# Video Summarization

**Dr. I. Mademlis, M. Kaseris, P. Alexoudi, C. Aslanidou,
Prof. Ioannis Pitas
Aristotle University of Thessaloniki
pitas@csd.auth.gr
www.aiia.csd.auth.gr
Version 1.5**

Artificial Intelligence &
Information Analysis Lab

# Contents

- Introduction
- Video summarization use-cases
- Video summary types
- Video summarization approaches
- Content selection algorithms
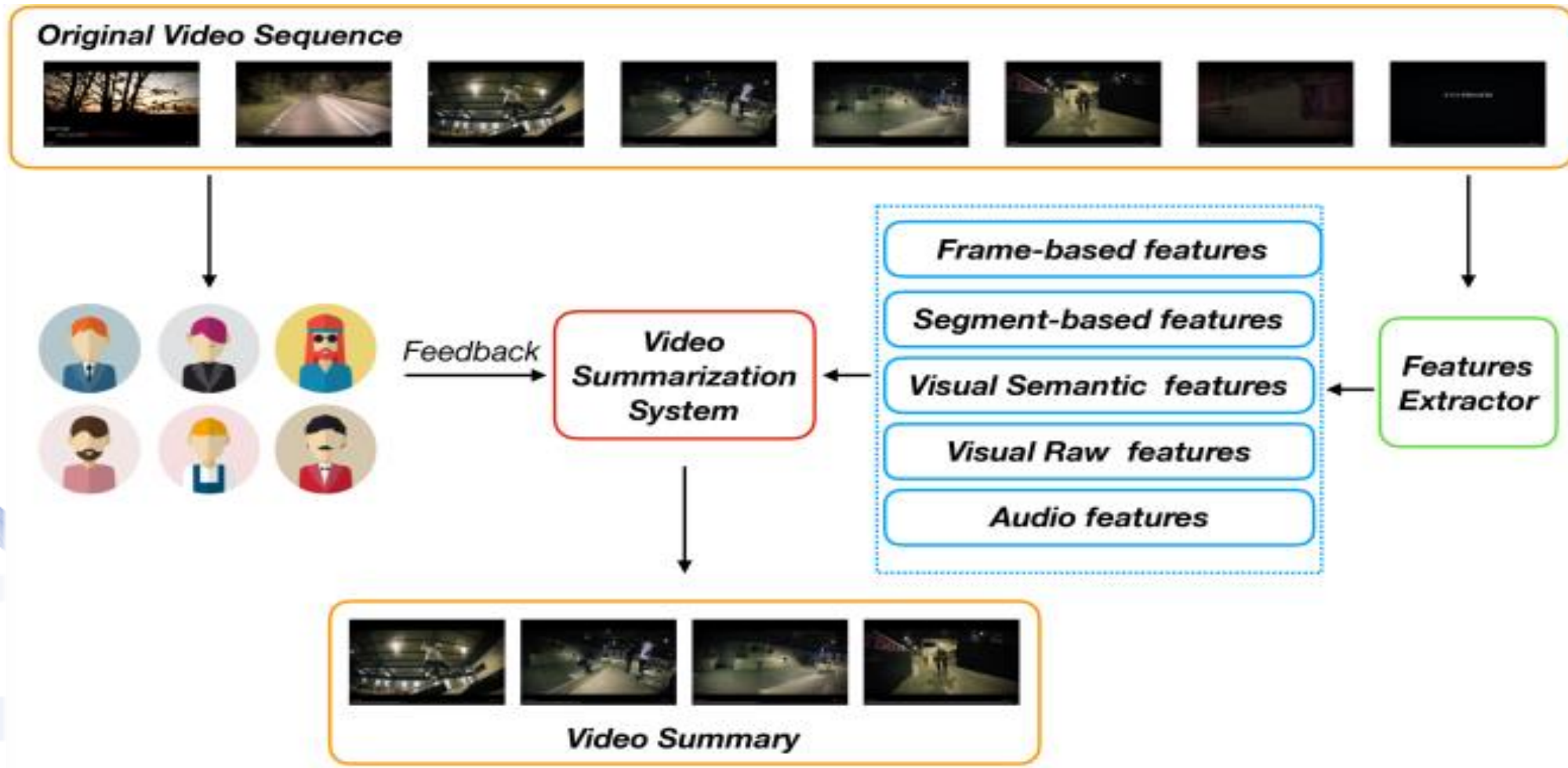- Video Summarization with Deep Neural Networks

# Contents

- GANs for video summarization
- SUM-GAN-AAE
- DTR-GAN
- Cycle-SUM
- Summary diversity
- DNNs and dictionary learning
- DNNs and deep reinforcement learning
- Evaluation datasets
- Bibliography

# Introduction

- ***Video summarization*** is the automated construction of a short version of an original full-length video.

- It is necessary in applications where videos are recorded, stored and accessed in abundance.

- Video summarization has various applications in several industries (media, surveillance, WWW, etc.).

- ***Example***: Users would ideally like to browse quickly through large video databases, to get an idea of the content.

# Introduction



(Image from Heartbeat Fritz AI)

# Introduction

- Video summarization algorithms result in a short **summary** of the video.

- The challenge is to automatically select which content will be retained and which will be discarded during the summarization process.

- Only the **most informative** and/or interesting parts should be kept.

Artificial Intelligence &
Information Analysis Lab

# Video summarization use-cases

- ***Movie trailers***

- Advertisement creation

- Sports highlights

- ***Anomaly detection from video surveillance***

- Redundancy removal

- Reduction of computational time, storage requirements

- Data visualization

- Search, Retrieval, Recommendation [WOR2020]

Artificial Intelligence &
Information Analysis Lab

# Video summarization use-cases

- **Summarization of personal videos** [DAR2014]
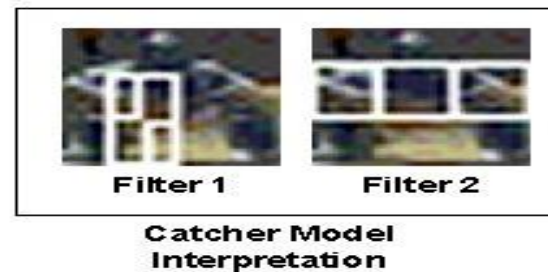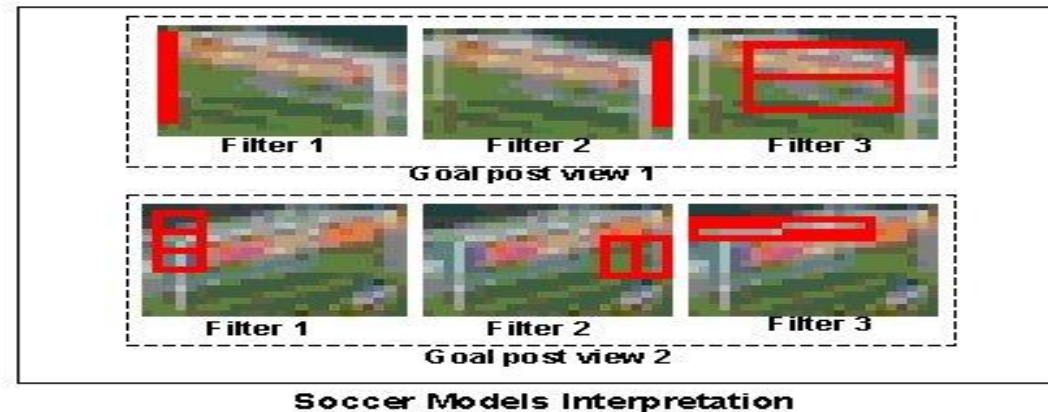
# Video summarization use-cases
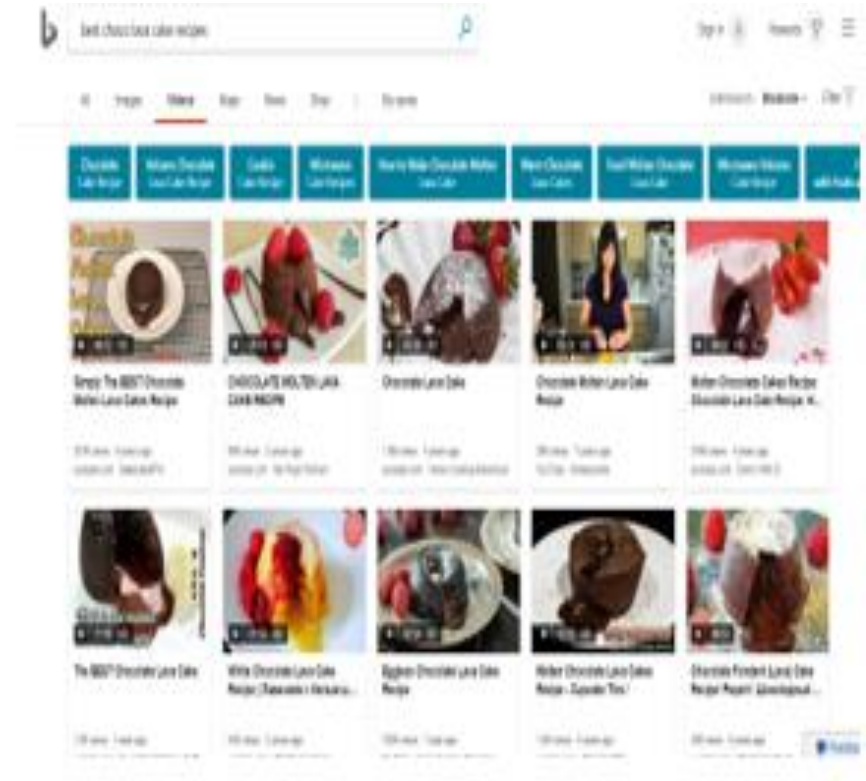
- **Sports highlights** [ZHA2006]

