# Learnable Optimal Sequential Grouping for Video Scene Detection

Daniel Rotman, Yevgeny Yaroker, Elad Amrani, Udi Barzelay, Rami Ben-Ari

[danieln@il.,yevgenyy@il.,elad.amrani@,udib@il.,ramib@il.]ibm.com

IBM Research

Haifa, Israel

## ABSTRACT

Video scene detection is the task of dividing videos into temporal semantic chapters. This is an important preliminary step before attempting to analyze heterogeneous video content. Recently, Optimal Sequential Grouping (OSG) was proposed as a powerful unsupervised solution to solve a formulation of the video scene detection problem. In this work, we extend the capabilities of OSG to the learning regime. By giving the capability to both learn from examples and leverage a robust optimization formulation, we can boost performance and enhance the versatility of the technology. We present a comprehensive analysis of incorporating OSG into deep learning neural networks under various configurations. These configurations include learning an embedding in a straight-forward manner, a tailored loss designed to guide the solution of OSG, and an integrated model where the learning is performed through the OSG pipeline. With thorough evaluation and analysis, we assess the benefits and behavior of the various configurations, and show that our learnable OSG approach exhibits desirable behavior and enhanced performance compared to the state of the art.

## KEYWORDS

Video Scene Detection, Deep Learning, Video Analysis, Temporal Segmentation, Dynamic Programming, Optimization

## 1 INTRODUCTION

With video content rapidly growing in quantity and availability, it becomes crucial to develop the relevant technologies to analyze, classify, and understand the content in videos. However, one of the biggest issues when dealing with videos is analyzing the temporal aspect. When dealing with heterogeneous video content, it is crucial to be able to partition a video into semantic scenes before performing any sort of algorithmic analysis. Besides contextual analysis, division to scenes can facilitate automatic construction of a table of contents, video summarization, chapter skimming, and more [12, 18, 19, 31].

Video scenes are an ingrained part of the hierarchical structure of videos. At the finest level of division, a video is composed of a series of images called *frames*. A sequence of frames captured from the same camera at the same time is called a *shot*. Identifying the shot transitions is considered somewhat a solved problem due to the relative uniformity of the frames in a shot [33], and established methods can be used off-the-shelf with impressive performance [1, 5]. A group of shots relating a specific event or narrative is called a *scene*. A formal definition of a scene is given by [29], as a sequence of semantically related and temporally adjacent shots depicting a high-level concept or story. Identifying the transition locations

between scenes is considered a much higher-level problem, and which is the focus of this work.

Recently, we proposed Optimal Sequential Grouping (OSG) [26] as an effective deterministic optimization formulation to solve the video scene detection problem. The approach takes the distance matrix of the shot representations and calculates the optimal division given a cost function on the intra-scene distances [28]. Despite its generality and strengths, the formulation leaves no room for learning from examples.

For learning from examples, deep learning has risen in popularity in recent years as a leading technology in many fields, and doubly so in the field of computer vision [38]. However, it can be beneficial to combine learning with analytical deterministic algorithms to gain the advantages of both learning from examples and incorporating designer knowledge and expertise [9, 10, 16, 25]. Merging learning with deterministic formulations can help arrive at more explainable technologies, ensure validity of performance, guide parameter learning, and support generalization as opposed to memorising.

Therefore, in this work, we present an approach to integrate the OSG formulation into a deep learning setting. Our model retains the original strengths of attaining an optimal division given the defined cost function, but additionally has the ability to learn better representations given annotated scene divisions.

We present a number of possible configurations for integrating OSG into the learning regime with different levels of integration. First, we present the use of the triplet loss [30], for a classical learning approach to train an embedding with valuable properties for division into scenes. This embodies the most straightforward and logical approach, but does not incorporate directly the properties of OSG. Next, we present a tailored loss directly on the distance matrix values. This loss is aimed to provide the input data in a representation which is favorable for the OSG formulation. Finally, we present an approach where learning is performed through the OSG pipeline and dynamic programming formulation to allow direct learning from results.

Figure 1 depicts our OSG model with the possible configurations. We analyze the performance and results of the different approaches and configurations. Besides out-performing the state of the art, we show how the different configurations function and analyze the behavior, benefits, and advantages.

## 2 PREVIOUS WORK

In this section we review some of the recent work on video scene detection where the task is focused on creating a complete partitioning of a motion-picture film using visual features. For a more complete
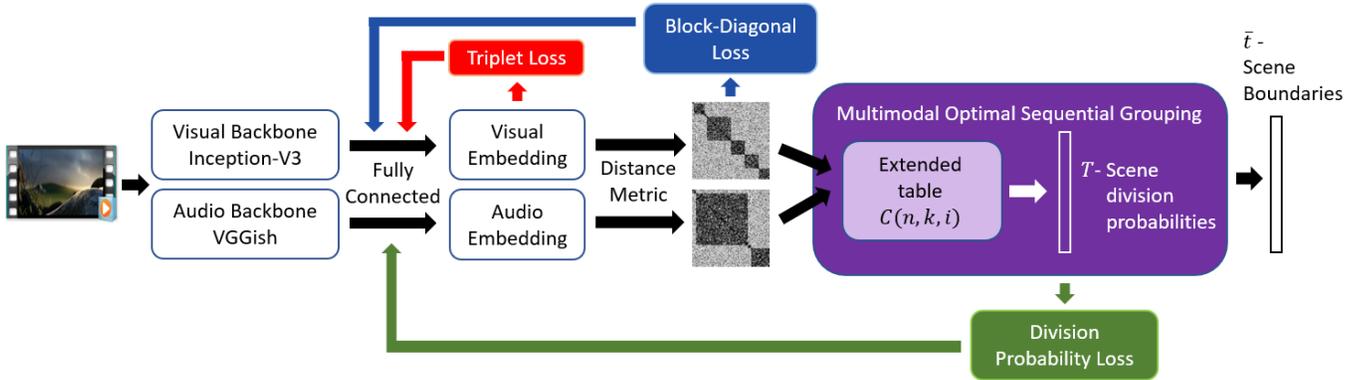
Figure 1: Our three configurations of OSG incorporated into a neural network: 1) the triplet loss (in red, OSG-Triplet), or 2) block-diagonal loss (in blue, OSG-Block), or 3) scene division probability loss (in green, OSG-Prob). The latter includes learning through the pipeline of the OSG dynamic programming algorithm, and aggregating the probabilities for division at specific locations. Video frame © Blender Foundation | gooseberry.blender.org.

review including, for example, transcript-based approaches, news segmentation, and scene retrieval, see [8].

## 2.1 Video Scene Detection

*2.1.1 Unsupervised Approaches.* The prior art of video scene detection consists of mostly unsupervised approaches even in the most recent works.

A prevalent approach for scene detection is to perform a variety of clustering techniques [2, 4, 21]. By representing video shots in some feature space the assumption is that shots from the same scene will cluster together. The weakness with such an approach is that the temporal aspect is not an inherent part of the formulation and is usually either enforced by post-processing or integrated into the feature space (as weighting, or as an additional dimension) instead of being an integral aspect of the problem.

Graph approaches [17, 23, 24] denote shots as nodes in a graph and perform graph analysis algorithms to determine the scene transitions using the graph cut algorithm. Additionally, [32, 39] construct Scene Transition Graphs by representing clusters of shots as nodes and calculating a cumulative confidence for the locations of scene divisions leveraging the primary set algorithm.

Regarding other advanced methods, [7] perform sequence alignment on shots categorized by clustering to identify recurring themes and production rules. [14] group shots with a bag of visual words descriptor and perform a sliding window for combining shots or short scenes together. [13] perform dynamic programming with a heuristic search scheme of boundaries calculated by linear discriminant analysis over shot similarities.

A specific brand of video scene detection which is of high interest focuses on egocentric videos [11, 20, 22]. Despite the overlap, the type of challenges and the level of variability in a scene when captured by an egocentric camera are not comparable to the complexity of a movie scene which demands a higher level of semantic understanding.

*2.1.2 Deep Learning for Video Scene Detection.* Regarding methods for video scene detection which incorporate deep learning, one less
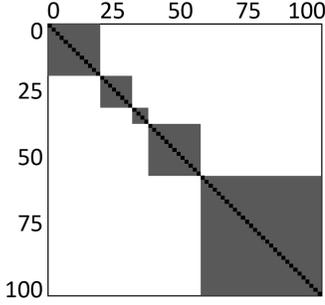
recent but relevant method [3] performs learning a distance measure using a deep siamese network and applies spectral clustering to approximate scene boundaries. They learn a joint representation of visual features and textual features (obtained from video transcription) for a similarity metric to represent the video. This method has the most in common with our approach. The main differences are that the authors do not incorporate the learning pipeline into the scene division. The learning is performed to train a distance metric, but not tailored to the spectral clustering segmentation or learned backward directly from the segmentation results. Therefore the learning could be seen as detached from the division stage, similar to the triplet configuration we present below (see Section 4.1).

In [36], the authors use deep visual, CSIFT, and MFCC audio features to represent shots. They apply a CNN architecture to each input modality and train an LSTM model to output whether each shot is a scene transition or not. This is one of the most advanced approaches with learning video scene detection in a straight-forward manner. However, results are comparable to [3] which means there is most likely room for improvement.

## 2.2 Optimal Sequential Grouping

In this work, we focus on our recently proposed method for video scene detection: Optimal Sequential Grouping (OSG).

In [26], we presented OSG as a dynamic programming algorithm to divide a video by finding the optimal solution to an additive cost function. The cost function sums the *block-diagonal* of the distance matrix which represents the intra-scene distances between shots. We additionally presented a *log-elbow* method to estimate the number of scenes directly from the distance matrix. We extended this work [27] to utilize multiple modalities in the OSG framework. By using an intermediate fusion approach, we merged the separate sequential grouping divisions into a single decision. In [28], we presented a new normalized cost function with analytically superior mathematical properties, and used deep features to represent the videos. Despite the beneficial mathematical properties, we chose to forgo using the normalized cost function in this work, because the normalization adds an additional computation complexity to

**Figure 2: A depiction of how an ideal distance matrix might look for a heterogeneous video. In this depiction, higher values are assigned brighter intensity. A dark block is a sequence of shots with low intra-distances which likely indicates a scene.**

the dynamic programming solution. We believe that when incorporating learning into OSG (as we detail in this paper), the resulting distance values will likely overcome the mathematical bias, making the choice of cost function less critical.

