video.

$l = 15\%$  as given in the paper) whose performance is achieved by using an additional supervised reward term that is informed by manually labeled diagnostic view planes.

In practical, the proposed video summarization method has the potential to save storage costs of ultrasound screening videos as well as to increase efficiency when browsing patient video data during retrospective analysis or audit without losing essential information.

#### H. Qualitative Results

We compare the qualitative video summarization results of the proposed 3DST-UNet method and our baseline method DSN [35]. In Fig. 8, the ground-truth frame-level importance scores of the video are indicated as light-green bars. The orange bars mark the intervals that have been selected by the summarization methods. Alongside, we show the selected key frames sampled which are the ones that have the highest prediction scores within a video shot.

We observe that our method preserves the temporal story of the videos by extracting intervals from different sections while focusing on key scenes. This implies that our method is able to preserve information essential for generating meaningful summaries.

## V. CONCLUSIONS

In this paper we have explored spatio-temporal CNN features for video frame representation in combination with 3DST-UNet, a spatio-temporal 3D U-Net, to model both the spatial and the temporal dependencies amongst video frames. We show with our experiments that 3D CNN features have stronger representation abilities than commonly used spatial 2D CNN features. The downstream RL framework learns from spatio-temporal latent spaces to predict actions for keeping/rejecting a video frame as being important for a summary. An RL agent takes the latent probability output from a spatio-temporal model to apply learned policies, which can maximize the accumulated reward of actions. We provide a comprehensive ablation study to investigate the contribution of the individual components and critically analyze the impact of hard-coded summary length constraints. The 3DST-UNet-RL model achieves competitive state-of-the-art performance on both general video summarization tasks and a medical video summarization task. The ultrasound video summarization method can be used for a variety of applications also when clinical annotations are unavailable. The proposed framework has the potential to save storage costs as well as to increase efficiency when browsing patient medical video data during retrospective analysis.

## ACKNOWLEDGMENT

We thank the volunteers and sonographers from routine fetal screening at St. Thomas' Hospital London. This work was supported by the Wellcome Trust IEH Award [102431] and EPSRC EP/S013687/1. The research was funded/supported by the National Institute for Health Research (NIHR) Biomedical Research Center based at Guy's and St Thomas' NHS Foundation Trust, King's College London and the NIHR Clinical Research Facility (CRF) at Guy's and St Thomas'. Data access only in line with the informed consent of the participants, subject to approval by the project ethics board and under a formal Data Sharing Agreement. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health.