# Video summarization use-cases

- **Automatic TV/film trailers** [BOR2018]

# Video summarization use-cases

- **Video search engines** [IRI2010]

# Video summarization use-cases

- **Egocentric Video storyboard**

Input: Egocentric video of the camera wearer's day

Output: Storyboard summary of important people and objects

Image from vision.cs.utexas.edu

# Video summarization use-cases

- **Medical Video summarization**



Image from E3S Web of Conferences

# Video summary types

- There are two main types of video summaries: [MAD2016]

  - **Static video summaries** (storyboard/gallery/key-frame set),

  - **Dynamic video summaries** (skims/trailers).

- A static summary is a temporally ordered set of selected ***key-frames***.

  - A collection of still images.

- A dynamic summary is a temporally ordered set of selected ***key-segments***.

  - A trailer.

# Video summary types



Static Summary

Original Video

Dynamic Summary

(Image from Semantic Scholar)

# Video summarization approaches

**VML**

- Several video summarization methods have been developed over the years.

- They can be classified into *four major categories*, based on their properties and characteristics. [BUR2020]

Artificial Intelligence &
Information Analysis Lab

# Video summarization approaches

- ***Feature-based summarization*** [BUR2020]

  - The original video content is represented by an aggregation of various features.

  - These features may capture properties such as visible objects, events, color, motion type, etc.

  - Feature extraction and aggregation is the most important step.

  - A machine learning method (e.g., clustering) processes these features, in order to select only a subset of the original content.

# Video summarization approaches

- The selection process may optionally be applied at different levels of detail.

- **First**, the original video is segmented into scenes and/or shots.

- **Then**, important *key-scenes* and/or *key-shots* are identified and retained, while the remaining ones are discarded.

- **Finally**, important *key-frames* and/or *key-segments* are identified within each of the selected scenes/shots.

Artificial Intelligence & Information Analysis Lab

# Video summarization approaches

- Multiple alternative algorithms exist both for temporal video segmentation and for content selection [KAI2012].

- All content selection algorithms for video summarization attempt to identify key-frames/shots/scenes, so that the final summary is:

  - *Representative* of the content of the full-length original video,

  - *Concise* in length (e.g., the number of key-frames may be 10% of the number of original video frames), and

  - *Complete*, in the sense that it covers the entire content of the original video (e.g., no sequence of a movie is completely left out of the summary).

Artificial Intelligence &
Information Analysis Lab

# Video summarization approaches



Original Frame Sequence (many frames are omitted for the simplicity of the illustration)

Scene Identification with Global Features

Keyframe Selection with Local Features

Video Summary

Video Summarization with Global and Local Features (Image from ResearchGate)

# Video summarization approaches



Video input → ① → Video frames pre-sampling → ② → Color feature extraction (color histogram, HSV, 16 bins) → ③ → Frames clustering (k-means, Euclidean distance) → ④ → Keyframe extraction, Elimination of similar keyframes → ⑤ → Static video summary composition (temporal order)

Image from Heartbeat-Fritz AI

# Video summarization approaches

- ***Event-based summarization*** [BUR2020]

  - Visually ***abnormal/rare events*** are considered interesting (e.g., a robbery or traffic accident scene in a film).

  - The nature of such events depends on the employed video frame representations:

    - Low-level features expressing perceived motion, colors, etc.

    - Higher-level semantic features expressing visible objects, activities, etc.

  - The selection algorithm retains in the summary only parts of the original video that seem to contain abnormal content.

# Video summarization approaches

- **Event-based video summarization**



Image from ResearchGate

23

# Video summarization approaches

- ***Object-based summarization*** [BUR2020].

  - There are cases where we are only interested in the parts of the video depicting a specific family of objects (e.g., people).

  - An object detector is required to analyze each scene.

  - Only parts of the original video (frames or segments) containing the desired object(s) are retained in the summary.

Artificial Intelligence &
Information Analysis Lab

# Video summarization approaches

- **Object-based video summarization**



Image from ScienceDirect

# Video summarization approaches

- ***Attention-based summarization*** [BUR2020]

  - There are various ways to identify which parts of an original video hold most of the users' interest when they view it.

  - The derived summary may only contain key-frames/shots that have been assigned a high attention score.

  - For example, motion attention models may be employed to measure each shot's interest.

# Content selection algorithms

- Various content selection algorithms have been employed for video summarization.

- Video frame/shot/scene **clustering** (e.g., K-Means) is the simplest approach.

- More sophisticated methods (e.g., **spectral clustering**) have also been employed.

- Dictionary learning approaches are a good alternative to clustering.