The technical details of the formulation and solution of OSG are expanded on in Section 3.

Our contributions are as follows: (1) We are one of the first to explore deep learning for video scene detection specifically on real-world motion-picture films (as opposed to egocentric videos, sports, news, etc.). (2) We present three configurations for combining learning into the OSG pipeline, with varying degrees of integration and tailored losses. (3) We evaluate the various approaches and analyze the advantages of the different techniques.

## 3 OPTIMAL SEQUENTIAL GROUPING

In this section we detail briefly the formulation and solution of Optimal Sequential Grouping (OSG). For more in-depth details see our previous publication [26]. Due to the generality of the approach, the description refers to a sequence of feature vectors undergoing partitioning into groups. For the task of video scene detection each feature vector describes a shot and the groups are the resulting scenes (see details in Section 5.1).

Intuitively, when representing a video containing scenes as a distance matrix, we expect to see a *block-diagonal* structure (see Figure 2). This structure is formed by the fact that shots belonging to the same scene will likely have lower distance values than shots belonging to different scenes. OSG is a dynamic programming algorithm which finds the block-diagonal with the lowest intra-scene distances.

We denote a sequence of $N$ feature vectors $X_1^N = (x_1, \ldots, x_N)$ where $x_i \in \mathbb{R}^d$, $d$ is the feature vector length. A partitioning of the sequence into $K \leq N$ groups is given by $\bar{t} = (t_1, \ldots, t_K)$, where $t_i \in \mathbb{N}$ denotes the index of the last feature vector in group $i$. A distance metric $\mathcal{D}(x_{j_1}, x_{j_2})$ measures the dissimilarity between two feature vectors. These distances guide a cost function $\mathcal{H}(\bar{t}) \in \mathbb{R}$ which measures the loss of a given division. The goal of OSG is to find $t^* = \arg\min(\mathcal{H})$ as the optimal division of $X_1^N$.

The additive cost function for a given division is defined as:

$$\mathcal{H}(\bar{t}) = \sum_{i=1}^{K} \sum_{j_1,j_2=t_{i-1}+1}^{t_i} \mathcal{D}(x_{j_1}, x_{j_2}), \tag{1}$$

where the abbreviated notation of the double sum indicates that $j_1$ and $j_2$ run from $t_{i-1} + 1$ to $t_i$ each. This cost function sums all of the intra-group distances over all of the groups in the division. Intuitively, this cost function finds a low-valued block diagonal as illustrated in Figure 2.

To find the optimal division $t^*$, we build the following recursive dynamic programming table:

$$C(n,k) = \min_i \left\{ \sum_{j_1,j_2=n}^{i} \mathcal{D}(x_{j_1}, x_{j_2}) + C(i+1, k-1) \right\}. \tag{2}$$

Here, $C(n,k)$ is the optimal cost when dividing $X_n^N$ into $k$ groups. Essentially, we find the best cost for dividing a sub-sequence which begins at index $n$, where $i$ is the location of the first point of division for this sub-sequence. The initialization:

$$C(n,1) = \sum_{j_1,j_2=n}^{N} \mathcal{D}(x_{j_1}, x_{j_2}), \tag{3}$$

is the cost of a sub-sequence starting at $n$ without any divisions. Building the table with ascending $k = 2 \ldots K$ (rising number of divisions) and descending $n = N \ldots 1$ (increasingly longer sequences) allows us to utilize the table to aggregate the partial solutions. Therefore we have that: $C(1,K) = \mathcal{H}(t^*)$, and we can reconstruct $t^*$ by storing the indexes of the chosen divisions from (2).

The number of divisions $K$ is estimated using the log-elbow approach [26, 35]. To this end, the singular values of the distance matrix are computed, and the plot of the log values is analyzed. The point of plateau ('elbow') in the plot was shown to correspond to the number of blocks with intuition from performing a low-rank matrix approximation. See Appendix B for details on how the elbow point is estimated.

When incorporating multiple modalities [27], the distance for each modality is used to build its own table $C_x$, $C_y$, where the subscript indicates the modality, and $Y_1^N = (y_1, \ldots, y_N)$ is an additional modality. Instead of choosing the point of division which yields the lowest cost for a single modality, the modality which has a more pronounced division point is chosen. We define:

$$G_x^{n,k}(i) = \sum_{j_1,j_2=n}^{i} \mathcal{D}_x(x_{j_1}, x_{j_2}) + C_x(i+1, k-1), \tag{4}$$

which is the argument of the minimum function in (2). $G_x^{n,k}$ is normalized to indicate the relative inclination for division:

$$\hat{G}_x(i) = \frac{G_x(i) - \text{mean}\{G_x\}}{\text{std}\{G_x\}}, \tag{5}$$

and the index is chosen as: $\arg\min_i \left\{ \min(\hat{G}_x(i), \hat{G}_y(i)) \right\}$ (superscripts were omitted for the sake of readability).

## 4 LEARNABLE OSG

Despite the strengths of OSG as an unsupervised optimization scheme, the main weakness is the dependency on choosing the representative features $X_1^N$ and distance metric $\mathcal{D}$. Here, deep learning as a data representation mechanism can be a powerful tool when joined with OSG. In this section we detail three possible configurations for joining learning with the OSG algorithm. In all of the sections below, we take the shot representations $X_1^N$ and feed them through a series of fully connected layers to learn a new representation $\widetilde{X}_1^N$ (in the notations below, we omit the tilde for simplicity). These parameters are what the network learns to better perform OSG.

### 4.1 Cluster Embedding (OSG-Triplet)

The most direct way to apply learning to the OSG problem would be with learning an embedding. Specifically, the triplet loss [30] learns a feature space embedding where samples from the same class are close in the feature space while samples from different classes are further apart. This is useful for a range of tasks, but for scene division this is doubly intuitive because the triplet loss causes samples (shots, in this case) to cluster together (see Appendix A).

These clusters will likely make scene detection a much simpler task, because often the task is approached as a variant of a shot clustering problem. In an embedding where shots are clustered into scenes, we can assume that the distance matrix will possess beneficial properties for OSG. Likely, the intra-scene distances will be reduced compared to the inter-scene distances causing the dynamic programming algorithm to arrive at the correct divisions.

Given a label $L(x_i) \in [1, K]$ indicating the number of the scene that feature vector $x_i$ belongs to, the neural network parameters are learned by minimizing the triplet loss:

$$\sum \min(\mathcal{D}(x_i, x_i^p) - \mathcal{D}(x_i, x_i^n) + \alpha, 0). \tag{6}$$

For anchor samples $x_i$, a positive and negative pair are chosen, where $L(x_i) = L(x_i^p)$ and $L(x_i) \neq L(x_i^n)$, and $\alpha$ is a margin parameter. The samples are chosen using the semi-hard approach, where the triplets that satisfy the condition $\mathcal{D}(x_i, x_i^p) < \mathcal{D}(x_i, x_i^n) < \mathcal{D}(x_i, x_i^p) + \alpha$ are chosen.

As stated above, this approach is intuitive and likely to aid OSG in division. In the next configurations we show how we go further to tailor the learning specifically for the OSG formulation.

### 4.2 Block-Diagonal Loss (OSG-Block)

In Section 3, we described the intuition behind OSG as identifying the block-diagonal structure in the distance matrix. In this configuration, we apply a loss designated to strengthen that block-diagonal structure.

If we present the distance values in a matrix $D$, where the $i$-th row and $j$-th column is $D_{i,j} = \mathcal{D}(x_i, x_j)$, then we can define an 'optimal' $D^*$ as:

$$D_{i,j}^* = \begin{cases} 0 & L(x_i) = L(x_j) \\ 1 & \text{else} \end{cases}. \tag{7}$$

Here, 0 is the minimal distance and is allocated for features from the same scene, and 1 is the maximal distance for features from different scenes (see Figure 3).



**Figure 3: $D^*$. For OSG-Block the entire matrix is used, while for OSG-Block-Adjacent only the gray (dark and light) portions are considered.**

OSG does not need an optimal $D$ matrix to perform well. The relative divisions are compared to each other so the correct solution only needs to have a slightly lower cost than any other solution. However, driving $D$ toward $D^*$ will likely help OSG find the the correct division. Therefore, the loss we use is the Frobenius norm of the subtraction:

$$\|D - D^*\|_F = \sqrt{\sum_i \sum_j \left| D_{i,j} - D_{i,j}^* \right|^2}. \tag{8}$$

A slight variant of this loss is to not consider the inter-scene distances between scenes which are not adjacent to each other (OSG-Block-Adjacent). The rational is that some scenes throughout a video might be quite similar to each other, but their temporal distance or an intervening scene will indicate their distinction. The cost function in OSG accumulates the inner values of the block-diagonal, while the far off-diagonal values do not impact the decision as long as the values in between are high enough.

In this case, the loss receives only a portion of the values. Specifically, in (8), we only consider values of $j$ that satisfy the constraint: $L(x_i) - 1 \leq L(x_j) \leq L(x_i) + 1$. I.e., only the intra-scene distances and distances between feature vectors belonging to neighboring scenes are considered (see Figure 3).

### 4.3 Scene Division Probabilities (OSG-Prob)

In this configuration, the learning process is performed through the OSG pipeline. The OSG formulation is altered slightly to allow division probabilities to be calculated, and this is contrasted to the ground truth divisions. The model then learns to raise the probability for division at the correct location.

In (2), the $C$ table is used to calculate optimal locations of division. As in (4), we retain the relative inclinations for division. Instead of (5), we output a probability vector with the established softmin operator and aggregate the values in a larger table:

$$C(n, k, i) = \frac{\exp(-G^{n,k}(i))}{\sum_j \exp(-G^{n,k}(j))}. \tag{9}$$

The values in this table retain the probability to divide the video at point $i$ when dividing $X_n^N$ into $k$ scenes. We average these probabilities in the $C$ table over $n$ and $k$ and arrive at a vector of 'scores'

for division at each location in the video:

$$T(i) = \frac{1}{N \cdot K} \sum_n \sum_k C(n, k, i).$$ (10)

Given the probabilistic nature of the values, we opt to use the cross-entropy loss on the probabilities at the indexes where a division is annotated:

$$-\sum_{i \in \bar{t}_{GT}} \log(T(i)).$$ (11)

Where $\bar{t}_{GT} = \{i | L(x_{i+1}) > L(x_i)\}$ is the ground truth division.

