

Modality-Aware Shot Relating and Comparing for Video Scene Detection

Jiawei Tan^{1,2}, Hongxing Wang^{1,2*}, Kang Dang³, Jiaxin Li^{1,2}, Zhilong Ou^{1,2}

¹Key Laboratory of Dependable Service Computing in Cyber Physical Society (Chongqing University), Ministry of Education, China

²School of Big Data and Software Engineering, Chongqing University, China

³School of AI and Advanced Computing, XJTLU Entrepreneur College (Taicang), Xi'an Jiaotong-Liverpool University, Suzhou, China

{jwta, ihwang}@cqu.edu.cn, Kang.Dang@xjtlu.edu.cn, jiaxin.li@cqu.edu.cn, zlou@stu.edu.cn

Abstract

Video scene detection involves assessing whether each shot and its surroundings belong to the same scene. Achieving this requires meticulously correlating multi-modal cues, *e.g.* visual entity and place modalities, among shots and comparing semantic changes around each shot. However, most methods treat multi-modal semantics equally and do not examine contextual differences between the two sides of a shot, leading to sub-optimal detection performance. In this paper, we propose the Modality-Aware Shot Relating and Comparing approach (MASRC), which enables relating shots per their own characteristics of visual entity and place modalities, as well as comparing multi-shots similarities to have scene changes explicitly encoded. Specifically, to fully harness the potential of visual entity and place modalities in modeling shot relations, we mine long-term shot correlations from entity semantics while simultaneously revealing short-term shot relations from place semantics. In this way, we can learn distinctive shot features that consolidate coherence within scenes and amplify distinguishability across scenes. Once equipped with distinctive shot features, we further encode the relations between preceding and succeeding shots of each target shot by similarity convolution, aiding in the identification of scene ending shots. We validate the broad applicability of the proposed components in MASRC. Extensive experimental results on public benchmark datasets demonstrate that the proposed MASRC significantly advances video scene detection.

Code — <https://github.com/ExMorgan-Alter/MASRC>

1 Introduction

Video scene detection involves identifying whether a shot is at the scene ending or not. Upon detection, the video is segmented into coherent sets of scenes. Temporal scenes provide structured information that serves as the foundation for downstream applications such as text-to-video retrieval (Bain et al. 2020) and human-centric storyline construction (Vicol et al. 2018).

Before determining whether a shot is an ending shot, one needs to examine the surrounding shots to grasp the context of the target shot from different visual modalities, especially in terms of imaging entities and places, in the video.

*Corresponding author: Hongxing Wang.
Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

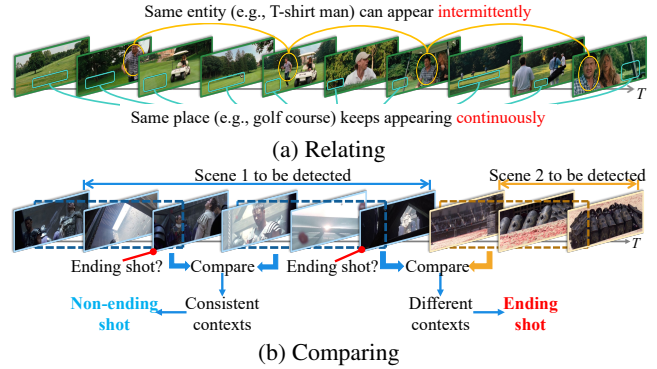


Figure 1: Intuition behind our proposed modality-aware shot relating and comparing. (a) Relating: We capture long-term shot relations by linking intermittently appearing identical entities and short-term shot relations by relating consecutive shots depicting the same place. (b) Comparing: Consistency drawn from comparisons between fore-and-aft contexts is an indicator to classify whether a target shot is an ending shot.

This process requires associating shared semantics between shots for each modality to model shot relations. As shown in Fig. 1a, we need to intermittently link shots featuring the same entity for capturing long-term shot relations, while relating consecutive shots depicting the same place for the short-term shot relations. Common methods (2020; 2023; 2024b) handle different modalities in indiscriminate manner for shot relations modeling. However, this modality-agnostic strategy fails to capture the diverse temporal relations exhibited by different modalities. To address this issue, we shift temporal relation modeling from modality-agnostic recipes to modality-aware scheme. As presented in Fig. 2, we design Modality-Aware Shot Relating (MASR) to enable capturing long-range and short-range temporal relations between shots respectively in entity and place modalities. In the entity modality, we construct an Entity Jumping Graph (EJG) to correlate shots with similar entities by which a long-term dependencies can be caught for similar but distant shots. With the place modality, we design a Place Continuity Graph (PCG) to connect time-continuous shots depicting the same place, modeling short-term dependencies between shots.