**Artificial Intelligence & Information Analysis Lab**

# Content selection algorithms
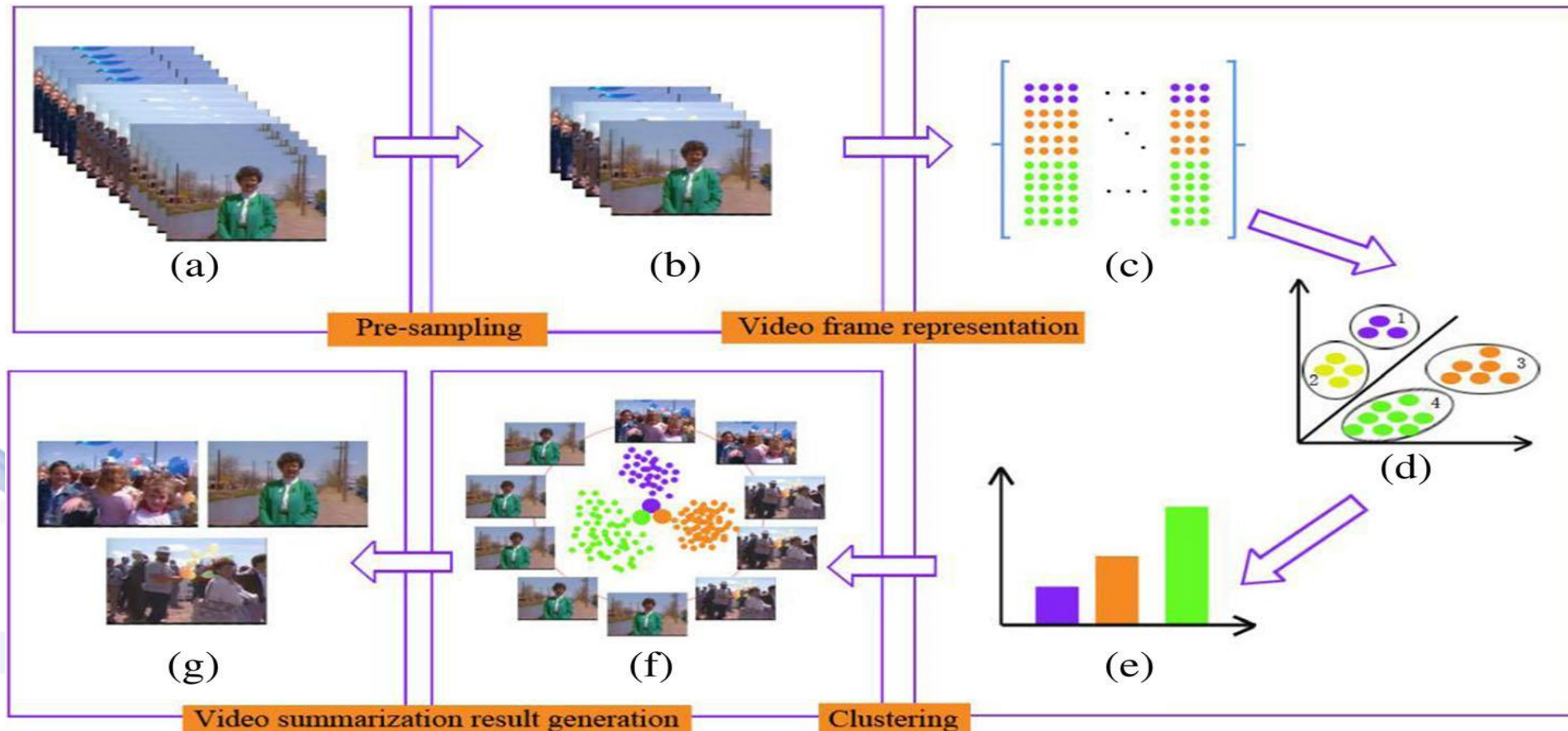
- All video frames are partitioned into clusters of similar properties and the medoid of each cluster is retained as a key-frame.

- Temporal subsampling may be applied before clustering, due to typically high similarities in the appearance of neighboring video frames.

- The exact same process may be applied at a shot or scene level.

Artificial Intelligence &
Information Analysis Lab

# Content selection algorithms

- **Clustering-based Video summarization.**



(a) Pre-sampling (b) Video frame representation (c) (d)

(g) Video summarization result generation (f) Clustering (e)

Image from Research Institute for Future Media Computing

# Content selection algorithms
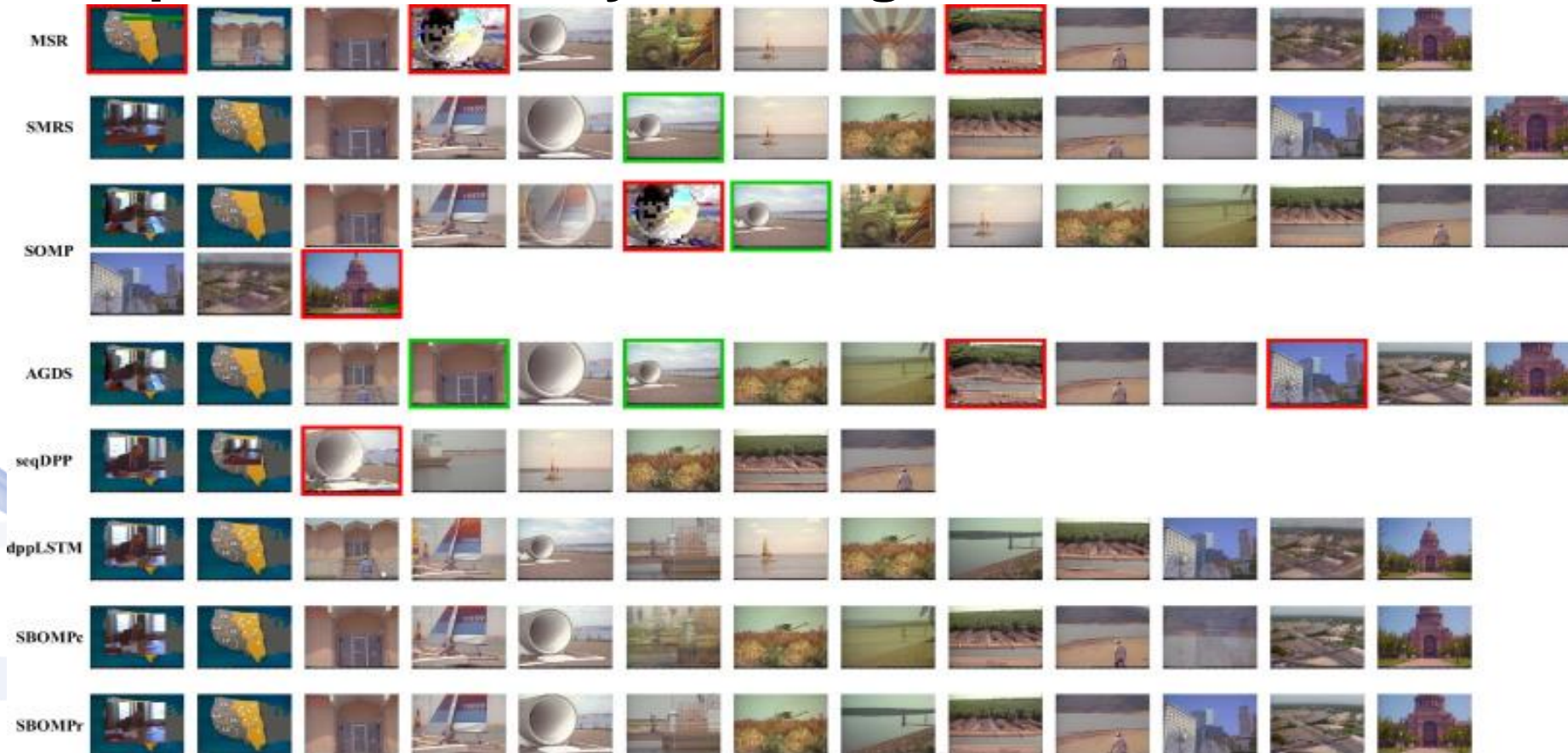
- ***Dictionary learning*** is an effective replacement for clustering algorithms.

- The extracted key-frames form a ***dictionary***.

- They should enable ***optimal reconstruction*** of the original video from the selected dictionary.

- Thus, the video summary is framed as the set of key-frames that can linearly reconstruct the full-length video in an algebraic sense [MAD2018].

# Content selection algorithms

- **Sparse dictionary learning.**



Image from sciencedirect.com

# Content selection algorithms

- Both clustering and dictionary learning are ***unsupervised learning*** approaches: no ground-truth summaries are required.

- The following approaches have also been proposed:
  - Reinforcement learning [WOR2020] or
  - supervised learning methods [DIN2019] .

- ***Supervised video summarization*** requires training of machine learning model using a manually annotated training dataset.

- The annotation may be an importance score assigned per video frame.

Artificial Intelligence & Information Analysis Lab

# Content selection algorithms



Image from IBM Developer

# Content selection algorithms

- The standard supervised approach has several disadvantages.

- ***Manual video annotation is quite expensive***, difficult and costly, especially if done at a per-frame level.

- Importance scores are quite subjective.

- The trained model may only perform well in test videos resembling the training dataset.

# Video Summarization with Deep Neural Networks

- In recent years, ***Deep Neural Networks*** (DNNs) have been employed for video summarization in various ways.

- The simplest approach is to exploit semantic video frame representations derived from pre-trained Convolutional Neural Networks (CNNs), as inputs to a traditional content selection algorithm.

**VML**

**Artificial Intelligence & Information Analysis Lab**

# Video Summarization with Deep Neural Networks

- A more sophisticated approach is to train a DNN under a supervised learning framework **to directly regress an importance score** for each original video frame.