We note that there is no inherent problem with evaluating $T$ only on the ground truth division indexes. The network cannot learn a trivial $T \equiv 1$, as high values in the $C$ table imply directly that other locations have lower values due to the softmin operation. For a location to arrive at a high average probability, it means there must be a comprehensive indication for a division at that index, and this is what the configuration attempts to create by learning.

## 5 EVALUATION AND ANALYSIS

In this section we evaluate and analyze the performance of the proposed configurations.

### 5.1 Technical Details

We use a pre-trained Inception-v3 architecture [34] as a 2048-dimension visual backbone feature extractor from images, and we use a pre-trained VGGish network [15] to encode the audio segments into 128-dimension vectors. In our experiments, we used four fully connected (FC) layers, $(3000, 3000, 1000, 100)$ for visual and $(200, 200, 100, 20)$ for audio. Batch normalization was applied on all layers, and ReLU activations were applied on all layers excluding output. The ADAM optimization algorithm was used to train the network with a learning rate of $5 \cdot 10^{-3}$. A stopping criteria to avoid overfitting was used, and aborted the learning process when the training loss decreased to 25% of its initial value. The cosine distance normalized between 0 and 1 was used as $\mathcal{D}$, the margin $\alpha$ was chosen as 0.5, and the log-elbow approach was used to estimate the number of scenes $K$.

Regarding runtime constraints, complexity of the OSG stage is unchanged compared to the original publication [26]. The embedding network is relatively light-weight and the addition to the forward pass compared to the backbone is negligible. Training took roughly 24 hours on a single GPU.

For video scene detection we used the OVSD dataset [27]. The dataset contains 21 full-length motion-picture films from a variety of genres with ground truth scene labeling (see Appendix C for dataset details). For each video, we perform shot boundary detection [5] and extract a center image for the visual representation fed into the Inception network. Audio for each 0.96 seconds was encoded using the VGGish network and average pooling was applied to encode each shot with its relevant audio representation. These features were used to provide a fair comparison to [28]. Better performance can likely be attained by incorporating advanced representations such as an I3D network [6].

In order to compare results on the entire OVSD dataset, we aimed to show 'test' performance on all of the videos. To accomplish this, we incorporated a 5-fold testing approach, where the videos were split into five groups of roughly equal size with regard to number of shots, i.e., some groups consisted of fewer but longer videos while others consisted of more videos, each with less shots (see Table 2 in the video name subscripts for the division to groups). Five separate models were each trained on four-fifths of the data, leaving out one group for testing. At test time, the five models were applied each to its test group, and scores were averaged over all of the videos.

### 5.2 Configuration Analysis

In this section, we analyze the behavior of the various configurations.

In Table 1, we present various stages of $D$ from visual features of the video Meridian from OVSD with the accompanied ground truth $D^*$. This video contains the smallest amount of shots, and offers the ability to visually and qualitatively inspect the structure of the matrix $D$ and behavior of the configurations. 'Orig' displays the distance metric applied directly to the backbone features. Since the features were chosen to provide a fair comparison, this is exactly the matrix that the method in [28] applies OSG on. 'Epoch 0' is the matrix from the features after applying an untrained embedding network and incurring a level of noise. On the right are the matrix after 20 epochs under the different configurations, with the gradients below each matrix.
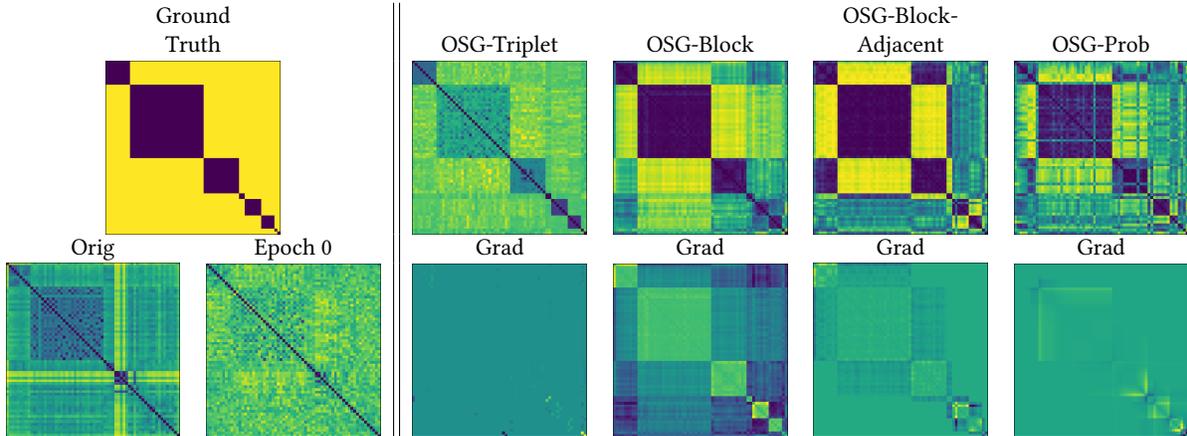
OSG-Triplet and OSG-Block both explicitly strive to minimize the distances between shots from the same scene and raise the distances between shots belonging to different scenes. The main difference between these approaches is strengthening the block-diagonal structure to better help OSG performance. While the triplet loss focuses on distinct samples, the block-diagonal loss concentrates on the complete structure and raises the values outside of the block-diagonal.

It is interesting to note the balance between how the configurations emphasize small versus large scenes. OSG-Triplet based off of distinct samples manages to emphasize the small scenes and has difficulty with long scenes. OSG-Block-Adjacent compared to OSG-Block dismisses the far off-diagonal blocks, focuses on the areas which are more important (the areas between adjacent scenes), and manages to accentuate both large and small scenes. OSG-Prob seems to converge more slowly but consistently over the scenes.

One interesting aspect to explore is how the losses affect the gradients of the distance matrix. In Table 1 (right, bottom), we see the map of the gradient values for the various configurations. As can be expected, OSG-Triplet is dependant on individual values of distance, OSG-Block puts emphasis on the entire block-diagonal, and OSG-Block-Adjacent on the relevant section of the block-diagonal as defined by the ground truth, while OSG-Prob has a much more 'local' impact focused around the points of division. For OSG-Prob, this is a direct outcome of the formulation which emphasizes the value of the average probability on the scene division.

OVSD [27], is one of the only freely-available video scene detection datasets. Despite the substantial length and variety, the amount of data is still very limited especially when considering other deep learning tasks. The reason the network is able to learn at all, we assume, is because the configurations we chose do not treat a complete video as a sample, but rather a shot as a sample. With each scene acting as a label instead of the entire division being

**Table 1: An example $D$ from OVSD. On the left: Ground Truth ($D^*$), Orig (without an applied embedding), Epoch 0 (embedding before learning). On the right, trained examples after 20 epochs for: OSG-Triplet, OSG-Block, OSG-Block-Adjacent, and OSG-Prob, with corresponding gradients (bottom row)**



considered a label, the models manage to generalize the important elements which represent shots belonging to the same scene.

Regarding the behavior of $T$, in Figure 4 we show the progression of the values of $T(i)$ over a number of iterations. We can see that as the iterations progress, $T$ raises the probability at the ground truth points of division. The probability is lowered for locations with no true division even though this is not specifically enforced by the loss but rather an outcome of the construction of $T$ (see Section 4.3). Additionally we see that the model has difficulty on the last small scenes which are more difficult to enforce.

## 5.3 Baselines

As explained in Section 1, the motivation to integrate learning into the OSG pipeline is to leverage the strengths of the deterministic formulations together with the ability to learn from examples. To emphasise this point, besides comparing to state-of-the-art methods, we implement two 'pure' deep learning baselines.

The first is a sliding window approach. Using a window of $W$ consecutive feature vectors (representing shots), the network is trained to identify when the scene transition is precisely in the middle of the window. This is a naive but straightforward way to apply learning to the problem but without leveraging a formulation such as OSG. To represent a fair comparison, we used all the same parameters as the embedding network detailed above. $W = 4$ embedded feature vectors were concatenated and a final FC layer to size 1 was added followed by a sigmoid output (additional values of $W$ were tested resulting in comparable or worse performance).

The second baseline we used is a Long Short-Term Memory (LSTM) recurrent neural network architecture. The LSTM is a classic choice for modeling sequence-to-sequence problems, and presents a more advanced baseline for comparison. A bi-directional LSTM component was used with a hidden state of length 1000, followed by two FC layers sized 100 and 1, the former with a ReLU activation and the latter with a sigmoid output. The network processes the sequence of feature vectors and for each time step outputs the probability for the current feature representing the end of a scene.

Both baselines were trained using the same methodology as our configurations. As an unfair advantage, sigmoid threshold was chosen as the value which maximized test performance. Table 2 under the 'Supervised' column shows the sliding window (SW) and LSTM results.

## 5.4 Scene Detection Evaluation

We measure the performance of our OSG configurations on the OVSD dataset [27]. For a metric, we use the widely accepted Coverage $C$ and Overflow $O$ [37], with a single value $F$-score for assessing the quality of division as the harmonic mean between $C$ and $1 - O$.

Figure 5 on the left presents the average $F$-score for the various configurations on the tested videos. We present the performance when using the visual or audio features separately, and when performing OSG with the multimodal fusion approach (see Section 3). The results show superior performance for OSG-Prob, and specifically the multimodal fusion approach, which is more preferable than using a single modality for most of the configurations.

As an analysis of the behavior of each modality, we divided the performance of OSG-Prob per genre of OVSD (Figure 5 on the right). It can be noted that for documentaries where it is characteristic for the visuals to change often, but for speakers to stay constant within a scene, the audio modality played a vital role in division. On the other hand, for comedy, crime, and animation, the visual aspect played a slightly more important role. In drama, both modalities contributed greatly, which coincides with the fact that often the changing visuals are accompanied with matching auditory ambiance in this genre.