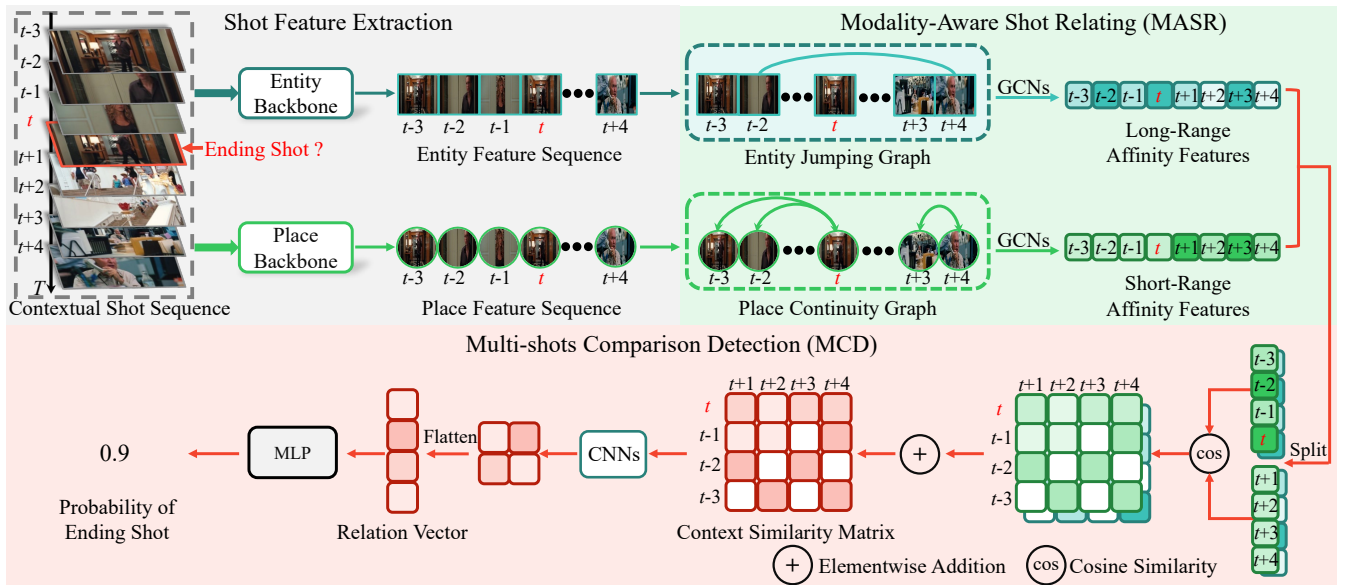


Figure 2: Diagram of how the proposed MASRC determines whether a target shot (in the red border) is an ending shot or not. We build an entity jumping graph and a place continuity graph for GCN message passing to separately embed long and short shot relations into shot representations. Relying on similarity comparison between the fore-and-aft shots of each target shot as well as further similarity change encoding by convolution, the probability of target shots being ending shots can be better predicted by a simple MLP classifier.

Once modeling shot relations through MASR, we are allowed to classify whether each shot is an ending shot via its own context-embedded shot feature, as previous efforts (Wu et al. 2022; Mun et al. 2022; Islam et al. 2023; Chen et al. 2023; Yang et al. 2023) have done. Unfortunately, a single context-embedded shot feature struggles to capture complex context changes, leading to numerous false positives. This is because these methods overlook the fact that an ending shot signifies a contextual difference between its two sides. As illustrated in Fig. 1b, it is crucial to compare the consistency between fore-and-aft contexts using context-embedded shot features. Fig. 2 shows that devised Multi-shots Comparison Detection (MCD), which compares multi-shot similarities and then encodes similarity changes by convolution, aiding the identification of ending shots.

Across the board, we propose a novel architecture for video scene detection, dubbed **Modality-Aware Shot Relating and Comparing (MASRC)**. We pursue MASR to capture diverse shot relations for entity and place modalities. Next, we design MCD to encode contextual different between shots to support the classification of ending shots.

In a nutshell, our contributions include:

- We propose MASR to capture long-range and short-range shot relations by considering the distinctive roles of entity and place modalities in video scene detection.
- Instead of relying solely on individual shots for video scene detection, we compare shot semantics on both sides of the target shot and propose MCD to better encode and detect scene transitions.
- We perform comprehensive evaluations on three video scene detection datasets including MovieNet (Huang

et al. 2020), BBC (Baraldi, Grana, and Cucchiara 2015), and OVSD (Rotman, Porat, and Ashour 2016). Our proposed method surpasses previous approaches by large margins in various learning settings.

2 Related Work

Video Scene Detection. Early methods directly employ shot features to cluster adjacent shots into scenes in an unsupervised manner. For instance, (Odobez, Gatica-Perez, and Guillemot 2003) clusters shots based color histograms and temporal closeness by a spectral clustering method. Similarly, (Ngo, Pong, and Zhang 2002) employs color histogram intersection for initial shot clustering, followed by merging the preliminary results using a sliding time window to derive scenes. (Rasheed and Shah 2005) utilizes color and motion information between shots to generate a similarity graph, partitioning the graph to identify video scenes. However, these methods achieve limited performance due to their manually modeling relations between shots, resulting in insufficient discriminability of shot features from different scenes. Recently, supervised methods focus on designing adaptive and flexible relating shot-correlation mechanisms. Some methods exploit single-modality visual shot semantics to explore shot relations. For instances, (Tan et al. 2023) and (Mun et al. 2022) employ a multi-head attention mechanism to emphasize long-term shot relationships. (Yang et al. 2023) additionally incorporates local window attention to capture short-term shot relationships. (Tan, Wang, and Yuan 2024) reweights each shot feature to enhance the ability of LSTM (Graves and Schmidhuber 2005) to correlate shots within long video scenes. Due to the limitations of

single modality in semantic representation, other methods employ multi-modality visual shot semantics, such as entity, place modalities. (Rao et al. 2020) designs multi-modal shot boundary features for capturing both difference and relations between neighborhood shots. (Wei et al. 2023) employs a multi-head self-attention mechanism on the similarity matrix generated by multi-modal shot features for exploring high-order relationships between shots. However, these methods tend to design a unified framework to handle different modalities, struggling to capture the diverse temporal relations exhibited by different modalities. In contrast, we propose MASR to capture long-range and short-range shot relations by considering the distinctive roles of entity and place modalities in video scene detection.