- During the test stage, any video frame which is assigned a score larger than a threshold can be selected as a key-frame.

- This approach has all the disadvantages of supervised summarization.

Artificial Intelligence & Information Analysis Lab

# Video Summarization with Deep Neural Networks

- Various deep neural architectures may be combined in a composite DNN for video summarization. For example:

  - ***Convolutional Neural Networks*** (CNNs)

  - ***Transformers***

  - ***3D CNNs***

  - ***Long Short-Term Memory Networks*** (LSTMs)

  - ***Generative Adversarial Networks*** (GANs).

# GANs for unsupervised video summarization

- ***GANs combined with LSTMs*** have recently been employed for unsupervised video summarization, using an end-to-end trainable DNN architecture.

- GANs are generative models which learn the distribution of the training data. They are composed of a Generator and a Discriminator involved in a minimax game.

  - The Generator learns to generate content that the Discriminator mistakes for real.

- After training, the Discriminator may be discarded.

**Artificial Intelligence & Information Analysis Lab**

# GANs for unsupervised video summarization



Image from LaptrinhX

# GANs for unsupervised video summarization

Examples of fake faces

Artificial Intelligence & Information Analysis Lab
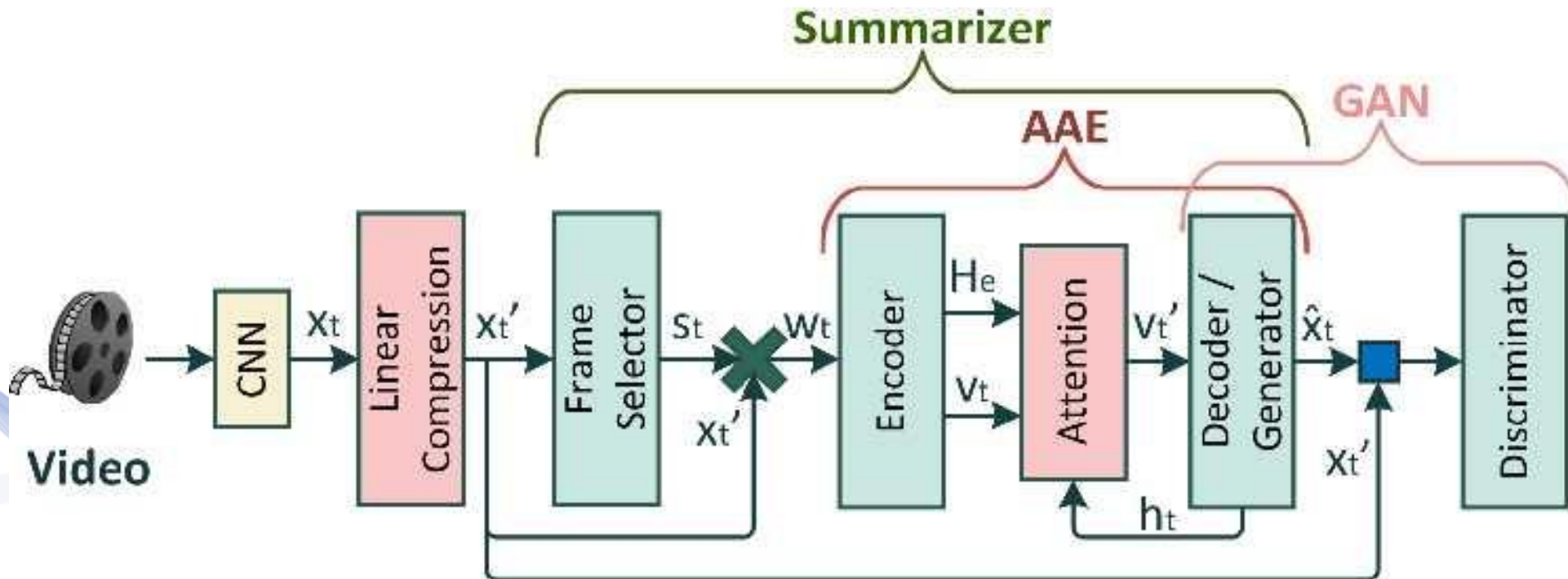
# GANs for unsupervised video summarization

- ***SUM-GAN-AA**E* [METS2020].

- ***Dilated Temporal Relational Adversarial Network*** for frame-level video summarization [DIN2019].

- ***Cycle-SUM***: Cycle-consistent Adversarial LSTM Networks for Unsupervised Video Summarization (Video Trailer) [PIN2019].

# SUM-GAN-AAE



The architecture of SUM GAN-AAE (Image from [METS2020])

# SUM-GAN-AAE

- SUM-GAN-AAE is a modification of SUM-GAN [MAH2017].

- The network architecture consists in a Summarizer subnetwork, which acts as a Generator, and a Discriminator subnetwork.

- The Summarizer is a pipeline of three smaller subnetworks:

  - Frame Selector, Encoder, Decoder.

- All subnetworks are LSTMs.

- After training, only the Frame Selector is required.

Artificial Intelligence &
Information Analysis Lab

# SUM-GAN-AAE

- The Frame Selector receives sequentially as input the original video frame representations.

- For each input video frame, it estimates and outputs an importance score.

- The original video frame representations and the importance scores are multiplied.

Artificial Intelligence &
Information Analysis Lab

# SUM-GAN-AAE

- The Encoder is sequentially fed the above products and produces a fixed-length representation for the entire video.

- The representation produced by the Encoder is fed to the Decoder, which is equipped with an attention mechanism.

- The Decoder is trained to sequentially output the original video frames.

- The Encoder-Decoder and the attention module jointly constitute the Attention Autoencoder subnetwork (AAE).

# SUM-GAN-AAE

- Both the original and the reconstructed video frame representations are then sequentially passed to the Discriminator, whose task is to determine whether each sequence is "real" (original) or "fake" (summary-based reconstruction).

- The Frame Selector and the AAE jointly constitute the Summarizer, which is trained to confuse the Discriminator.

  - This forces the Frame Selector to learn how to extract representative key-frames, jointly capable of accurately reconstructing the full-length video.

# SUM-GAN-AAE

- $\mathbf{X} \in \mathbb{R}^{M \times N}$ : The input video data matrix.

- Each column $\mathbf{x}_i \in \mathbb{R}^M$ of the matrix $\mathbf{X}$, is the feature representation of the $i$-th frame.

- The baseline summarization architecture includes:

  - An LSTM-based **Frame Selector** $S$ parameterized by weights $\mathbf{w}_s$.

  - An LSTM-based **Encoder** $E$ parameterized by weights $\mathbf{w}_e$.

  - An LSTM-based **Decoder** $D$ parameterized by weights $\mathbf{w}_d$.

  - An LSTM-based **Discriminator** (binary classifier) $C$ parameterized by weights $\mathbf{w}_c$.

Artificial Intelligence &
Information Analysis Lab

# SUM-GAN-AAE

- $S$ is fed $\mathbf{x}_i$ as input and outputs a corresponding scalar importance factor $s_i \in [0,1]$ .

- The product $s_i \mathbf{x}_i$ is fed to $E$ resulting in a state vector $\mathbf{e} \in \mathbb{R}^H$ encoding the summary.

- Subsequently, $\mathbf{e}$ is fed to $D$ which attempts to reconstruct the original $\mathbf{X}$, by outputting a reconstructed $\hat{\mathbf{x}}_i \in \mathbb{R}^M$ , $1 \le i \le N$ .

- Finally, the video reconstruction $\hat{\mathbf{X}}$ is forwarded to the Discriminator $C$ as a "fake" training example, while the original video $\mathbf{X}$ is used as a "real" training example.

# SUM-GAN-AAE

- The following loss functions are employed during training:

    - **_Reconstruction loss_**:

$$\mathcal{L}_{recon} = \left\| \phi(\mathbf{X}) - \phi(\widehat{\mathbf{X}}) \right\|_2^2,$$

    - $\phi(\mathbf{X})$ is the last hidden LSTM state when it is fed $\mathbf{X}$ as input

    - $\phi(\widehat{\mathbf{X}})$ the corresponding hidden LSTM state when $C$ is fed $\widehat{\mathbf{X}}$.