Table 2 presents the $F$-scores over all of the videos in the OVSD dataset with our OSG configurations leveraging multimodal fusion of visual and audio features. As a comparison, we show the 'Prior Art' column which are two state-of-the-art unsupervised methods [2, 32]. The 'OSG Prior Art' column [26–28], are the prior art on OSG. Specifically, as stated above in Section 5.1, our backbone features are the same as [28]. Therefore, this can be seen as directly comparable to the case when no learning is performed.
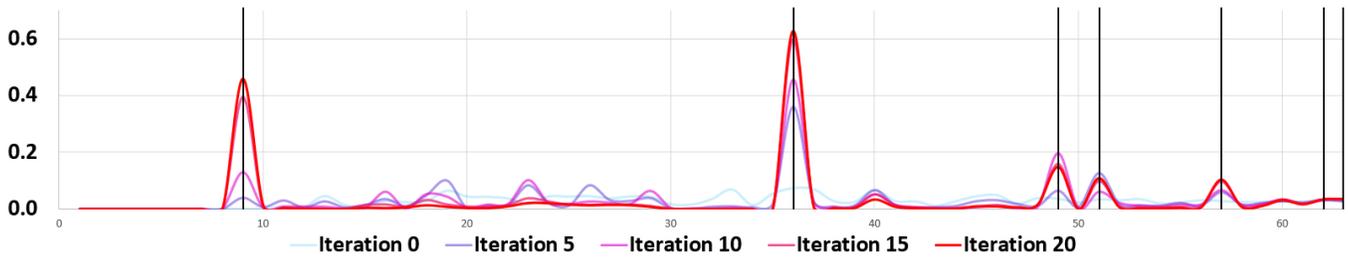
**Figure 4: Progression of $T$ as a function of $i$ (shot number) over a number of iterations. Graphs go from translucent blue to opaque red as iterations progress (best viewed in color). Vertical black lines indicate ground truth divisions.**
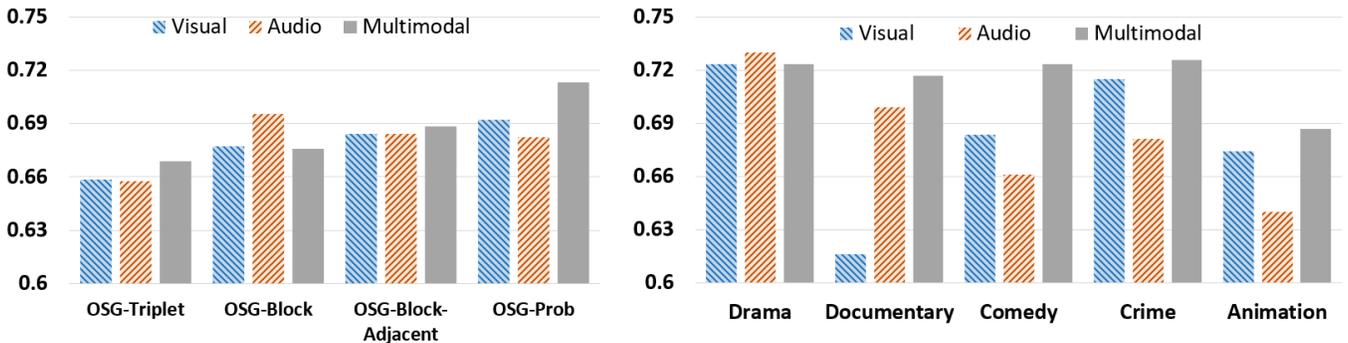


**Figure 5: Modality analysis. Average $F$-score when using visual, audio, or a multimodal fusion. Performance per configuration (left), and performance of OSG-Prob consolidated per genre (right).**

Additionally, the 'Supervised' column presents the results of our implemented baselines as detailed in Section 5.3. These represent applying learning to the problem without leveraging a strong deterministic formulation such as OSG.

It is interesting to note the substantial increase in performance of the configurations compared to the prior art on OSG. But even more so, when comparing to the performance of the supervised baselines, we can clearly see the benefit of leveraging both learning and OSG. Despite OSG-Prob attaining the best results, it is not directly our intention to promote it as the only viable option. Indeed, with closely comparable results, we felt that it would be beneficial for the advancement of future work to present multiple options, as opposed to promoting a single architecture. When extending or applying our work to future problems, it can be beneficial to choose the relevant configuration for the problem while being able to understand the behavior and trade-offs.

Figure 6 shows results on part of a video from the OVSD dataset (additional results in Appendix G). In general, we can see divisions which result in reasonable and often precise scene divisions. Using these divisions for applying video understanding and classification technologies will undoubtedly be superior over applying them on the entire video or on naive uniform divisions.

## 6 CONCLUSION

In this work, we have presented a novel approach for incorporating learning into OSG for the task of video scene detection. We presented a number of different configurations, evaluated their performance, and analyzed their behavior and various advantages.

Overall, our goal was to explore the possibility of integrating learning into the powerful formulation of OSG, so as to merge learning models with this effective analytic algorithm. We demonstrated this ability with varying amounts of model complexity and dependence on the OSG pipeline. Beyond creating a more precise and robust video scene detection technology, we believe this approach can enable constructing a model with designated performance on specific content or genres. Additionally, this model could be integrated into other video models or leveraged for additional video understanding tasks. We hope this work encourages continued research in the field of temporal video analysis.

## REFERENCES

[1] Evlampios Apostolidis and Vasileios Mezaris. 2014. Fast shot segmentation combining global and local visual descriptors. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 6583–6587.

[2] Lorenzo Baraldi, Costantino Grana, and Rita Cucchiara. 2015. Analysis and Re-Use of Videos in Educational Digital Libraries with Automatic Scene Detection. In *11th Italian Research Conference on Digital Libraries*. Springer, 155–164.

[3] Lorenzo Baraldi, Costantino Grana, and Rita Cucchiara. 2015. A Deep Siamese Network for Scene Detection in Broadcast Videos. In *Proceedings of the 23rd ACM International Conference on Multimedia* (Brisbane, Australia) *(MM '15)*. ACM, New York, NY, USA, 1199–1202. https://doi.org/10.1145/2733373.2806316

[4] Lorenzo Baraldi, Costantino Grana, and Rita Cucchiara. 2015. Measuring scene detection performance. In *Iberian Conference on Pattern Recognition and Image Analysis*. Springer, 395–403.

[5] Lorenzo Baraldi, Costantino Grana, and Rita Cucchiara. 2015. Shot and scene detection via hierarchical clustering for re-using broadcast video. In *International Conference on Computer Analysis of Images and Patterns*. Springer, 801–811.

[6] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6299–6308.

**Table 2: *F*-score results on OVSD. Best score per video in bold. Subscript on video name indicates the group for 5-fold testing**

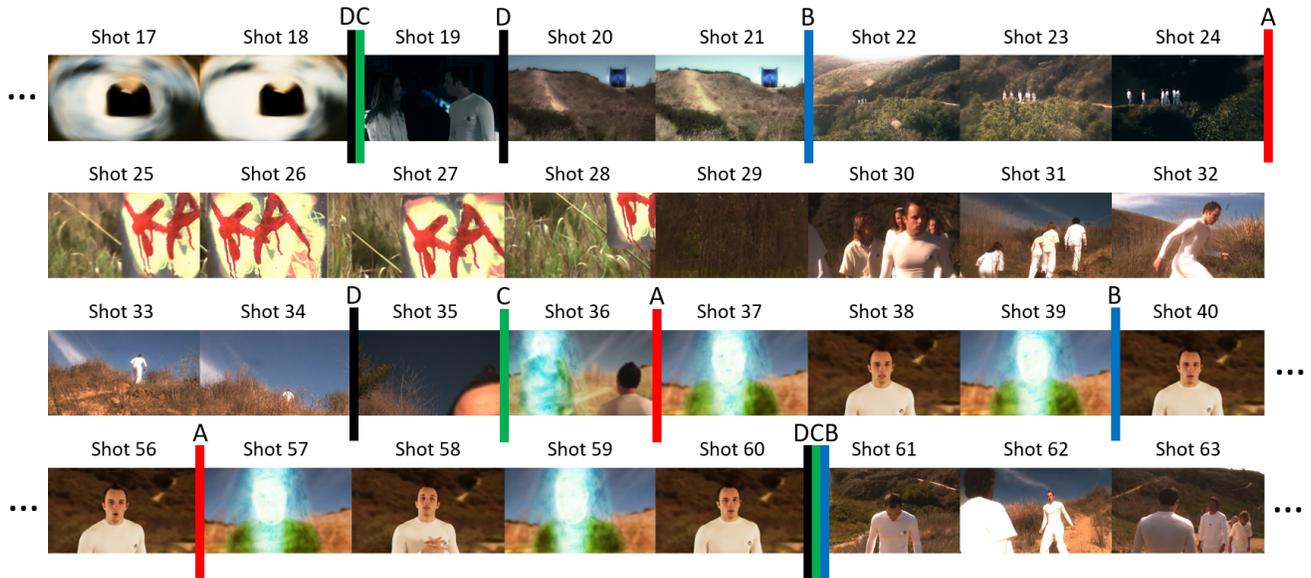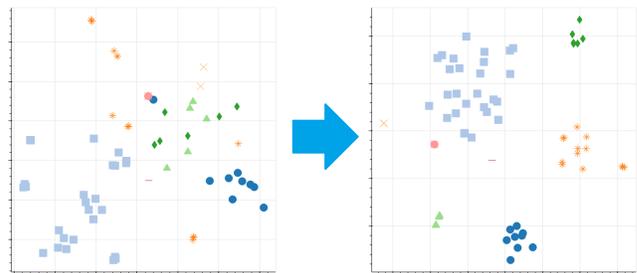| Video Name | This Work | | | | Supervised | | OSG Prior Art | | | Prior Art | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | OSG-Prob | OSG-Block-Adjacent | OSG-Block | OSG-Triplet | SW | LSTM | [28] | [27] | [26] | [32] | [2] |
| $1000_1$ | 0.70 | **0.74** | 0.67 | 0.68 | 0.49 | 0.38 | 0.60 | 0.38 | 0.57 | 0.50 | 0.39 |
| $BBB_5$ | 0.76 | 0.74 | 0.77 | 0.81 | 0.66 | 0.54 | 0.63 | **0.83** | 0.69 | 0.49 | 0.46 |
| $BWNS_1$ | **0.80** | 0.74 | 0.71 | 0.75 | 0.65 | 0.62 | 0.70 | 0.63 | 0.20 | 0.61 | 0.43 |
| $CH7_1$ | 0.72 | **0.74** | 0.66 | 0.67 | 0.56 | 0.45 | 0.60 | 0.63 | 0.49 | 0.52 | 0.26 |
| $CL_2$ | **0.68** | 0.55 | 0.57 | 0.49 | 0.60 | 0.43 | 0.51 | 0.53 | 0.53 | 0.45 | 0.07 |
| $ED_2$ | 0.70 | 0.68 | 0.64 | **0.73** | 0.55 | 0.31 | 0.61 | 0.6 | 0.69 | 0.56 | 0.55 |
| $FBW_5$ | 0.76 | 0.77 | **0.78** | 0.76 | 0.59 | 0.63 | 0.59 | 0.57 | 0.14 | 0.61 | 0.52 |
| $Honey_2$ | **0.74** | 0.66 | 0.67 | 0.73 | 0.58 | 0.26 | 0.63 | 0.58 | 0.38 | 0.38 | 0.36 |
| $JW_2$ | **0.79** | 0.65 | 0.65 | 0.65 | 0.62 | 0.29 | 0.63 | 0.75 | 0.64 | 0.28 | 0.22 |
| $LCDP_5$ | **0.73** | 0.60 | 0.60 | 0.61 | 0.61 | 0.45 | 0.72 | 0.53 | 0.42 | 0.18 | 0.22 |
| $LM_5$ | 0.73 | **0.74** | 0.64 | 0.65 | 0.60 | 0.60 | 0.64 | 0.69 | 0.25 | 0.71 | 0.28 |
| $Meridian_2$ | 0.64 | 0.69 | 0.68 | 0.69 | 0.47 | 0.66 | 0.79 | 0.63 | 0.71 | 0.63 | **0.82** |
| $Oceania_3$ | 0.73 | 0.77 | **0.78** | 0.68 | 0.40 | 0.62 | 0.61 | 0.67 | 0.45 | 0.51 | 0.26 |
| $Pentagon_3$ | 0.65 | 0.65 | 0.64 | **0.80** | 0.57 | 0.61 | 0.65 | 0.73 | 0.18 | 0.48 | 0.16 |
| $Route 66_3$ | 0.67 | 0.71 | 0.66 | **0.72** | 0.44 | 0.19 | 0.60 | 0.54 | 0.31 | 0.36 | 0.05 |
| $SDM_4$ | **0.82** | 0.72 | 0.80 | **0.82** | 0.51 | 0.53 | 0.70 | 0.68 | 0.67 | 0.81 | 0.55 |
| $Sintel_5$ | 0.59 | 0.59 | 0.60 | 0.48 | 0.43 | 0.54 | 0.58 | 0.46 | **0.66** | 0.59 | 0.51 |
| $SStB_3$ | **0.66** | 0.51 | 0.65 | 0.44 | 0.34 | 0.51 | 0.62 | 0.46 | 0.48 | 0.43 | 0.22 |
| $SW_4$ | **0.72** | 0.71 | 0.67 | 0.66 | 0.52 | 0.59 | 0.61 | 0.55 | 0.36 | 0.40 | 0.13 |
| $ToS_5$ | 0.66 | **0.78** | 0.56 | 0.53 | 0.65 | 0.55 | 0.73 | 0.5 | 0.62 | 0.75 | 0.23 |
| $Valkaama_5$ | 0.74 | 0.71 | **0.78** | 0.69 | 0.66 | 0.46 | 0.70 | 0.63 | 0.73 | 0.73 | 0.16 |
| Average | **0.71** | 0.69 | 0.68 | 0.67 | 0.55 | 0.49 | 0.64 | 0.60 | 0.52 | 0.46 | 0.33 |