On the other hand, when detecting ending shots, some methods (Wu et al. 2022; Mun et al. 2022; Islam et al. 2023; Chen et al. 2023) make a prediction by a single context-embedded shot representation. Other methods (Wei et al. 2023; Tan et al. 2023) design learnable boundary class vectors and use it to query each shot feature to detect whether each shot is an ending shot. However, these methods overlook the fact that an ending shot signifies a semantic difference between its two sides and do not explicitly compare the shot contexts for detection. In this paper, we compare the shot semantics on both sides of the target shot to identify the ending shots.

Graph Neural Network. In recent years, an increasing number of non-Euclidean data have been represented as graphs, posing significant challenges to existing neural network methods. To effectively deal with such data, Graph Neural Networks (GNNs) have attracted much attention. GNNs have been applied in various fields, including recommendation systems (He et al. 2020), computer vision (Yan, Xiong, and Lin 2018), and natural language processing (Yao, Mao, and Luo 2019). Graph convolutional networks, a type of GNN, can be divided into spectral-based and spatial-based approaches (Wu et al. 2021). Spectral-based approaches implement graph convolution by defining filters similar to graph signal processing, while spatial-based approaches define graph convolution by information propagation and have gained momentum for their efficiency, flexibility, and generality. Widely used GNN techniques include Graph-SAGE (Hamilton, Ying, and Leskovec 2017), GAT (Velickovic et al. 2018), and GCN (Kipf and Welling 2017).

3 Methodology

3.1 Problem Formulation

Given a video divided into shots, video scene detection aims to learn a classifier f , which holds $f(s) = 1$ if shot s ends a scene, $f(s) = 0$ otherwise. Since a shot itself cannot form the concept of the end or non-end of a scene, we have to put a shot in its temporal context towards this end. To be explicit, we rewrite $f(s)$ as $f(c(s))$, where $c(\cdot)$ truncates the contextual sequence of shots including s . In this study, we propose a Modality-Aware Shot Relating and Comparing solution (MASRC). As illustrated in Fig. 2, it allows Modality-Aware Shot Relating (MASR) for Multi-shot Comparison

Detection (MCD) to model $f(c(s))$.

3.2 Modality-Aware Shot Relating

To perform prediction on shot s_t at time t , we cut its context $c(s_t)$ as a time sliding window of length T centered around s_t . Since the frames within a shot are similar, we represent each shot with a randomly selected frame, as is common practice (Mun et al. 2022; Yang et al. 2023; Tan et al. 2024a). Without loss of generality, we let T be even, and denote $c(s_t)$ by $\{s_{t-T/2+1}, \dots, s_t, \dots, s_{t+T/2}\}$. Given the importance of actors, objects, and places in composing visual-centric shot semantics (Rao et al. 2020), we directly employ different pre-trained ResNet-50 models (He et al. 2016) to extract diverse visual semantics, as done in (Tan et al. 2024b). One model, pre-trained on the ImageNet dataset (Russakovsky et al. 2015), extracts visual entity features related to actors and general objects, *i.e.*, $\mathbf{X}^E = \{\mathbf{x}_{t-T/2+1}^E, \dots, \mathbf{x}_{t+T/2}^E\}$. Another ResNet-50 model, pre-trained on the Places dataset (Zhou et al. 2018), is used to obtain place features $\mathbf{X}^P = \{\mathbf{x}_{t-T/2+1}^P, \dots, \mathbf{x}_{t+T/2}^P\}$ for shot sequence $c(s_t)$.

Entity-based Long-Term Dependency. The reappearance of the same actors or objects in a shot after several shots reflects that entities carry long-range shot relations. To relate shots with similar entity semantics, we construct an entity jumping graph (EJG) $G^E = \langle \mathbf{X}^E, \mathbf{E}^E \rangle$. The shots with features \mathbf{X}^E act as nodes and edges $\mathbf{E}^E = \{E_{ij}^E\}$ between nodes are determined by the semantic similarity between shots, given by:

$$E_{ij}^E = \begin{cases} S_{ij}^E, & S_{ij}^E \in \text{top-}k(S_i^E), \\ 0, & \text{otherwise,} \end{cases} \quad (1)$$

where cosine similarity $S_{ij}^E = \cos(\mathbf{x}_i^E, \mathbf{x}_j^E)$ is computed based on entity features \mathbf{x}_i^E and \mathbf{x}_j^E . EJG establishes connections between each shot and its k most similar shots in the entity feature space, enabling each shot to be linked with distant but similar shots in the temporal dimension.