    - $\mathcal{L}_{recon}$ is used to update $\mathbf{w}_s, \mathbf{w}_e, \mathbf{w}_d$.

**Artificial Intelligence & Information Analysis Lab**

# SUM-GAN-AAE

- **Original video loss**:

$$\mathcal{L}_{orig} = \left(1 - C(\mathbf{X})\right)^2.$$

- It is the MSE between the original video label (i.e., 1) and the discriminator output (in [0,1]) when $C$ is fed $\mathbf{X}$ as input.

- $\mathcal{L}_{orig}$ updates $\mathbf{w}_c$.

- *Summary loss*:

$$\mathcal{L}_{sum} = \left(C(\widehat{\mathbf{X}})\right)^2$$

- is the MSE between the summary label (i.e., 0) and the computed probability when $C$ is fed $\widehat{\mathbf{X}}$ as input.

- $\mathcal{L}_{sum}$ updates $\mathbf{w}_c$.

**Artificial Intelligence & Information Analysis Lab**

# SUM-GAN-AAE

- **_Generator loss_**:

$$\mathcal{L}_{gen} = \left(1 - C(\widehat{\mathbf{X}})\right)^2.$$

- It is the MSE between the original video label (i.e., 1) and the discriminator output, when $C$ is fed $\widehat{\mathbf{X}}$ as input. $\mathcal{L}_{gen}$ updates the Decoder parameters $\mathbf{w}_d$.
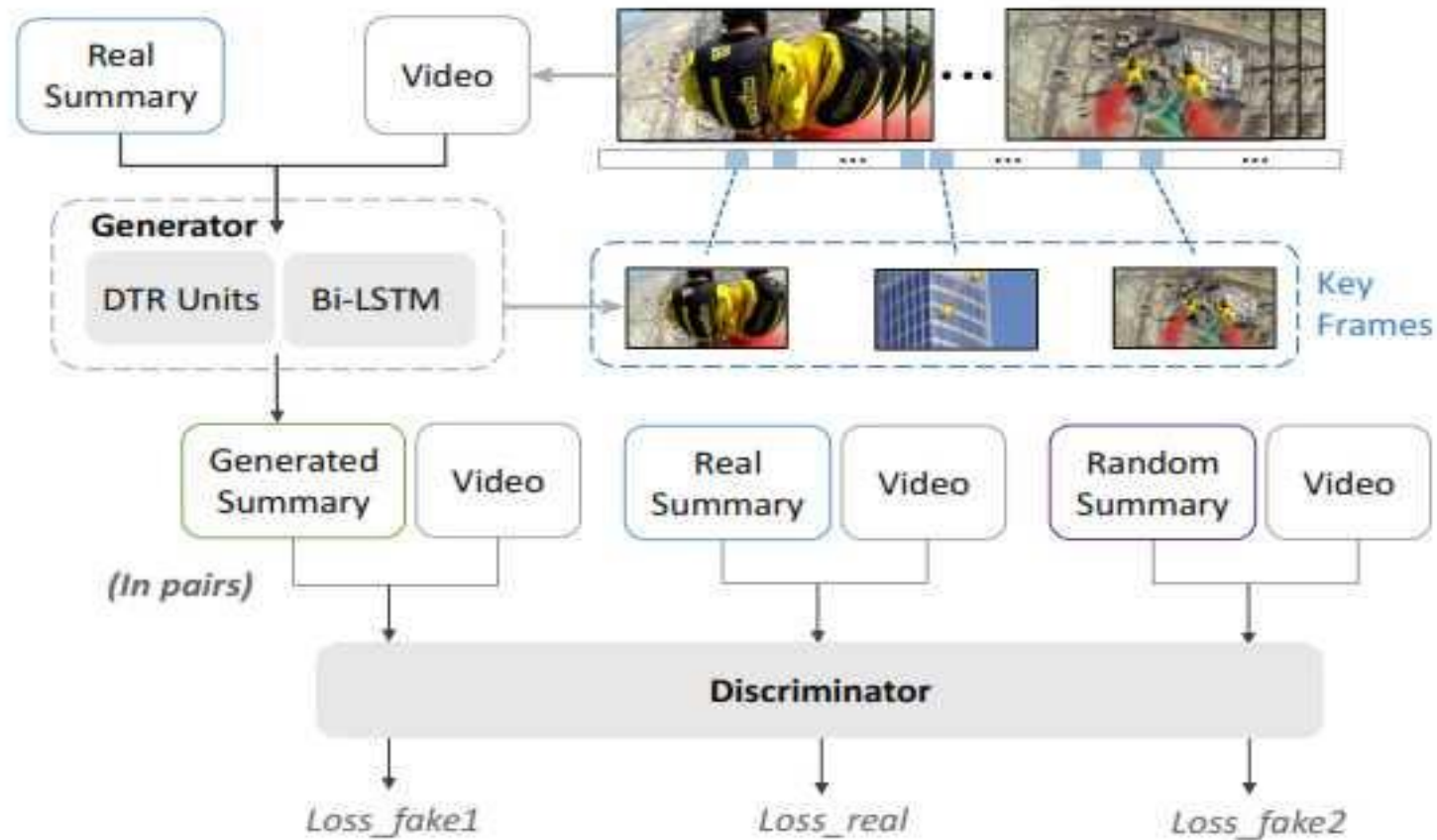
- **_Sparsity Loss_**:

$$\mathcal{L}_{sparsity} = \left\|\frac{1}{N}\sum_{t=1}^{N} s_t - \sigma\right\|_2.$$

- It pushes the Selector towards assigning high importance (i.e., key-frame status probability) to a specific (**_small_**) percentage of the total number of original video frames, defined by a scalar hyperparameter $\sigma \in [0,1]$.

- Typically $\sigma \in [0.1, 0.2]$.

- The sparsity loss updates $\mathbf{w}_s$.

Artificial Intelligence & Information Analysis Lab

# DTR-GAN



DTR-GAN (Image from [DIN2019])

# DTR-GAN

- The **Dilated Temporal Relational Generative Adversarial Network** (DTR-GAN) is an architecture slightly similar to SUM-GAN, but it is **supervised**.

- The Discriminator in DTR-GAN is trained with a composite three-part loss function, that takes jointly into account the generated summary, the ground-truth summary and a random summary.

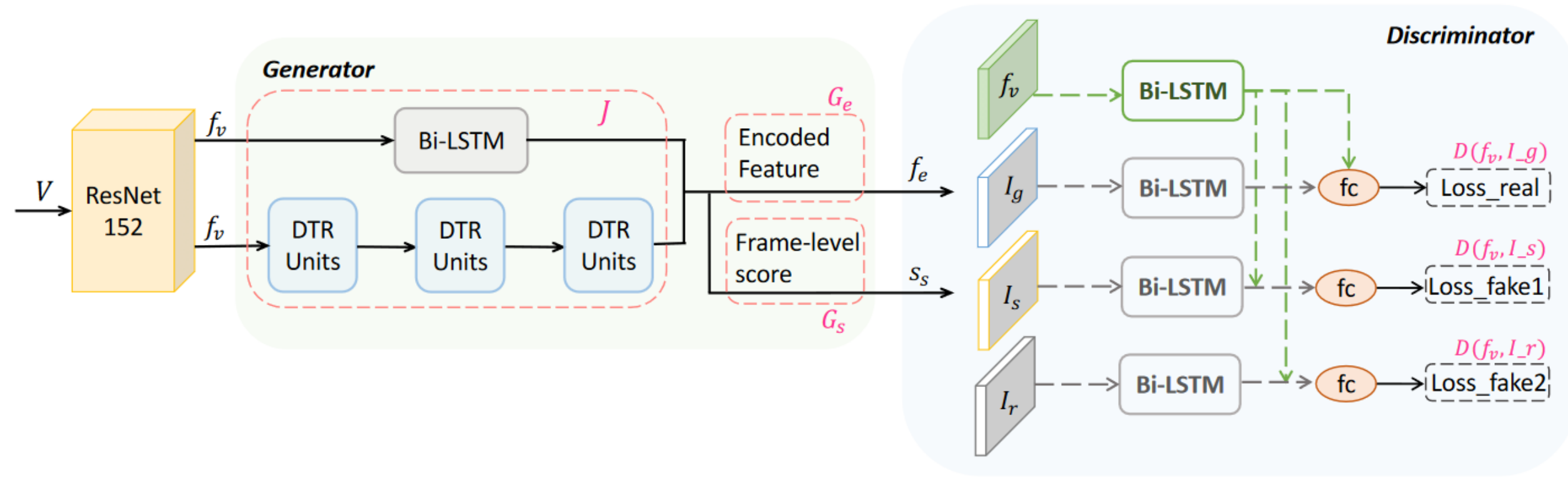- This provides better regularization.