**Figure 6: Qualitative results on sections of Jathia's Wager from OVSD dataset. Division marked by A OSG-Triplet (red) B OSG-Block (blue) C OSG-Prob (green) and D Ground truth (black) - See Appendix G. Images © 2009, Solomon Rothman.**

[7] Vasileios T Chasanis, Aristidis C Likas, and Nikolaos P Galatsanos. 2008. Scene detection in videos using shot clustering and sequence alignment. *IEEE transactions on multimedia* 11, 1 (2008), 89–100.

[8] Manfred Del Fabro and Laszlo Böszörmenyi. 2013. State-of-the-art and future challenges in video scene detection: a survey. *Multimedia systems* 19, 5 (2013), 427–454.

[9] Diego Didona, Francesco Quaglia, Paolo Romano, and Ennio Torre. 2015. Enhancing performance prediction robustness by combining analytical modeling and machine learning. In *Proceedings of the 6th ACM/SPEC international conference on performance engineering*. ACM, 145–156.

[10] Alex Endert, William Ribarsky, Cagatay Turkay, BL William Wong, Ian Nabney, I Díaz Blanco, and Fabrice Rossi. 2017. The state of the art in integrating machine learning into visual analytics. In *Computer Graphics Forum*, Vol. 36. Wiley Online Library, 458–486.

[11] Antonino Furnari, Giovanni Maria Farinella, and Sebastiano Battiato. 2016. Temporal segmentation of egocentric videos to highlight personal locations of interest. In *European Conference on Computer Vision*. Springer, 474–489.

[12] Jiyang Gao, Runzhou Ge, Kan Chen, and Ram Nevatia. 2018. Motion-appearance co-memory networks for video question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6576–6585.

[13] Bo Han and Weiguo Wu. 2011. Video scene segmentation using a novel boundary evaluation criterion and dynamic programming. In *2011 IEEE International conference on multimedia and expo*. IEEE, 1–6.

[14] Muhammad Haroon, Junaid Baber, Ihsan Ullah, Sher Muhammad Daudpota, Maheen Bakhtyar, and Varsha Devi. 2018. Video Scene Detection Using Compact Bag of Visual Word Models. *Advances in Multimedia* 2018 (2018).

[15] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. 2017. CNN architectures for large-scale audio classification. In *Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on*. IEEE, 131–135.

[16] Alina Kloss, Stefan Schaal, and Jeannette Bohg. 2017. Combining learned and analytical models for predicting action effects. *arXiv preprint arXiv:1710.04102* (2017).

[17] Chao Liang, Yifan Zhang, Jian Cheng, Changsheng Xu, and Hanqing Lu. 2009. A novel role-based movie scene segmentation method. In *Pacific-Rim Conference on Multimedia*. Springer, 917–922.

[18] Debabrata Mahapatra, Ragunathan Mariappan, and Vaibhav Rajan. 2018. Automatic Hierarchical Table of Contents Generation for Educational Videos. In *Companion Proceedings of the The Web Conference 2018*. International World Wide Web Conferences Steering Committee, 267–274.

[19] Bernd Münzer and Klaus Schoeffmann. 2018. Video Browsing on a Circular Timeline. In *International Conference on Multimedia Modeling*. Springer, 395–399.

[20] Alessandro Ortis, Giovanni M Farinella, Valeria D'Amico, Luca Addesso, Giovanni Torrisi, and Sebastiano Battiato. 2017. Organizing egocentric videos of daily living activities. *Pattern Recognition* 72 (2017), 207–218.

[21] Rameswar Panda, Sanjay K Kuanar, and Ananda S Chowdhury. 2017. Nyström Approximated Temporally Constrained Multisimilarity Spectral Clustering Approach for Movie Scene Detection. *IEEE Transactions on Cybernetics* (2017).

[22] Yair Poleg, Chetan Arora, and Shmuel Peleg. 2014. Temporal segmentation of egocentric videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2537–2544.

[23] Stanislav Protasov, Adil Mehmood Khan, Konstantin Sozykin, and Muhammad Ahmad. 2018. Using deep features for video scene detection and annotation. *Signal, Image and Video Processing* 12, 5 (2018), 991–999.

[24] Zeeshan Rasheed and Mubarak Shah. 2005. Detection and representation of scenes in videos. *IEEE transactions on Multimedia* 7, 6 (2005), 1097–1105.

[25] Paramita Ray and Amlan Chakrabarti. 2019. A Mixed approach of Deep Learning method and Rule-Based method to improve Aspect Level Sentiment Analysis. *Applied Computing and Informatics* (2019).

[26] Daniel Rotman, Dror Porat, and Gal Ashour. 2016. Robust and efficient video scene detection using optimal sequential grouping. In *2016 IEEE International Symposium on Multimedia (ISM)*. IEEE, 275–280.

[27] Daniel Rotman, Dror Porat, and Gal Ashour. 2017. Robust video scene detection using multimodal fusion of optimally grouped features. In *2017 IEEE 19th International Workshop on Multimedia Signal Processing (MMSP)*. IEEE, 1–6.

[28] Daniel Rotman, Dror Porat, Gal Ashour, and Udi Barzelay. 2018. Optimally Grouped Deep Features Using Normalized Cost for Video Scene Detection. In *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval*. ACM, 187–195.

[29] Yong Rui, Thomas S Huang, and Sharad Mehrotra. 1999. Constructing table-of-content for videos. *Multimedia systems* 7, 5 (1999), 359–368.

[30] Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 815–823.

[31] Yair Shemer, Daniel Rotman, and Nahum Shimkin. 2019. ILS-SUMM: Iterated Local Search for Unsupervised Video Summarization. *arXiv preprint arXiv:1912.03650* (2019).

[32] Panagiotis Sidiropoulos, Vasileios Mezaris, Ioannis Kompatsiaris, Hugo Meinedo, Miguel Bugalho, and Isabel Trancoso. 2011. Temporal video segmentation to scenes using high-level audiovisual features. *IEEE Transactions on Circuits and Systems for Video Technology* 21, 8 (2011), 1163–1177.

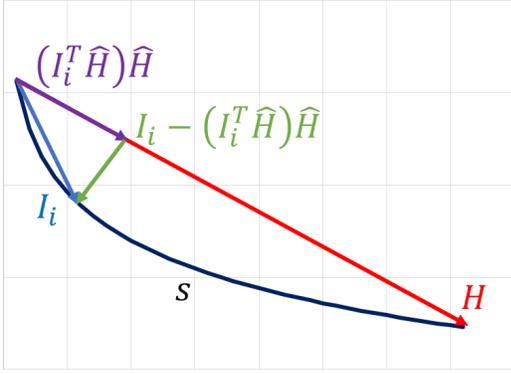[33] Alan F Smeaton, Paul Over, and Aiden R Doherty. 2010. Video shot boundary detection: Seven years of TRECVid activity. *Computer Vision and Image Understanding* 114, 4 (2010), 411–418.

[34] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2818–2826.

[35] Makarand Tapaswi, Martin Bauml, and Rainer Stiefelhagen. 2014. Storygraphs: visualizing character interactions as a timeline. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 827–834.

[36] Tiago H. Trojahn, Rodrigo M. Kishi, and Rudinei Goularte. 2018. A New Multimodal Deep-learning Model to Video Scene Segmentation. In *Proceedings of the 24th Brazilian Symposium on Multimedia and the Web* (Salvador, BA, Brazil) *(WebMedia '18)*. ACM, New York, NY, USA, 205–212. https://doi.org/10.1145/3243082.3243108

[37] Jeroen Vendrig and Marcel Worring. 2002. Systematic evaluation of logical story unit segmentation. *IEEE Transactions on Multimedia* 4, 4 (2002), 492–499.