After that, we perform message passing for the shot nodes on G^E with shot features \mathbf{X}^E . In this process, the shot nodes will receive messages from its connected shots in G^E , facilitating information exchange between distant shots. Specifically, we stack two graph convolution network (GCN) (Kipf and Welling 2017) layers on EJG to model long-term shot relations for obtaining long-range affinity features \mathbf{X}^{LR} :

$$\begin{cases} \mathbf{X}_0^{\text{LR}} = \mathbf{X}^E, \\ \mathbf{X}_i^{\text{LR}} = \text{LN}(\mathbf{X}_{i-1}^{\text{LR}} + \sigma(\text{GCN}_i(\mathbf{X}_{i-1}^{\text{LR}}, \tilde{\mathbf{E}}^E))), \end{cases} \quad (2)$$

where, in each layer $i \in \{1, 2\}$, we apply vanilla GCN smoothing to current features $\mathbf{X}_{i-1}^{\text{LR}}$ on the edge affinities $\tilde{\mathbf{E}}^E$ normalized from \mathbf{E}^E (Kipf and Welling 2017), followed by activation mapping with $\sigma(\cdot)$. The resultant activated features will be added to original features $\mathbf{X}_{i-1}^{\text{LR}}$ as smoothly rectified features for layer normalization (LN) (Ba, Kiros, and Hinton 2016).

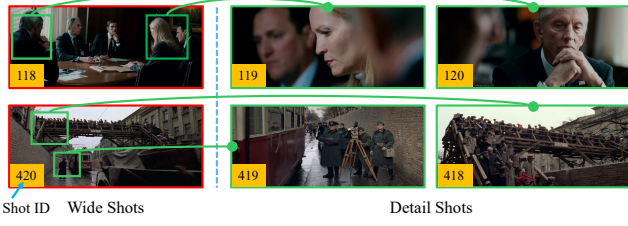


Figure 3: Each wide shot depicts an overview of a place, and its detail shots zoom in on specific details within the same place.

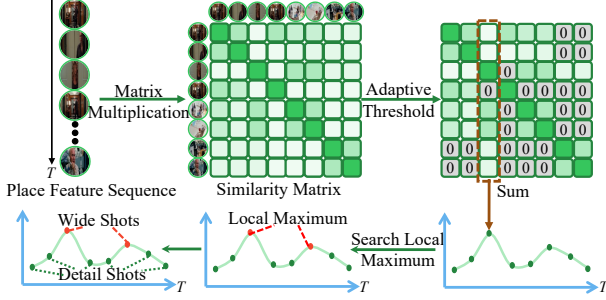


Figure 4: Identification of wide shots and detail shots.

Place-based Short-Term Dependency. Consecutive shots in the same scene often share the same place. Establishing correlations among these shots allows us to effectively model short-range shot relations. However, the challenge lies in grouping shots that depict the same place. As shown in Fig. 3, some shots provide an overview of the layout for place, while others focus on on specific details within the same place. For clarity, we refer shots portraying the overall layout of the place to “wide shots” and the remaining shots to “detail shots”. That is to say, wide shot can have similarities to its detail shots, but detail shots are usually very diverse. Hence, a wide shot has more similar shots than its detail shots. With this observation, we count the number of similar shots to each involved shot to identify wide and detail shots. As shown in Fig. 4, we leverage an adaptive threshold to obtain the number of similar shots to shot $s_i \in c(s_t)$ as follow,

$$n_i = \sum_{j \in [t - \frac{T}{2} + 1, t + \frac{T}{2}]} \mathbb{I}(S_{ij}^P > \bar{S}^P), \quad (3)$$

where $S_{ij}^P = \cos(\mathbf{x}_i^P, \mathbf{x}_j^P)$ denotes cosine similarity between features from the place modality, \bar{S}^P denotes the averaged similarity among shot features of place modality, and $\mathbb{I}(\cdot)$ refers to an indicator function that takes 1 for non-zero inputs. According to $\{n_i\}_{i \in [t - T/2 + 1, t + T/2]}$, we select those shots $\{s_W\}$ with $W = \{i \in [t - T/2 + 1, t + T/2] | n_i > n_{i-1} \cap n_i > n_{i+1}\}$ as wide shots, and the rest shots $\{s_D\}_{D=[t - T/2 + 1, t + T/2] \setminus W}$ as detail shots.

We then use the cosine similarity between detail shots and wide shots as well as the temporal proximity between them to obtain detail-wide shots affiliation between $i \in \mathcal{D}, j \in \mathcal{W}$

as follow,

$$j_i^* = \operatorname{argmax}_j (S_{ij}^P + D_{ij}), \quad (4)$$

where $D_{ij} = \frac{1}{|i-j|}$ measures the temporal proximity between shots s_i and s_j . The motivation for D_{ij} is to prevent linking detail shots with wide shots that are far apart, thus ensuring that the associated wide and detail shots describe the same place. It is worth noting that due to the inability to access the focal length of the shots, the wide shots we found by the similarity between shots are sometimes different from the wide shots in reality.

With the detail-wide shots affiliation established, we proceed to capture the short-term shot dependency based on the place features \mathbf{X}^P . Considering the hierarchical structure of the space depicted by the wide shot and its corresponding detail shot, we devise a two-stage graph reasoning process.