# DTR-GAN

- The Frame Selector is enhanced in DTR-GAN: besides the LSTMs, it also contains **Dilated Temporal Relational** (DTR) units.

- DTR units aim to exploit **long-range temporal dependencies**, complementing LSTMs.

- They integrate context among video frames at multi-scale time spans, in order to enlarge the model's temporal field-of-view and, thus, effectively model temporal inter-frame relations.

# DTR-GAN

- There is no LSTM auto-encoder in the DTR-GAN Summarizer, because the Discriminator is given **video + summary pairs** as inputs.

- Thus, the Discriminator learns to evaluate the correspondence between an input video and its summary,

  - rather than whether its input video has been reconstructed from a generated summary or not, as is the case in SUM-GAN-AAE.

# DTR-GAN



DTR-GAN (Image from [DIN2019])

# Cycle-SUM

- ***Cycle-SUM*** is an unsupervised end-to-end trainable DNN for key-frame extraction, which extends the original SUM-GAN.

- During training, it replaces the unidirectional reconstruction of SUM-GAN/SUM-GAN-AAE (the original video is reconstructed from the generated summary) with a ***"circular" bidirectional video reconstruction***.

- A ***cyclic consistency loss term*** is added to the training objectives of the overall framework.

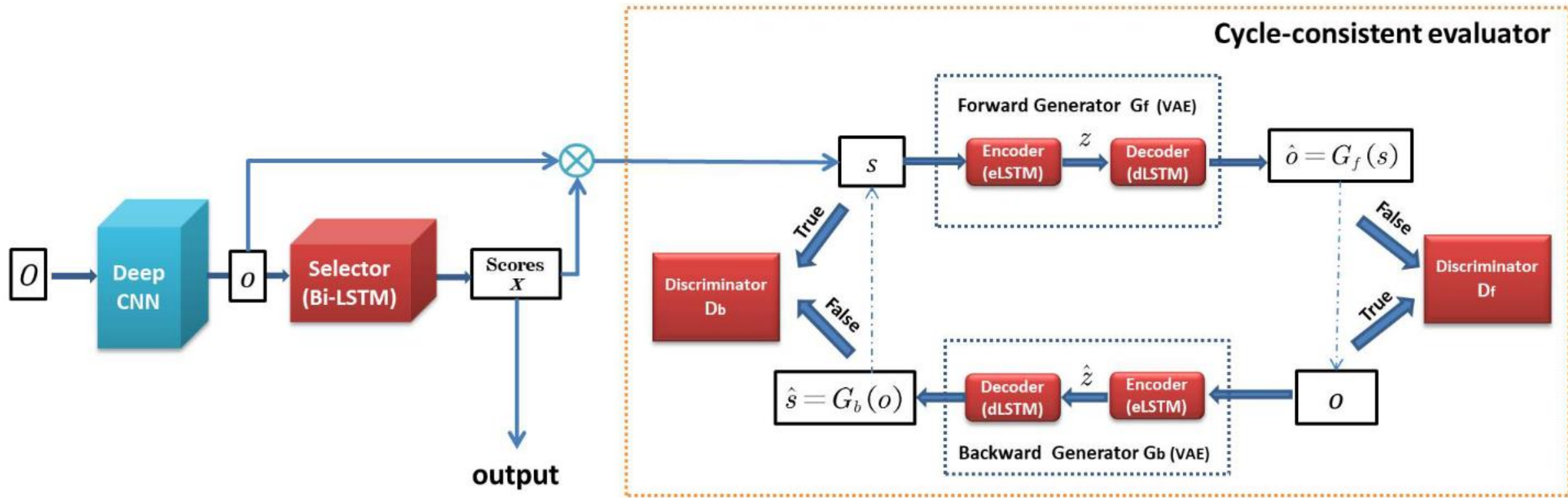**Artificial Intelligence & Information Analysis Lab**

# Cycle-SUM

- Cycle-SUM is composed of an initial Frame Selector, two autoencoders (instead of one) and two Discriminators (instead of one).

- The **forward autoencoder** and Discriminator reconstruct the original video from the generated summary and evaluate it, respectively.

- The **backward autoencoder** and Discriminator reconstruct the summary from the original video and evaluate it, respectively.

# Cycle-SUM

- The closed training loop which enforces the cyclic consistency aids the DNN to **maximize mutual information between the summary and the original**, full-length video.

- Explicitly enforcing the reconstruction cycle original → summary → original → summary, better guarantees summary completeness and representativeness.

# Cycle-SUM



Cycle-SUM architecture. (Image from [PIN2019])

# Summary diversity

- Most DNN-based methods for video summarization emphasize representativeness, conciseness and completeness of the summary.

- However, it may be equally important that the selected key-frames are ***diverse in visual content***.

- Summary variety makes it summary more interesting and reduces redundancy.

# Summary diversity

- A straightforward way to achieve summary diversity with DNNs is to add the so-called ***Determinantal Point Process*** (DPP) loss term in the pool of training objectives.

- In frameworks similar to SUM-GAN, the DPP loss directs the training process so that the Frame Selector learns to assign importance scores so that ***the overall summary is diverse***.

- This diversity pertains to the semantic content captured in the input video frame representations (e.g., visible objects).

# Summary diversity

- The DPP loss operates by:

  - Quantifying the variance of the set of video frame representations.

  - Penalizing candidate key-frame sets/summaries that do not capture significant percentage of the original video variance.

- Consider a matrix $\mathbf{L} \in \mathbb{R}^{T \times T}$ by computing the pairwise cosine similarity for time step $t$ and $t'$ that is, $L_{ij} = \mathbf{e}_t^T \mathbf{e}_{t'}$.

- $\mathbf{e}_t$ and $\mathbf{e}_{t'}$ are the Encoder's hidden states at time step $t$ and $t'$, respectively.

# Summary diversity

- ***DPP loss***:

$$\mathcal{L}_{dpp} = -\log\left(\frac{\det(\mathbf{L}_y)}{\det(\mathbf{L}+\mathbf{I})}\right).$$

- $\mathbf{L}_y$ is a submatrix whose rows and columns are dictated by the indices of the selected keyframes and $\mathbf{I}$ is the identity matrix.

- Recently, the DPP loss was extended so that it also captures diversity of additional modalities, besides the CNN-derived representations expressing visible objects in each video frame.

- By enforcing diversity in the textual descriptions of each video frame, scene context and visible activities are also considered [KAS2022].

# Summary diversity

- SUM-GAN-AAE is employed as a baseline and a pre-trained image captioner $P$ is required.

- Then, the *DPP-caption loss* exhorts the video summary to be more diverse in terms of textual semantic content.

- During training, each video frame is forwarded to $P$, in parallel to feeding it to the Encoder.

- The following cost is used for Frame Selector weight update:

$$\mathcal{L}_{dpp-c} = -\log \frac{\det(\mathbf{P}_y)}{\det(\mathbf{P+I})}.$$

# DNNs and dictionary learning

- Integrating dictionary learning into unsupervised deep neural frameworks such as SUM-GAN-AAE, has also been attempted [KAS2021].

- Using SUM-GAN-AAE as a baseline, an additional pre-trained autoencoder is employed to pre-encode the entire video sequence into a single fixed-length vector $\mathbf{h}$.