[38] Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. 2018. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience* 2018 (2018).

[39] Minerva Yeung, Boon-Lock Yeo, and Bede Liu. 1998. Segmentation of video by clustering and graph analysis. *Computer vision and image understanding* 71, 1 (1998), 94–109.

Figure 7: A reduced 2-dimensional representation of shot feature vectors (using TSNE) from the video Meridian from OVSD. Different marker types and colors represent scenes. After applying the triplet loss (right) the scenes are better separated into clusters.

## A TRIPLET LOSS FOR VIDEO SCENE DETECTION

As mentioned in the paper, the triplet loss [30] learns a feature space embedding where samples from the same class are close in the feature space while samples from different classes are further apart. This is useful for a range of tasks, but for scene division this is doubly intuitive because the triplet loss causes samples (shots, in this case) to cluster together. In Figure 7 is a reduced 2-dimensional representation of shot feature vectors (using TSNE) from the video Meridian from the OVSD dataset. This video contains the smallest amount of shots, and offers the ability to visually and qualitatively inspect the distribution of the shot representations.

Despite the success of separating the shot representations into clusters, we can see that classic clustering algorithms might have trouble dividing correctly. Specifically we are referring to the single-shot scenes surrounding the large light blue square scene. In this instance, the OSG algorithm will likely be beneficial given the order and locations of the scenes, and the ability to make a decision based on the temporal order of the shots.

**Figure 8: A depiction of the estimation of the log-elbow plateau point in the log graph of the singular values of the distance matrix.**

## B ESTIMATING THE NUMBER OF SCENES $K$

The number of divisions $K$ is estimated using the log-elbow approach [26, 35]. To this end, the singular values of the distance matrix are computed, and the plot of the log values is analyzed. The point of plateau ('elbow') in the plot was shown to correspond to the number of blocks with intuition from performing a low-rank matrix approximation. The mathematical intuition is that given a distance matrix with a block-diagonal structure, we can see the rows of the matrix which belong to a specific block as being roughly linearly dependant. If the matrix were ideal (zeros on the block diagonal and ones outside), the rank of the matrix would be exactly the number of blocks in the block diagonal. Given a real noisy matrix, we expect the noise to act as the high frequency and low energy additions to the underlying inherent structure of the matrix. By identifying the plateau point of the singular values we can estimate the rank of the fundamental structure of the matrix.

Practically, this plateau point is located with an elbow estimation, as the point farthest from the diagonal running over the graph. Formally, if $s$ is the log singular values of length $N$ and we consider the index of each value as the first dimension, then the vector $I_i = [i, s_i]^T$ represents the values of the graph. The diagonal would be: $H = [N - 1, s_N - s_1]^T$, with $\hat{H} = H/\|H\|$ the unit vector in the same direction, and using the euclidean distance to each point and projecting the vector $I$, we can identify the index of the plateau point:

$$\text{log-elbow} = \arg\max_i \left\{ \|I_i - (I_i^T \hat{H})\hat{H}\| \right\}. \tag{12}$$

See Figure 8 for an illustration.

## C OVSD DATASET

For video scene detection we used the OVSD dataset [27]. OVSD, is one of the only freely-available video scene detection datasets allowing both academic and industrial research use (creative commons licenses). To the extent of our knowledge, this dataset is the only video scene detection dataset that has entire movies and is freely available with only minimal legal restrictions.

The dataset contains 21 full-length motion-picture films from a variety of genres with ground truth scene labeling. Table 3 presents

the details of the OVSD dataset. 'Short Name' is the name presented in the results table in the paper to conserve space, and '# shots' is the number of shots as estimated using a shot boundary detection method [5]. Some videos are defined by a number of genres (as is acceptable with films). For the analysis per genre in the paper, the first genre was used to aggregate results, where Meridian was added to Crime (being the genre closest to Mystery).

## D EVALUATION METRIC

We measure the performance of our OSG configurations on the OVSD dataset. For a metric, we use the widely accepted Coverage $C$ and Overflow $O$ [37], with a single value $F$-score for assessing the quality of division as the harmonic mean between $C$ and $1 - O$.

Formally, as in [4], we denote $s_1, s_2, \ldots, s_m$ as the series of detected scenes, and $\tilde{s}_1, \tilde{s}_2, \ldots, \tilde{s}_n$ as the series of ground truth scenes, where each element $s$ is a set of shots. The coverage $C_t$ of ground truth scene $\tilde{s}_t$ is computed as:

$$C_t = \frac{\max_{i=1,\ldots,m} \#(s_i \cap \tilde{s}_t)}{\#\tilde{s}_t}, \tag{13}$$

where $\#(s)$ is the number of shots in scene $s$. Essentially, this is the relative amount of the ground truth scene that was allocated to a single scene in the proposed division. The overflow $O_t$ for ground truth scene $\tilde{s}_t$ is computed as:

$$O_t = \frac{\sum_{i=1}^{m} [\#(s_i \setminus \tilde{s}_t) \cdot \min(1, \#(s_i \cap \tilde{s}_t))]}{\#(\tilde{s}_{t-1}) + \#(\tilde{s}_{t+1})}. \tag{14}$$

Essentially, $\min(1, \#(s_i \cap \tilde{s}_t))$ is a binary indicator whether scene $s_i$ shares at least one shot with $\tilde{s}_t$, and $\#(s_i \setminus \tilde{s}_t)$ are the shots of these scenes which are not part of $\tilde{s}_t$. Therefore this measures how much the overlapping proposed scenes extend beyond the ground truth scene normalized by the number of shots in the neighboring scenes.

These measures for each ground truth scene are aggregated into video-wide metrics as the weighted average:

$$C = \sum_{t=1}^{n} C_t \cdot \frac{\#(\tilde{s}_t)}{\sum_i \#(\tilde{s}_i)}, \qquad O = \sum_{t=1}^{n} O_t \cdot \frac{\#(\tilde{s}_t)}{\sum_i \#(\tilde{s}_i)}. \tag{15}$$

Finally, as a single score for the quality of the scene detection, we compute the harmonic mean:

$$F = 2 \cdot \frac{C \cdot (1 - O)}{C + (1 - O)}. \tag{16}$$

## E ADDITIONAL $D$ EXAMPLES

In Table 4 we present various stages of $D$ from visual features of the video La Chute D'une Plume from OVSD with the accompanied ground truth $D^*$, and in Table 5 the same for the video Big Buck Bunny. In Tables 6 and 7 we show how the $D$ matrices and gradients evolve over a number of epochs for the videos La Chute D'une Plume and Big Buck Bunny respectively.

In general our observations are that OSG-Triplet manages to emphasize the small scenes better than large scenes, while OSG-Block is the reverse. OSG-Block-Adjacent gives a good trade-off of emphasizing the immediate off-diagonal, but results in some low distances in the far off-diagonal. In practice, these shouldn't affect the OSG algorithm if the intervening distances are large enough. OSG-Prob converges more slowly, learns from the boundary edges, and gives a good trade-off as well.

Regarding this last point, part of our motivation for OSG-Prob is to have a configuration which is specifically reliant on division locations as opposed to the block-diagonal. Such a structure would allow OSG to be integrated into a larger learning pipeline. For example, there are other temporal analysis tasks where division is only a part of the process. In the weakly-supervised regime there might not be ground truth divisions with which to perform OSG-Triplet or OSG-Block. OSG-Prob on the other hand, could be configured to perform backpropagation on a loss which reflects on the locations of division, and is inferred back to the distance values. In this respect, our continued research involves having this component as a plug-and-play module for other tasks which can act as temporal region proposal networks (see Figure 9).

## F    ADDITIONAL $T$ EXAMPLES

In Figure 10 and 11 we show the progression of the values of $T(i)$ over a number of iterations for videos La Chute D'une Plume and Big Buck Bunny respectively. We can see that as the iterations progress, $T$ raises the probability at the ground truth points of division. The probability is lowered for locations with no true division even though this is not specifically enforced by the loss but rather an outcome of the construction of $T$.

Specifically we note that in these instances the small beginning scenes proved difficult for the network to emphasize $T$ on. We speculate that this is due to the formulation, where smaller values of $n$ inspect longer and longer sequences (see the formulation in the paper). Possible future work could be to formulate an additional mirrored OSG which inspects the $D$ matrix backwards.

## G    ADDITIONAL VISUAL RESULTS

Figures 12, 13, and 14, show results on sections of videos from the OVSD dataset. In general, we can see divisions which result in reasonable and often precise scene divisions. Using these divisions for applying video understanding and classification technologies will undoubtedly be superior over applying them on the entire video or on naive uniform divisions. Specifically, Figure 14 is a single scene from the video Tears of Steel which includes intricate character and setting changes. This portrays the complexity of the task and the challenges that the method needs to overcome. Despite the fact that all of the proposed scene divisions in this case are technically false, we note that they present a plausible division to story-units, and can be useful for a variety of downstream tasks.