On the one hand, each detail shot presents the spatial details of its corresponding wide shots. For message passing from detail shots to wide shots, we build detail-to-wide graph $G^{\mathcal{D}2\mathcal{W}} = \langle \mathbf{X}^P, \mathbf{E}^{\mathcal{D}2\mathcal{W}} \rangle$. The edge weight $\mathbf{E}^{\mathcal{D}2\mathcal{W}} = \{E_{ij}^{\mathcal{D}2\mathcal{W}}\}$ is formulated as:

$$E_{ji}^{\mathcal{D}2\mathcal{W}} = \begin{cases} (\mathbf{W}_1^{\mathcal{D}2\mathcal{W}} \mathbf{x}_j^P)^T \mathbf{W}_2^{\mathcal{D}2\mathcal{W}} \mathbf{x}_i^P, & j = j_i^*, \\ -\infty, & \text{otherwise,} \end{cases} \quad (5)$$

where $\mathbf{W}_1^{\mathcal{D}2\mathcal{W}}$ and $\mathbf{W}_2^{\mathcal{D}2\mathcal{W}}$ are learnable matrices.

Next, we perform message passing for the shot nodes on the $G^{\mathcal{D}2\mathcal{W}}$, where wide-shot nodes integrates information from its detail-shot nodes. We apply a one-layer GCN message passing and inference to obtain short-range detail to wide shot features $\mathbf{X}^{\mathcal{D}2\mathcal{W}}$:

$$\mathbf{X}^{\mathcal{D}2\mathcal{W}} = \text{LN}(\mathbf{X}^P + \sigma(\text{GCN}(\mathbf{X}^P, \text{softmax}(\mathbf{E}^{\mathcal{D}2\mathcal{W}})))), \quad (6)$$

where we employ softmax for normalization to make the edge weights learnable (Velickovic et al. 2018).

On the other hand, we expect the updated information from wide shots can be conveyed back to the detail shots. Hence, based on short-range detail to wide shot features $\mathbf{X}^{\mathcal{D}2\mathcal{W}}$, we build the wide-to-detail graph $G^{\mathcal{W}2\mathcal{D}} = \langle \mathbf{X}^{\mathcal{D}2\mathcal{W}}, \mathbf{E}^{\mathcal{W}2\mathcal{D}} \rangle$. The edge weight $\mathbf{E}^{\mathcal{W}2\mathcal{D}} = \{E_{ij}^{\mathcal{W}2\mathcal{D}}\}$ is formulated as:

$$E_{ij}^{\mathcal{W}2\mathcal{D}} = \begin{cases} (\mathbf{W}_1^{\mathcal{W}2\mathcal{D}} \mathbf{x}_i^{\mathcal{D}2\mathcal{W}})^T \mathbf{W}_2^{\mathcal{W}2\mathcal{D}} \mathbf{x}_j^{\mathcal{D}2\mathcal{W}}, & j = j_i^*, \\ -\infty, & \text{otherwise.} \end{cases} \quad (7)$$

By learning $\mathbf{W}_1^{\mathcal{W}2\mathcal{D}}$ and $\mathbf{W}_2^{\mathcal{W}2\mathcal{D}}$ on $G^{\mathcal{W}2\mathcal{D}}$ with a one-layer GCN message passing and inference, we have the ultimate short-range affinity features \mathbf{X}^{SR} :

$$\mathbf{X}^{\text{SR}} = \text{LN}(\mathbf{X}^{\mathcal{D}2\mathcal{W}} + \sigma(\text{GCN}(\mathbf{X}^{\mathcal{D}2\mathcal{W}}, \text{softmax}(\mathbf{E}^{\mathcal{W}2\mathcal{D}})))). \quad (8)$$

3.3 Multi-shots Comparison Detection

While each shot enriched with diverse temporal-scale information becomes more discriminative, detecting the ending shots remains a challenge based on each shot via its own context-embedded shot feature. This is because the ending shot implies a kind of semantic difference between its two

sides. It is essential to compare the shot semantics on both sides of the target shot to make predictions. In pursuit of this objective, we propose Multi-shots Comparison Detection (MCD). For shot s_t , we compare shot semantics on its two sides, having context similarity matrix M ,

$$M_{ij} = \cos(\mathbf{x}_i^{\text{LR}}, \mathbf{x}_j^{\text{LR}}) + \cos(\mathbf{x}_i^{\text{SR}}, \mathbf{x}_j^{\text{SR}}) \quad (9)$$

where $i \in [t - T/2 + 1, t]$, $j \in [t + 1, t + T/2]$.

To capture intricate patterns of semantic variation, we employ a CNN-based network on M , producing the relation vector r ,

$$r = \text{Flatten}(\text{CNNs}(M)), \quad (10)$$

where $\text{Flatten}(\cdot)$ is to flatten its input matrix into a one-dimension vector, and $\text{CNNs}(\cdot)$ denotes a 4-layer VGG (Simonyan and Zisserman 2015) Network.

Finally, we can utilize a fully connected MLP and a sigmoid layer on the flatten R to predict the probability \hat{y}_t of the shot s_t being an ending shot.

3.4 Training and Objective Functions

We employ two loss functions in our model training: the self-supervised loss and the supervised loss. The details of these two loss functions are as follows.

Self-supervised loss. Same as previous methods (Mun et al. 2022; Islam et al. 2023), the self-supervised loss aligns predictions with pseudo-scene boundaries using binary cross-entropy loss:

$$L_f = -y_t^b \log(\hat{y}_t) + (1 - y_t^b) \log(1 - \hat{y}_t), \quad (11)$$

where y_t^b denotes the pseudo label of shot s_t generated by the Modified Dynamic Warping algorithm (Mun et al. 2022).

Supervised loss. It is an ending shot prediction loss in form of the binary cross-entropy loss:

$$L_f = -y_t \log(\hat{y}_t) + (1 - y_t) \log(1 - \hat{y}_t), \quad (12)$$

where $y_t \in \{0, 1\}$ denotes the ground-truth binary label of shot t .