- During training, a novel loss term is added to the framework:

$$\mathcal{L}_{dict} = \|\mathbf{h} - \mathbf{A}\mathbf{e}\|_2.$$

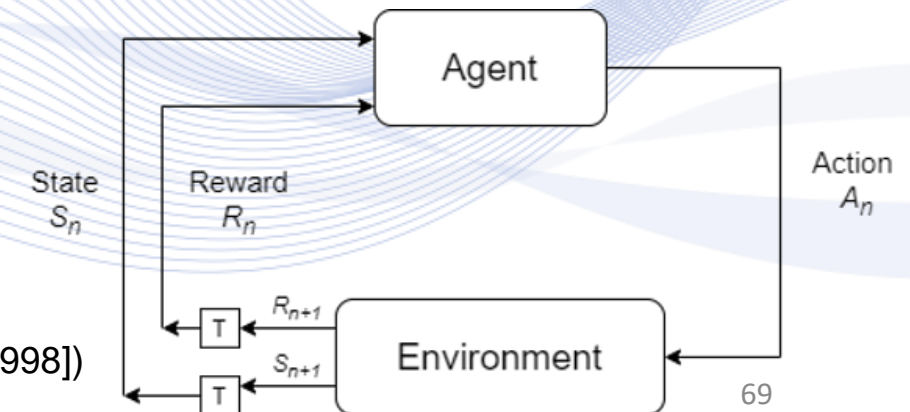- Vector $\mathbf{e}$ is given by the Encoder, while $\mathbf{A}$ is learnt.

**Artificial Intelligence & Information Analysis Lab**

# DNNs and dictionary learning

- Matrix **A** transforms the current summary representation to a vector space being simultaneously learnt from all the original videos.

- **A** essentially serves as a ***global visual dictionary***.

- Thus, each summary representation is exhorted towards being a set of linear reconstruction coefficients that are jointly able to reproduce the corresponding original video representation.

- This is on top of the non-linear reconstruction objective enforced by the baseline SUM-GAN-AAE.
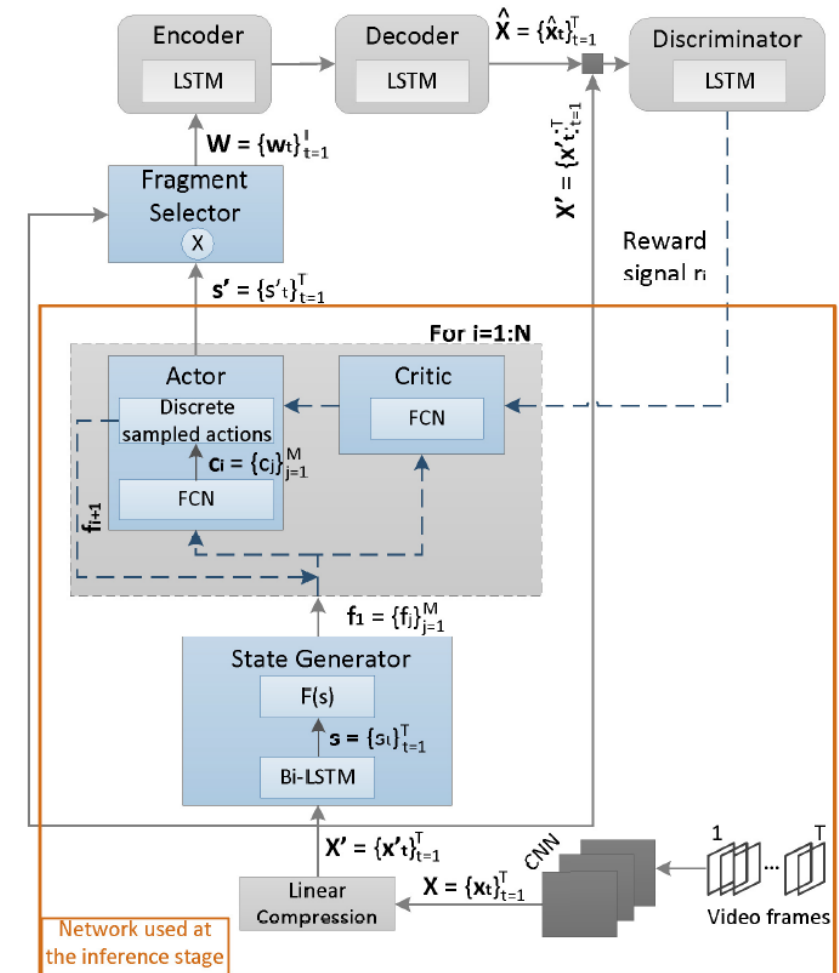
# DNNs and reinforcement learning

- Reinforcement learning (RL) has also been integrated into unsupervised deep neural frameworks for video summarization.

- In RL, a cognitive agent is trained through interaction: it interacts with its environment, in order to find a policy that maximizes a cumulative reward.

- The reward is a numerical measure that determines how good the agent's action was.

- The learned policy maps states to actions.

Environment-action interaction (Image from [SUT1998])

State $S_n$   Reward $R_n$   Agent   Action $A_n$

$R_{n+1}$

$S_{n+1}$

Environment

Artificial Intelligence & Information Analysis Lab

# DNNs and reinforcement learning

- AC-SUM-GAN is a good example of combining SUM-GAN with RL. [APO2020]

- A neural Actor-Critic architecture is embedded into SUM-GAN.

- During training, it learns the optimal policy for key-frame extraction.

- During inference, the RL agent modifies/adjusts the video frame importance scores outputted by the Frame Selector.



The architecture of AC-SUM GAN (Image from [APO2020])

# DNNs and reinforcement learning

- ***The Actor generates sequences incrementally***, based on a set of discrete sampled actions over a group of video fragments.

- ***The Critic evaluates the Actor's choices*** and returns a value for scoring each choice, according to its impact on the action-state space.

- The Discriminator acts as the RL environment and returns a reward that is used to train the Actor-Critic model, which learns a value function (Critic) and a policy for key-fragment selection (Actor).

- The Critic can be discarded after training.

Artificial Intelligence & Information Analysis Lab

# DNNs and reinforcement learning

- The Actor plays an "N-picks" game to explore the action-state space.

- For every step $i$, $(1 \leq i \leq N)$:

  - It receives the current state $\mathbf{f}_i = \{f_j\}_{j=1}^{M}$, where $M$ is the number of non-overlapping fragments into which the video is segmented.

    - At time $i = 1$, $\mathbf{f}_1$ is derived from the vector of importance scores outputted by the Frame Selector.

**Artificial Intelligence & Information Analysis Lab**

# DNNs and reinforcement learning

- (continued)

  - It produces a distribution of actions $\mathbf{c}_i = \{c_j\}_{j=1}^{M}$.

  - It takes an action by sampling the computed distribution $\mathbf{c}_i$, thus, picking a video fragment $k$ for inclusion in the summary.

  - This action modifies the state and produces $\mathbf{f}_{i+1}$.

  - During training, the reward is the Discriminator's classification decision.

# Evaluation Datasets

- There are several public datasets for evaluating video summarization algorithms.

- Typically, these datasets provide a collection of videos with associated per-frame ground truth importance scores.

- The most common ones are TVSum and SumMe.

  - *SumMe* includes 25 videos of 1 to 6 minutes duration with diverse video contents, captured both from first and third-person view.

  - *TVSum* consists of 50 videos of 1 to 11 minutes duration, containing video content from 10 categories of the TRECVid MED dataset.

Artificial Intelligence &
Information Analysis Lab

# Evaluation Datasets

- Every video of the dataset **is annotated by multiple users** in the form of key fragments (SumMe) or frame-level importance scores (TVSum)

  - Single ground-truth summaries are also provided.

- To evaluate a video summarization algorithm, the generated summary for a given video is compared with the users' summary, separately per user.
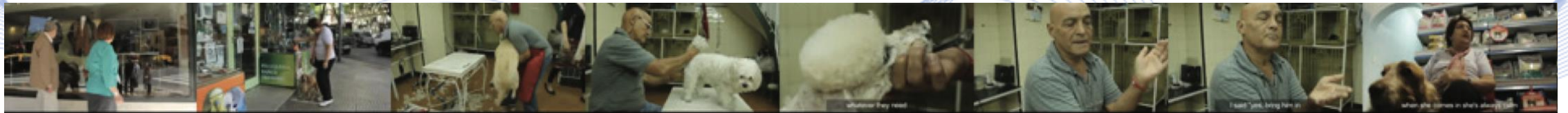
# Evaluation Datasets

- An F-Score (F-measure) is computed for each pair of compared summaries.