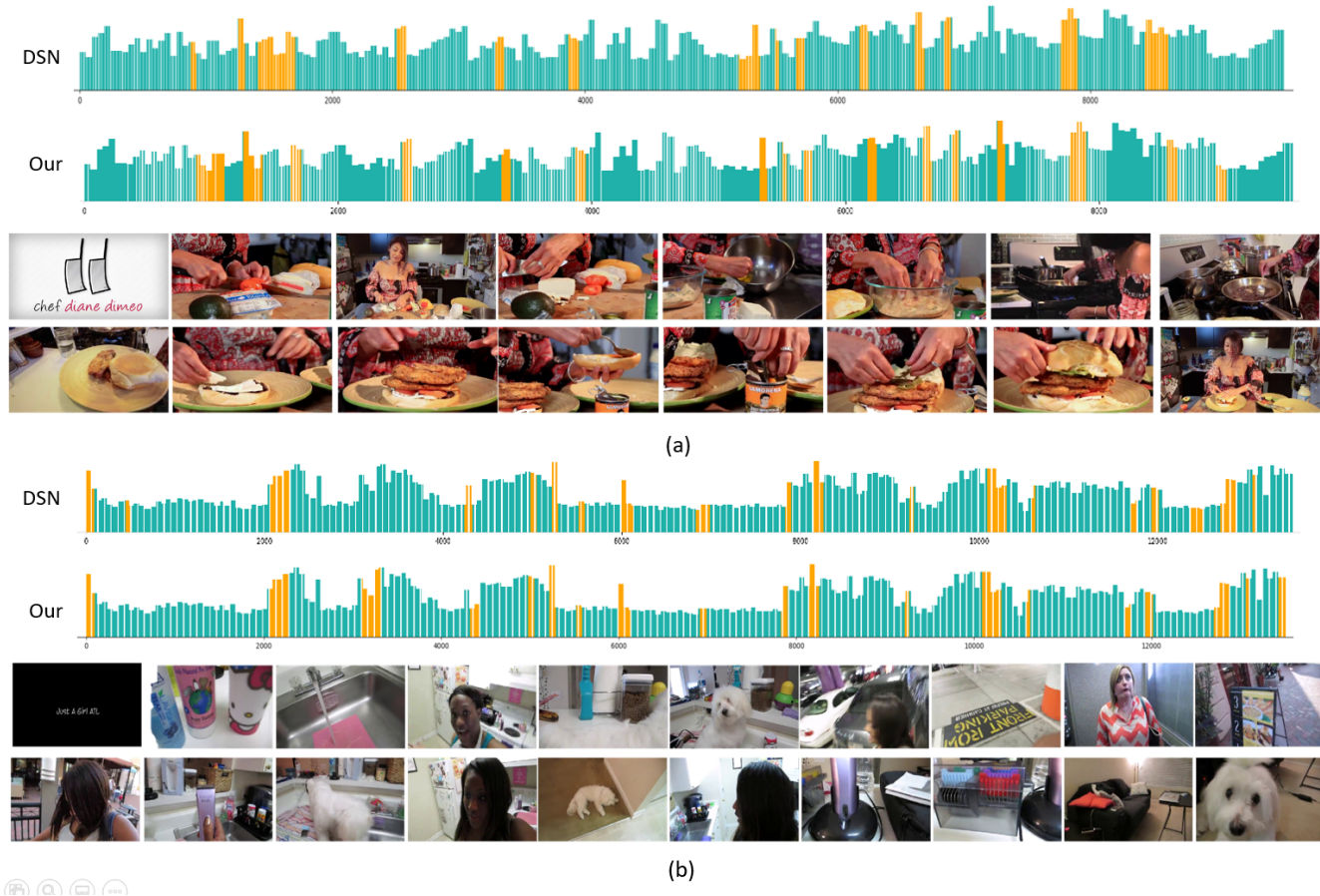


Fig. 8. Qualitative video summarization result of an exemplar video from the TVSum dataset. In the bar charts, the light-green bars correspond to ground-truth importance scores; the orange bars mark the intervals that have been selected by the proposed 3DST-UNet-RL and the baseline method DSN [35], respectively.

## REFERENCES

- [1] S. K. Kuanar, R. Panda, and A. S. Chowdhury, "Video key frame extraction through dynamic Delaunay clustering with a structural constraint," *Journal of Visual Communication and Image Representation*, vol. 24, no. 7, pp. 1212–1227, 2013.
- [2] M. Gygli, H. Grabner, H. Riemenschneider, and L. Van Gool, "Creating summaries from user videos," in *ECCV*, 2014.
- [3] Y. Cong, J. Yuan, and J. Luo, "Towards scalable summarization of consumer videos via sparse dictionary selection," *IEEE Transactions on Multimedia*, vol. 14, no. 1, pp. 66–75, 2011.
- [4] M. Yang, D. Dai, L. Shen, and L. Van Gool, "Latent dictionary learning for sparse representation based classification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [5] Y. Li, T. Zhang, and D. Treter, "An overview of video abstraction techniques," Technical Report HPL-2001-191, HP Laboratory, Tech. Rep., 2001.
- [6] T. Liu and S. Chan, "Automatic shot boundary detection algorithm using structure-aware histogram metric," in *2014 19th International Conference on Digital Signal Processing*, Aug 2014, pp. 541–546.
- [7] B. T. Truong and S. Venkatesh, "Video abstraction: A systematic review and classification," *ACM transactions on multimedia computing, communications, and applications (TOMM)*, vol. 3, no. 1, pp. 3–es, 2007.
- [8] D.-A. Huang, V. Ramanathan, D. Mahajan, L. Torresani, M. Paluri, L. Fei-Fei, and J. Carlos Niebles, "What makes a video a video: Analyzing temporal information in video understanding models and datasets," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7366–7375.
- [9] C. Li, Q. Zhong, D. Xie, and S. Pu, "Collaborative spatiotemporal feature learning for video action recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7872–7881.
- [10] K. Zhang, W.-L. Chao, F. Sha, and K. Grauman, "Video summarization with long short-term memory," in *European conference on computer vision*. Springer, 2016, pp. 766–782.
- [11] B. Mahasseni, M. Lam, and S. Todorovic, "Unsupervised video summarization with adversarial LSTM networks," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 202–211.
- [12] B. Zhao, X. Li, and X. Lu, "Hierarchical recurrent neural network for video summarization," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 863–871.
- [13] —, "HSA-RNN: Hierarchical structure-adaptive RNN for video summarization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7405–7414.
- [14] M. Rochan, L. Ye, and Y. Wang, "Video summarization using fully convolutional sequence networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 347–363.
- [15] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [16] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [17] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4489–4497.
- [18] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler, "Convolutional learning of spatio-temporal features," in *European conference on computer vision*. Springer, 2010, pp. 140–153.
- [19] J. Carreira and A. Zisserman, "Quo vadis, action recognition? A new