Figure 9: OSG-Prob as a plug-and-play temporal region proposal network.
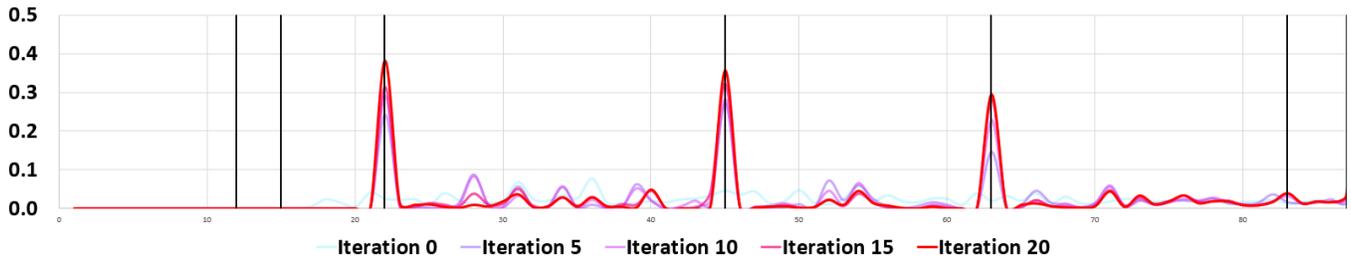


**Figure 10: Progression of $T$ as a function of $i$ (shot number) over a number of iterations for video La Chute D'une Plume. Graphs go from translucent blue to opaque red as iterations progress (best viewed in color). Vertical black lines indicate ground truth divisions.**
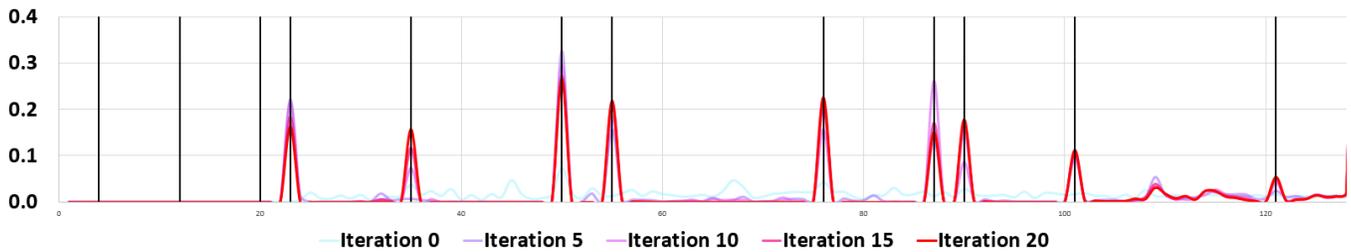


**Figure 11: Progression of $T$ as a function of $i$ (shot number) over a number of iterations for video Big Buck Bunny. Graphs go from translucent blue to opaque red as iterations progress (best viewed in color). Vertical black lines indicate ground truth divisions.**

### Table 3: OVSD dataset details

| Video Name | Short Name | Duration (minutes) | # Scenes | # Shots | Genre |
|---|---|---|---|---|---|
| 1000 Days | 1000 | 43 | 23 | 404 | Drama |
| Big Buck Bunny | BBB | 8 | 13 | 129 | Animation |
| Boy Who Never Slept | BWNS | 69 | 23 | 336 | Comedy, Romance |
| CH7 | CH7 | 86 | 45 | 1293 | Crime |
| Cosmos Laundromat | CL | 10 | 6 | 94 | Animation |
| Elephants Dream | ED | 9 | 8 | 128 | Animation |
| Fires Beneath Water | FBW | 76 | 63 | 411 | Documentary |
| Honey | Honey | 86 | 21 | 326 | Drama |
| Jathia's Wager | JW | 21 | 16 | 177 | Drama, Sci-Fi |
| La Chute D'une Plume | LCDP | 10 | 11 | 88 | Animation |
| Lord Meia | LM | 37 | 28 | 333 | Crime, Comedy |
| Meridian | Meridian | 12 | 10 | 64 | Mystery, Sci-Fi |
| Oceania | Oceania | 54 | 32 | 253 | Drama, Mystery |
| Pentagon | Pentagon | 50 | 32 | 305 | Comedy, Drama |
| Route 66 | Route 66 | 103 | 56 | 1357 | Documentary |
| Seven Dead Men | SDM | 57 | 35 | 167 | Crime |
| Sintel | Sintel | 12 | 7 | 198 | Animation |
| Sita Sings the Blues | SStB | 81 | 53 | 1384 | Animation, Comedy |
| Star Wreck | SW | 103 | 56 | 1439 | Comedy, Sci-Fi |
| Tears of Steal | ToS | 10 | 6 | 136 | Drama, Sci-Fi |
| Valkaama | Valkaama | 93 | 49 | 714 | Drama |