- The computed F-Scores for TVSum are averaged or the maximum of them is kept for SumMe and a final F-Score is obtained for this video.

- The computed F-Scores for the entire set of testing videos are finally averaged to quantify the algorithm's performance.

# Evaluation Datasets



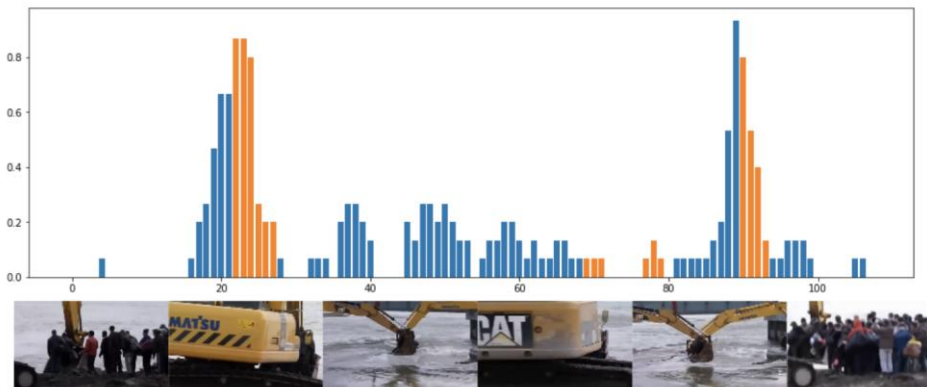Video frames from the sequence "Cooking" of the SumMe dataset.



Video frames from the sequence "Dog grooming in Buenos Aires" of the TVSum dataset.

**Artificial Intelligence & Information Analysis Lab**

# Evaluation Datasets



Video frames from the sequence "Excavators road crossing"
of the SumMe dataset.



The video frame importance scores and the extracted
summary using SUM-GAN-AAE in combination with $\mathcal{L}_{dict}$ +
$\mathcal{L}_{dpp}$

# Bibliography

[PIT2017] I. Pitas, "Digital video processing and analysis" , China Machine Press, 2017 (in Chinese).

[PIT2013] I. Pitas, "Digital Video and Television" , Createspace/Amazon, 2013.

[PIT2021] I. Pitas, "Computer vision", Createspace/Amazon, in press.

[NIK2000] N. Nikolaidis and I. Pitas, "3D Image Processing Algorithms", J. Wiley, 2000.

[PIT2000] I. Pitas, "Digital Image Processing Algorithms and Applications", J. Wiley, 2000.

Artificial Intelligence & Information Analysis Lab

# Bibliography

[DAR2014] K. Darabi and G. Ghinea, "Personalized video summarization by highest quality frame", IEEE International Conference on Multimedia and Expo Workshops (ICMEW), 2014.

[ZHA2006] Z. Zhao, S. Jiang, Q. Huang and G. Zhu, "Highlight summarization in sports video based on replay detection", IEEE International Conference on Multimedia and Expo, 2006.

[BOR2018] A. Bora and S. Sharma, "A review on video summarization approaches: Recent advances and directions", International Conference on Advances in Computing, Communication Control and Networking (ICACCCN), 2018.

[IRI2010] G. Irie, T. Satou, A. Kojima, T. Yamasaki and K. Aizawa, "Automatic trailer generation", ACM International Conference on Multimedia, 2010.

[KAS2022] M. Kaseris, I. Mademlis, and I. Pitas, "Exploiting caption diversity for unsupervised video summarization", Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2022.

[MAD2018] I. Mademlis, A. Tefas, and I. Pitas, "A salient dictionary learning framework for activity video summarization via key-frame extraction", Elsevier Information Sciences, 432, 319-331, 2018.

Artificial Intelligence & Information Analysis Lab

# Bibliography

[BUR2020] H. B.U. Haq, M. Asif and M. B. Ahmad, "Video summarization techniques: A review", International Journal of Scientific Technology Research", volume 9, issue 11, 2020.

[KAI2012] G. Guan, Z. Wang, K. Yu, S. Mei, M. He and D. Feng, "Video summarization with global and local features", IEEE International Conference on Multimedia and Expo Workshops, 2012.

[SAB2012] W. Sabbar, A. Chergui, A. Bekkhoucha, "Video summarization using shot segmentation and local motion estimation", Innovative Computing Technology, pp. 190–193, 2012.

[MAD2016] I. Mademlis, A. Tefas, N. Nikolaidis and Ioannis Pitas, "Multimodal stereoscopic movie summarization conforming to narrative characteristics", IEEE Transactions on Image Processing, 25.12: 5828-5840, 2016.

[SUT1998] R. S. Sutton and A. G. Barto, "An Introduction to Reinforcement Learning", MIT Press, 1998.

Artificial Intelligence &
Information Analysis Lab

# Bibliography

[WOR2020] A. Workie and R. Sharma, "Digital video summarization techniques: A survey", International Journal of Engineering Research Technology, vol. 9, issue 01, pp. 81-85, 2020.

[SUP2017] J. Supancic III, D. Ramanan. "Tracking as online decision-making: Learning a policy from streaming videos with reinforcement learning", Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017.

[TRU2007] B. T. Truong, Venkatesh s., "Video Abstraction: A Systematic Review and Classification", ACM Transactions on Multimedia Computing, Communications, and Applications, 3:3, 2007.

[XIA2021] H. Xiao and J. Shi, "Diverse video captioning through latent variable expansion", arXiv preprint arXiv:1910, 2021.

[KAS2021] M. Kaseris, I. Mademlis and I. Pitas, "Adversarial Unsupervised Video Summarization Augmented With Dictionary Loss", Proceedings of the IEEE International Conference on Image Processing (ICIP), 2021.

Artificial Intelligence &
Information Analysis Lab

# Bibliography

[SHE2015] C. V. Sheena, N. K. Narayanan, "Key-frame extraction by analysis of histograms of video frames using statistical methods", International Conference on Eco-friendly Computing and Communication Systems, 2015.

[BAL2019] D. Sen and B. Raman, "Video skimming: Taxonomy and comprehensive survey", arXiv preprint arXiv:1909.12948, 2019.

[METS2020] I. A. Metsai, V. Mezaris, E. Apostolidis, E. Adamantidou and I. Patras, "Unsupervised video summarization via attention-driven adversarial learning", International Conference on Multimedia (MMM), 2020.

[DIN2019] D. Zhang, M. Tan, E. P. Xing ,Y. Zhang, X. Liang, "Dilated temporal relational adversarial network for generic video summarization", Springer Multimedia Tools and Applications, 2019, 78.24, pp. 35237-35261.

[MAH2017] B. Mahasseni, M. Lam and S. Todorovic, "Unsupervised video summarization with adversarial LSTM networks", Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2017.

Artificial Intelligence & Information Analysis Lab

# Bibliography

[PIN2019] P. Li, L. Zhou, J. Feng, L. Yuan, F. EH Tay, "Cycle-sum: Cycle-consistent adversarial lstm networks for unsupervised video summarization", Proceedings of the AAAI Conference on Artificial Intelligence, 2019.

[APO2020] E. Apostolidis, E. Adamantidou, A. I. Metsai, V. Mezaris, I. Patras, "AC-SUM-GAN: Connecting Actor-Critic and Generative Adversarial Networks for Unsupervised Video Summarization", IEEE Transactions on Circuits and Systems for Video Technology, 2020, 31.8: 3278-3292.

**Artificial Intelligence & Information Analysis Lab**

# Q & A

**Thank you very much for your attention!**

**More material in**
**http://icarus.csd.auth.gr/cvml-web-lecture-series/**

**Contact: Prof. I. Pitas**
**pitas@csd.auth.gr**