- model and the kinetics dataset,” in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [20] Yale Song, J. Vallmitjana, A. Stent, and A. Jaimes, “Tvsun: Summarizing web videos using titles,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015, pp. 5179–5187.
- [21] “Open video project,” <https://open-video.org/>.
- [22] S. E. F. de Avila, A. P. B. Lopes, A. da Luz, and A. de Albuquerque Araújo, “Vsumm: A mechanism designed to produce static video summaries and a novel evaluation method,” *Pattern Recognition Letters*, vol. 32, no. 1, pp. 56 – 68, 2011, image Processing, Computer Vision and Pattern Recognition in Latin America.
- [23] K. Zhang, K. Grauman, and F. Sha, “Retrospective encoders for video summarization,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 383–399.
- [24] M. Rochan and Y. Wang, “Video summarization by learning from unpaired data,” in *Proceedings of the AAAI Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7902–7911.
- [25] L. Yuan, F. E. Tay, P. Li, L. Zhou, and J. Feng, “Cycle-SUM: cycle-consistent adversarial LSTM networks for unsupervised video summarization,” in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 9143–9150.
- [26] R. Panda, A. Das, Z. Wu, J. Ernst, and A. K. Roy-Chowdhury, “Weakly supervised summarization of web videos,” in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [27] A. Khosla, R. Hamid, C.-J. Lin, and N. Sundaresan, “Large-scale video summarization using web-image priors,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 2698–2705.
- [28] S. Cai, W. Zuo, L. S. Davis, and L. Zhang, “Weakly-supervised video summarization using variational encoder-decoder and web prior,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 184–200.
- [29] D. Potapov, M. Douze, Z. Harchaoui, and C. Schmid, “Category-specific video summarization,” in *European conference on computer vision*. Springer, 2014, pp. 540–555.
- [30] K. Zhou, T. Xiang, and A. Cavallaro, “Video summarisation by classification with deep reinforcement learning,” *arXiv preprint arXiv:1807.03089*, 2018.
- [31] J. Supancic III and D. Ramanan, “Tracking as online decision-making: Learning a policy from streaming videos with reinforcement learning,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 322–331.
- [32] Y. Rao, J. Lu, and J. Zhou, “Attention-aware deep reinforcement learning for video face recognition,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 3931–3940.
- [33] Y. Chen, S. Wang, W. Zhang, and Q. Huang, “Less is more: Picking informative frames for video captioning,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 358–373.
- [34] N. Zeng, H. Li, Z. Wang, W. Liu, S. Liu, F. E. Alsaadi, and X. Liu, “Deep-reinforcement-learning-based images segmentation for quantitative analysis of gold immunochromatographic strip,” *Neurocomputing*, vol. 425, pp. 173–180, 2021.
- [35] K. Zhou, Y. Qiao, and T. Xiang, “Deep reinforcement learning for unsupervised video summarization with diversity-representativeness reward,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [36] Y. Li, L. Wang, T. Yang, and B. Gong, “How local is the local diversity? reinforcing sequential determinantal point processes with dynamic ground sets for supervised video summarization,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 151–167.
- [37] S. Lan, R. Panda, Q. Zhu, and A. K. Roy-Chowdhury, “FFNet: Video fast-forwarding via reinforcement learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6771–6780.
- [38] A. Vlontzos, A. Alansary, K. Kamnitsas, D. Rueckert, and B. Kainz, “Multiple landmark detection using multi-agent reinforcement learning,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 262–270.
- [39] A. Alansary, O. Oktay, Y. Li, L. Le Folgoc, B. Hou, G. Vaillant, K. Kamnitsas, A. Vlontzos, B. Glocker, B. Kainz *et al.*, “Evaluating reinforcement learning agents for anatomical landmark detection,” *Medical image analysis*, vol. 53, pp. 156–164, 2019.
- [40] T. Liu, Q. Meng, A. Vlontzos, J. Tan, D. Rueckert, and B. Kainz, “Ultrasound video summarization using deep reinforcement learning,” in *Proceedings of International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2020.
- [41] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, “3D U-Net: Learning dense volumetric segmentation from sparse annotation,” in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 424–432.
- [42] S. Qamar, H. Jin, R. Zheng, P. Ahmad, and M. Usama, “A variant form of 3d-unet for infant brain segmentation,” *Future Generation Computer Systems*, vol. 108, pp. 613–623, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167739X18332291>
- [43] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [44] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [45] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [46] T. Liu, W. Luo, L. Ma, J.-J. Huang, T. Stathaki, and T. Dai, “Coupled network for robust pedestrian detection with gated multi-layer feature extraction and deformable occlusion handling,” *IEEE Transactions on Image Processing*, vol. 30, pp. 754–766, 2021.
- [47] M. Gygli, H. Grabner, and L. Van Gool, “Video summarization by learning submodular mixtures of objectives,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3090–3098.
- [48] ., “Fetal anomaly screening programme:,” *handbook for ultrasound practitioners April 2015*, 2015.
- [49] K. Hara, H. Kataoka, and Y. Satoh, “Can spatiotemporal 3D CNNs retrace the history of 2D CNNs and ImageNet?” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 6546–6555.
- [50] B. Gong, W.-L. Chao, K. Grauman, and F. Sha, “Diverse sequential subset selection for supervised video summarization,” in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 2069–2077.
- [51] H. Wei, B. Ni, Y. Yan, H. Yu, X. Yang, and C. Yao, “Video summarization via semantic attended networks,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [52] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [53] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, “The kinetics human action video dataset,” *arXiv preprint arXiv:1705.06950*, 2017.
- [54] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris, and I. Pateras, “Unsupervised video summarization via attention-driven adversarial learning,” in *International Conference on Multimedia Modeling*. Springer, 2020, pp. 492–504.
- [55] M. Otani, Y. Nakashima, E. Rahtu, and J. Heikkilä, “Rethinking the evaluation of video summaries,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7596–7604.