Table 4: An example $D$ from the video La Chute D'une Plume from OVSD. On the left: Ground Truth ($D^*$), Orig (without an applied embedding), Epoch 0 (embedding before learning). On the right, trained examples after 20 epochs for: OSG-Triplet, OSG-Block, OSG-Block-Adjacent, and OSG-Prob, with corresponding gradients (bottom row)
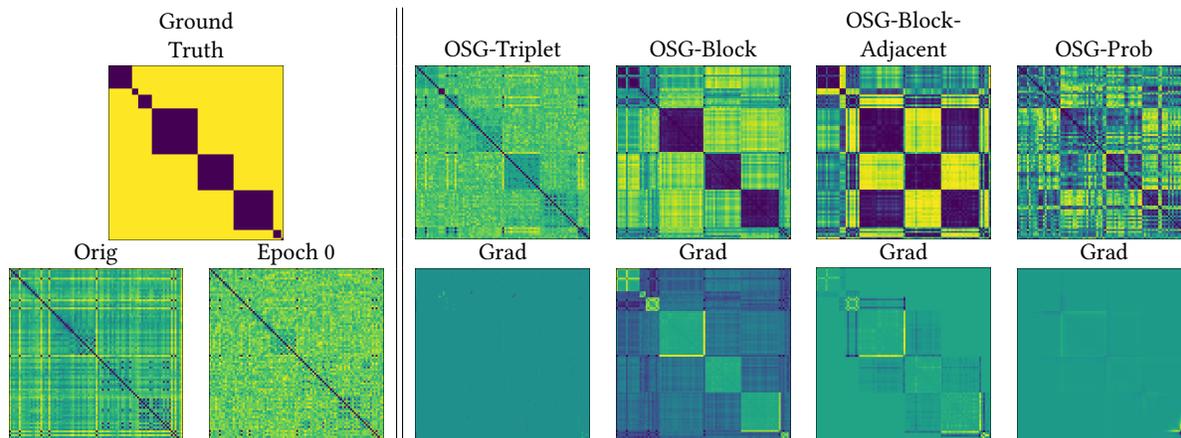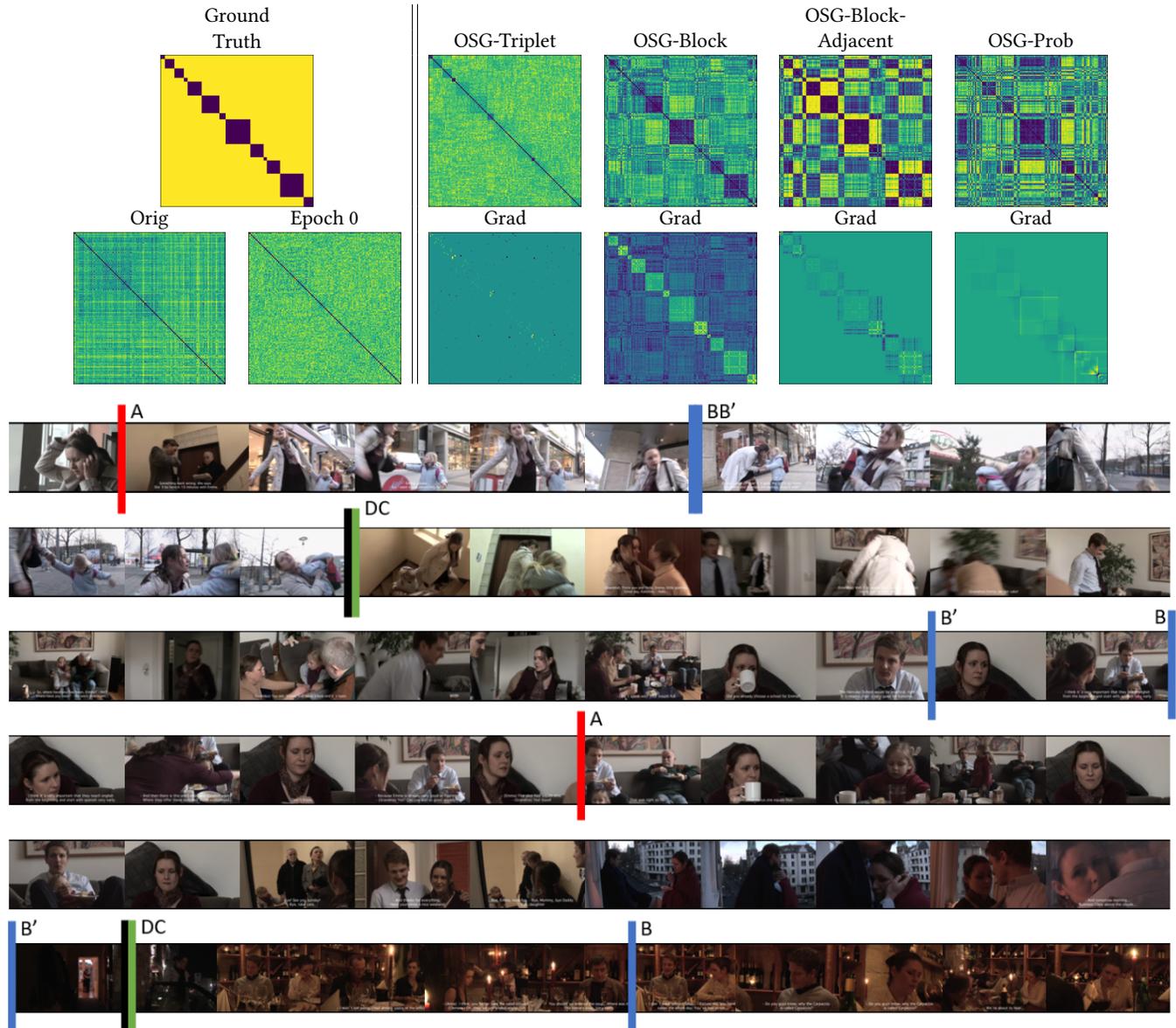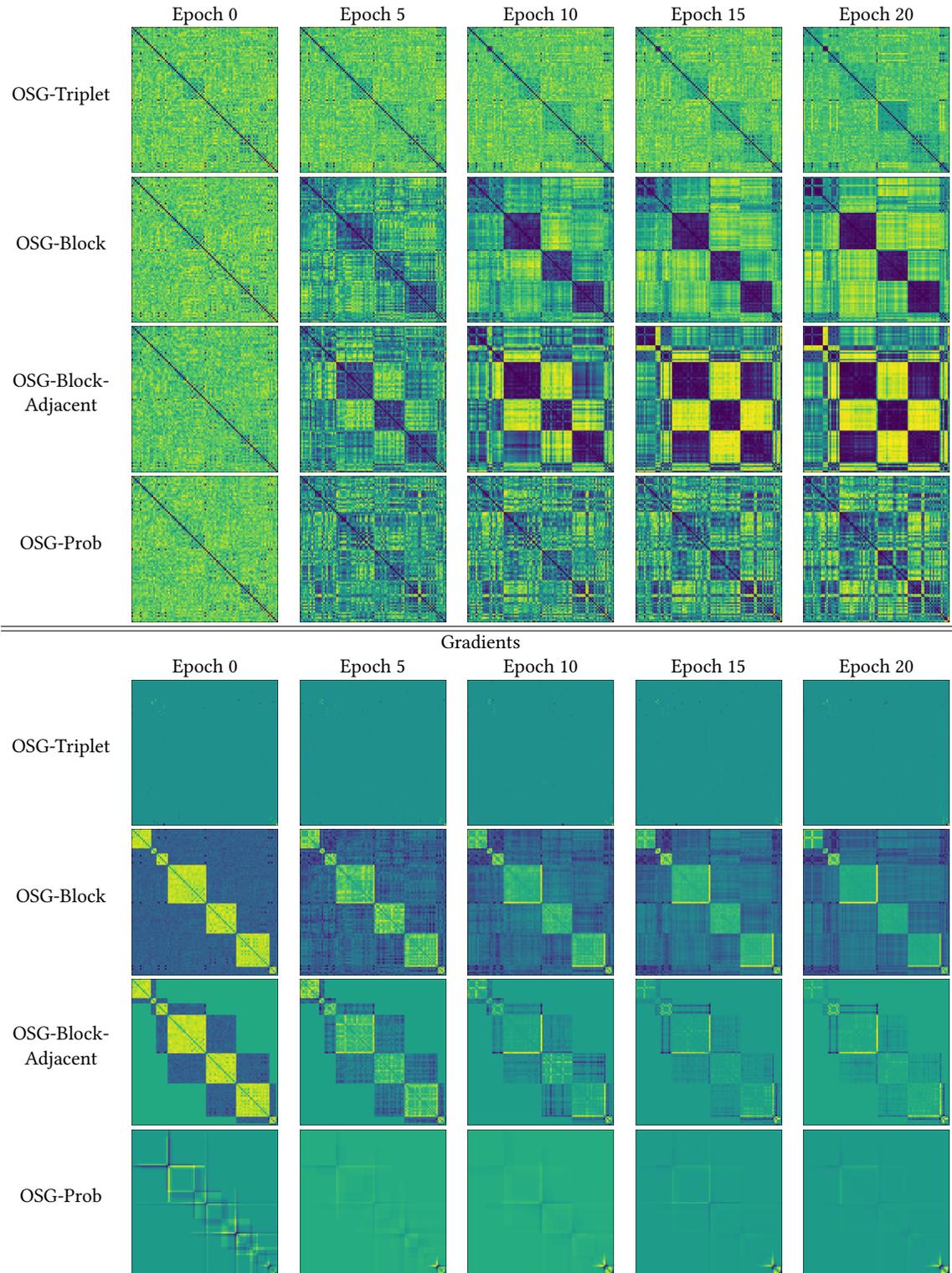
**Table 5: An example $D$ from the video Big Buck Bunny from OVSD. On the left: Ground Truth ($D^*$), Orig (without an applied embedding), Epoch 0 (embedding before learning). On the right, trained examples after 20 epochs for: OSG-Triplet, OSG-Block, OSG-Block-Adjacent, and OSG-Prob, with corresponding gradients (bottom row)**
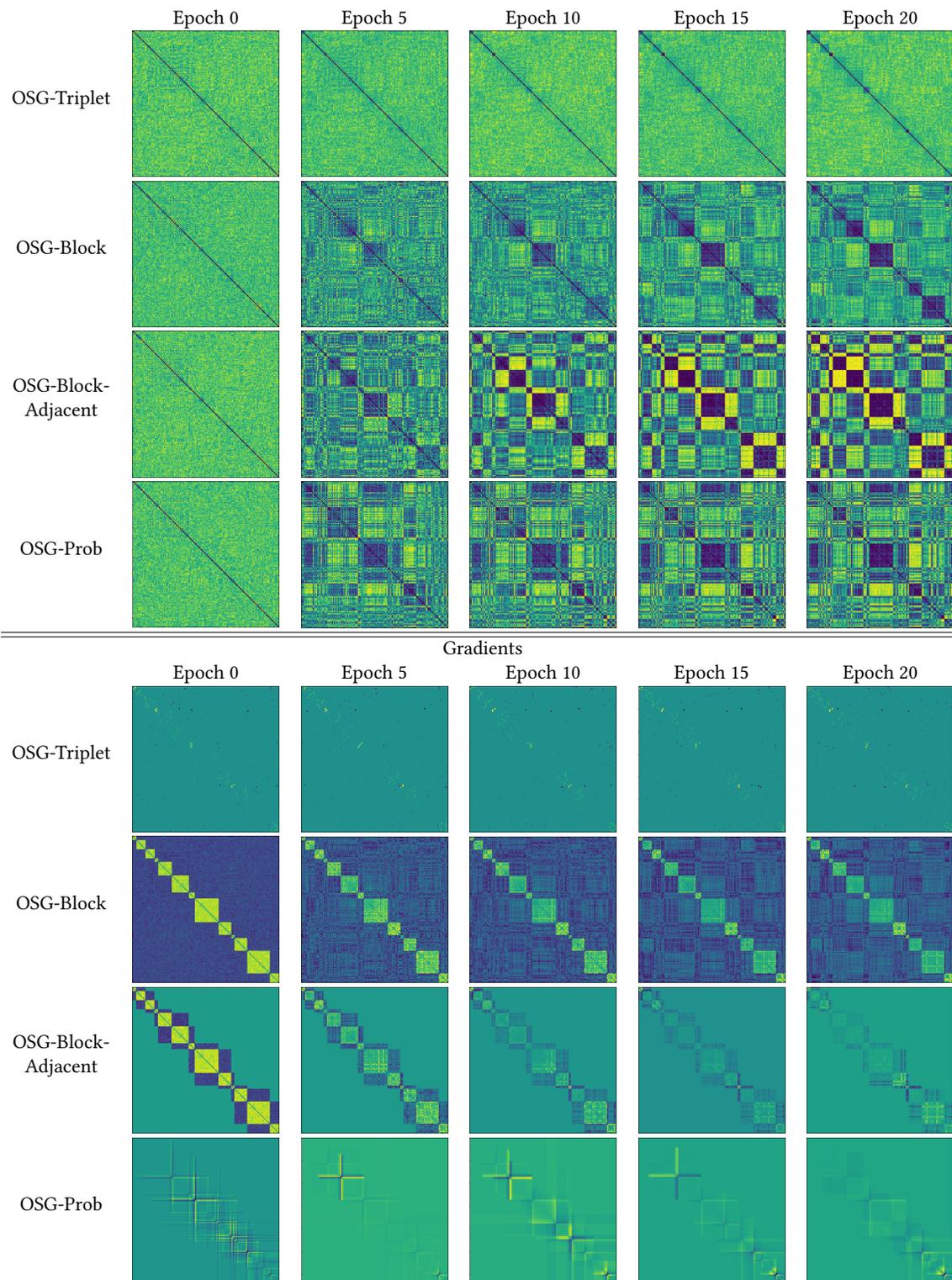


**Figure 12: Qualitative results of configurations on shots 68 through 128 of the video 1000 Days from the OVSD dataset. Points of division marked by A. OSG-Triplet (red) B. OSG-Block (blue) B'. OSG-Block-Adjacent (blue) C. OSG-Prob (green) and D. Ground truth (black).**

**Table 6:** $D$ and gradients from the video La Chute D'une Plume from OVSD evolving over a number of Epochs

**Table 7:** $D$ and gradients from the video Big Buck Bunny from OVSD evolving over a number of Epochs
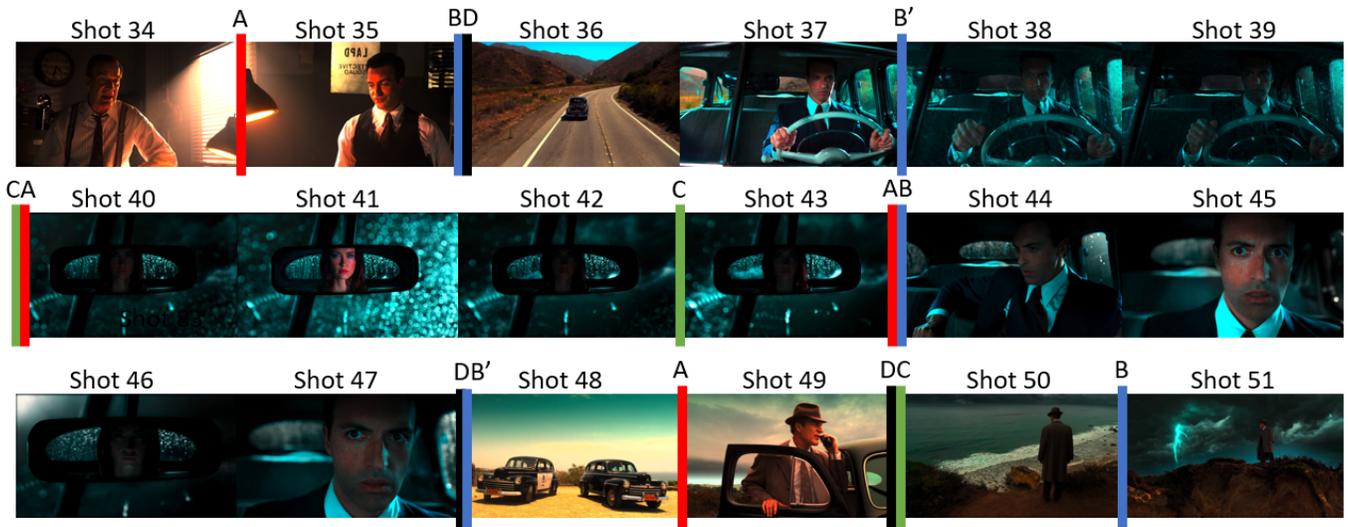
Figure 13: Qualitative results of configurations on a section of the video Meridian from the OVSD dataset. Points of division marked by A. OSG-Triplet (red) B. OSG-Block (blue) B'. OSG-Block-Adjacent (blue) C. OSG-Prob (green) and D. Ground truth (black).



Figure 14: Qualitative results of configurations on a section of the video Tears of Steel from the OVSD dataset. Points of division marked by A. OSG-Triplet (red) B. OSG-Block (blue) B'. OSG-Block-Adjacent (blue) C. OSG-Prob (green). The shots are part of a single complex ground truth scene.