Combining the aforementioned losses, our MASRC supports various learning approaches: **self-supervised learning**, **fully supervised learning**, and **self-supervised transfer learning**. In self-supervised learning, we employ the self-supervised loss for model training. In fully supervised learning, we utilize a supervised loss. Self-supervised transfer learning involves two phases, where the self-supervised loss is applied for pre-training, followed by the use of the supervised loss for model fine-tuning.

4 Experiments

4.1 Settings

Datasets: We assess the performance of our method on three widely used video scene detection datasets, *i.e.*, MovieNet (Huang et al. 2020), BBC (Baraldi, Grana, and Cucchiara 2015), and OVSD (Rotman, Porat, and Ashour 2016).

MovieNet. It is a vast dataset with 1,100 movies and 1.6 million shots. 318 movies are annotated with scene boundaries, forming the MovieScenes dataset (Rao et al. 2020) for video scene detection. MovieScenes is further divided into subsets of 190 movies for training, 64 for validation, and 64 for testing. For different learning ways, we remain consistent with the settings of the previous methods (Wu et al. 2022; Yang et al. 2023) and always evaluate our model on the test split of MovieScenes. In the self-supervised scenario, we utilize the 660 unlabeled videos from MovieNet for pre-training. In the supervised setting, we utilize 190 training videos from MovieScenes for training. For self-supervised transfer learning, we utilize the 660 unlabeled videos from MovieNet for pre-training, and 190 training videos from MovieScenes for fine-tuning.

OVSD. It consists of 21 short films, each lasting approximately 30 minutes. It contains a total of 10,000 shots and 300 scenes, extracted from movie scripts. Due to lacking predefined splits, we follow prior studies (Wu et al. 2022; Mun et al. 2022; Islam et al. 2023), training our model using the MovieNet dataset and then assessing its performance on OVSD without additional fine-tuning.

BBC. It comprises of 11 episodes from the BBC educational TV series *Planet Earth* (BBC 2006). These videos have an average duration of 50 minutes and include a total of 670 scenes and 4.8K shots. As with our evaluation on the OVSD dataset, we train our model using the MovieNet dataset and assess its performance on the BBC dataset without additional fine-tuning, following established research practices (Wu et al. 2022; Mun et al. 2022; Islam et al. 2023).

Metrics: To measure the performances, we use the same evaluation metrics used in prior methods (Mun et al. 2022; Wu et al. 2022; Islam et al. 2023), which include the Average Precision (AP), the mean Intersection over Union (mIoU), and the F1-score (F1). These metrics serve to evaluate the effectiveness of video scene detection, with higher values indicating better performance.

Implementation Details: We take $T = 14$ neighboring shots as input to our model. In MASRC, we set the activation functions $\sigma(\cdot)$ in Eqs. (2), (6) and (8) as ReLU (He et al. 2015). In Eq. (1), we specify the the number of top similar shots as $k = 4$. For model training, we employ the Adam (Kingma and Ba 2015) optimizer with a mini-batch size of 512. For fully supervised learning and self-supervised learning, we initialize the learning rate at 10^{-4} . In the case of self-supervised transfer learning, we set the initial learning rate to 10^{-3} for pre-training and reduce it to 10^{-5} for fine-tuning. Across all training stages, we apply a linear warm-up strategy during the initial epoch, followed by a learning rate decay according to a cosine schedule (He et al. 2019). We train our MASRC on a NVIDIA RTX 3060 GPU. In all experiments, we report the average of metrics across five different random seeds.

4.2 Comparison with State-of-the-Art Methods

Results on MovieNet. Table 1 displays the quantitative results on MovieNet (Huang et al. 2020). Given that nu-

Methods	Modalities	Training Par.	AP	mIoU	F1
Self Supervision					
BaSSL (Mun et al. 2022)	Entity	43.8 M	31.6	39.4	32.6
SSM (Gu, Goel, and Ré 2022)	Entity	32.5 M	33.3	38.1	32.2
TranS4mer (Islam et al. 2023)	Entity	32.0 M	34.5	39.6	33.4
VSMDB (Tan et al. 2024b)	Entity, Place	26.5 M	38.3	42.7	37.9
NeighborNet (Tan et al. 2024a)	Entity	35.5 M	51.2	52.9	46.4
MASRC (Ours)	Entity	10.7 M	47.3	48.6	42.0
	Place	12.9 M	48.9	47.0	41.9
	Entity, Place	21.4 M	52.8	56.8	53.0
Fully Supervision					
Temporal Perceiver (Tan et al. 2023)	Entity	52.1 M	53.3	53.2	-
MHRT (Wei et al. 2023)	Entity, Place, Audio	47.1 M	54.8	51.2	46.3
CANet (Tan, Wang, and Yuan 2024)	Face, Body	15.3 M	56.8	55.7	-
NeighborNet (Tan et al. 2024a)	Entity	35.5 M	64.0	61.2	57.8
MASRC (Ours)	Entity	10.7 M	59.2	58.5	55.0
	Place	12.9 M	59.6	57.9	52.9
	Entity, Place	21.4 M	67.4	65.8	63.8
Self-Supervised Transfer					
ShotCoL (Chen et al. 2021)	Entity	26.3 M	53.4	-	51.4
SCRL (Wu et al. 2022)	Entity	26.3 M	54.6	-	51.4
BaSSL (Mun et al. 2022)	Entity	43.8 M	57.4	50.7	47.0
SSM (Gu, Goel, and Ré 2022)	Entity	32.5 M	59.7	51.3	48.4
CAT (Yang et al. 2023)	Entity	43.8 M	59.6	53.7	51.9
TranS4mer (Islam et al. 2023)	Entity	32.0 M	60.8	51.9	48.4
VSMDB (Tan et al. 2024b)	Entity, Place	26.5 M	63.7	56.4	55.3
NeighborNet (Tan et al. 2024a)	Entity	35.5 M	71.9	64.5	62.7
MASRC (Ours)	Entity	10.7 M	65.5	62.6	64.0
	Place	12.9 M	65.9	61.4	61.3
	Entity, Place	21.4 M	73.2	70.7	67.3

Table 1: Performance comparison between our proposed method and recent baselines on the MovieNet dataset (Huang et al. 2020). The best results for each category of methods are indicated in **bold**.

Methods	OVSD	BBC
ShotCoL (Chen et al. 2021)	25.5	28.0
SCRL (Wu et al. 2022)	38.8	30.2
BaSSL (Mun et al. 2022)	28.7	40.0
TranS4mer (Islam et al. 2023)	36.0	43.6
NeighborNet (Tan et al. 2024a)	47.3	50.6
MASRC	48.3	53.2

Table 2: Cross dataset transfer result (AP) on OVSD (Rotman, Porat, and Ashour 2016) and BBC (Baraldi, Grana, and Cucchiara 2015) without fine-tuning.

merous methods rely solely on single-modal features, we showcase performance of our method in the corresponding modality for a fair comparison. Results show that, in contrast to counterparts, our method consistently excels with fewer training parameters across various experimental scenarios. Compared with multi-modal methods, our method employs less modality knowledge and outperforms the state-of-the-art MHRT (Wei et al. 2023). The superior performance of our approach arises from diversely modeling inter-shot relations and detecting scenes through comparisons between multi-shots relations.

Transfer Evaluation. We demonstrate the generalization capability of our MASRC in comparison with recent methods (Chen et al. 2021; Wu et al. 2022; Mun et al. 2022; Islam et al. 2023) in Table 2. All models used have undergone self-supervised pre-training and fine-tuning on MovieNet (Huang et al. 2020). We test these MovieNet-trained model without any additional fine-tuning on BBC (Baraldi, Grana, and Cucchiara 2015) and OVSD (Rotman, Porat, and Ashour 2016). The results

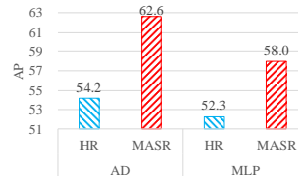


Figure 5: Comparison between our MASRC and HR (Wei et al. 2023) under the same detector, AD (Wei et al. 2023) or MLP (Islam et al. 2023).

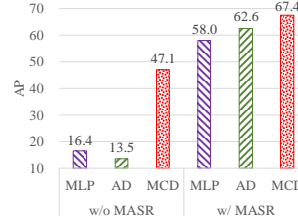


Figure 7: Comparison of detectors, including MLP (Islam et al. 2023), AD (Wei et al. 2023), and our MCD.

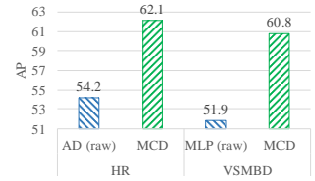


Figure 6: Comparison of our MCD with their raw detectors on temporal modeling methods, HR (Wei et al. 2023) and VSMDB (Tan et al. 2024b).

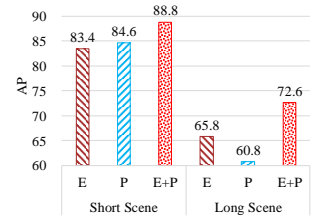


Figure 8: Influence of Entity detectors, including MLP (Islam et al. 2023), AD (Wei et al. 2023), and our MCD.

demonstrate that our proposed method achieves the best performance among all the comparisons on the OVSD and BBC datasets, which verifies the strong generalization capability of MASRC.

4.3 Ablation Studies

The following results are obtained in a fully supervised manner. We have also included more ablation experiments and visualizations in the supplementary.

Pluggability of MASRC. Fig. 5 provides the performance of HR (Wei et al. 2023) and our MASRC combined with different ending shot detectors, including MLP (Mun et al. 2022; Yang et al. 2023; Islam et al. 2023) and AD (Wei et al. 2023). The results show that combined with the same detector, the performance of our proposed MASRC is better than HR, highlighting the contribution and plug-and-play ability of our proposed MASRC.

Pluggability of MCD. Fig. 6 shows the performance of our proposed MCD combined with different temporal relation modeling methods, including HR (Wei et al. 2023) and VSMDB (Tan et al. 2024b). The results highlight that the combination of our proposed MCD with different methods is better than their raw alternatives, highlighting the advantage and versatility of MCD.

Different Ending Shot Detectors. Fig. 7 compares results about different detectors, where MLP (Mun et al. 2022; Yang et al. 2023; Islam et al. 2023) and AD (Wei et al. 2023) are two competitors of our MCD detector. The former makes predictions through single-shot representation, and the latter uses learnable class vectors to query each shot feature to decide whether each shot is a boundary. Our proposed MCD outperforms the competitors under evaluation metrics,

Entity		Place		AP
Short Term	Long Term	Short Term	Long Term	
✓	-	-	-	56.0
-	✓	-	-	59.2
-	-	-	✓	57.3
-	-	✓	-	59.6
✓	-	-	✓	63.2
-	✓	✓	-	67.4

Table 3: Comparison of different combinations of modalities and shot relations modeling.

ELD	PSD	MCD	AP
-	-	-	16.4
-	-	✓	47.1
✓	-	✓	59.2
-	✓	✓	59.6
✓	✓	✓	67.4

Table 4: Ablation study on MASRC in different inclusions of entity-based long-term dependency (ELD), place-based short-term dependency (PSD), and multi-shots comparison detection (MCD). ✓ signifies “included”, while - “excluded”.

mainly because MCD explicitly uses shot contexts to predict boundaries.

Impact of Each Modality on Detecting Video Scenes of Different Lengths. Fig. 8 provides an in-depth analysis of the impact of different modal inputs on video scene detection. Short scenes comprise fewer than 12 shots, while long scenes consist of a higher number of shots, based on the fact that the average number of scenes in MovieNet (Rao et al. 2020) is 12.6. As can be seen from Fig. 8, compared to the entity modality, the place modality shows a greater advantage in short scene detection. Conversely, the entity modality outperforms the place modality in detecting long scenes. This observation validates that the place modality is well-suited for capturing short-term shot relations, whereas the entity modality excels in capturing long-range shot relations. Notably, combined modalities consistently perform optimally under either scene length, confirming the complementarity of the two modalities.

Each Modality and Its Optimal Modeling Temporal Relations. Table 3 displays the performance of temporal relation modules with inputs from various modalities. The long-term relations module is modeled by Eqs. (1) and (2), and the short-term relations module is built upon Eqs. (3)-(8). For the entity modality, switching from short-term to long-term relations modeling improves AP by 3.2%. Conversely, for the place modality, switching to short-term relations modeling improves AP by 3.6%. The best performance is achieved by combining long-term entity modeling with short-term place modeling. In summary, the entity modality is suited for modeling long-term shot relations, whereas the place modality is adept at capturing short-term shot relations.

Network Components. Table 4 presents the results of assessing the contributions of different components of our pro-

Alternatives	Max Pooling	Self Attention	CNNs (used)	Unireplknet
AP	56.7	59.9	67.4	68.9

Table 5: Comparison of alternatives to our used CNNs.

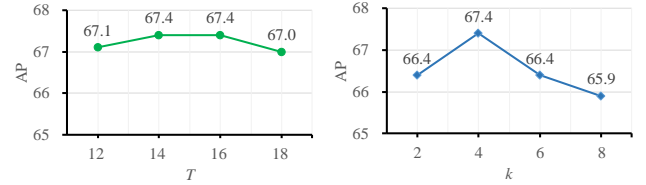


Figure 9: Impact of hyperparameters. (a) Temporal window scale T as defined in Sec. 3.2. (b) Entity feature neighbor scale k as defined by top- k in Eq. (1).

posed method to overall performance. The first row reports baseline results of feeding the concatenation of entity and place features into the MLP to predict shot classes. Replacing the MLP detector with the proposed MCD results in a significant improvement across all metrics, emphasizing the effects of contrasting shot contexts to detect ending shots. The addition of the ELD or PSD module alone leads to increment in AP, which arises from our consideration of the distinct roles of various visual semantics in video scene detection. The combination of ELD and PSD achieves peak performance, demonstrating their complementary effects.

Alternatives to CNNs in MCD. Table 5 presents performance of alternatives to CNNs in MCD. As shown, CNN-based models, such as UnireplkNet (Ding et al. 2024) and our used CNNs in Eq. (10), achieve better results compared to other alternatives, likely because CNN models are better suited to capturing the variation patterns of gridded data.

Temporal Window Scale T . Fig. 9a illustrates the impact of the time sliding window scale T , as defined in Sec. 3.2. It is evident that all metrics reach their peaks when $T = 14$. This observation aligns with the average number of shots per scene in the MovieNet dataset which is 12.6.

Entity Feature Neighbor Scale k . Fig. 9b presents the effects of the amount of most similar shots k defined in Eq. (1). The proposed method achieves optimal performance when $k = 4$. This selection strikes the best trade-off between capturing sufficient similar entity variations and avoiding the introduction of extraneous noisy information.

5 Conclusion

We propose a novel multi-modal shot relationship modeling and comparison framework for video scene detection. It reasons on the established long-range entity dependency graph and short-range place dependency graph to capture long-short dependencies between shots. For detecting ending shots, predictions are made by comparing the semantic relationships of surrounding shots to the target shot. Experimental results in public datasets show the effectiveness and superiority of our MASRC over previous SOTAs.

Acknowledgements

This work is supported in part by the National Natural Science Foundation of China under Grant 61976029 and the Key Project of Chongqing Technology Innovation and Application Development under Grant cstc2021jcsx-gksbX0033.

References

- Ba, L. J.; Kiros, J. R.; and Hinton, G. E. 2016. Layer Normalization. *CoRR*, abs/1607.06450.
- Bain, M.; Nagrani, A.; Brown, A.; and Zisserman, A. 2020. Condensed Movies: Story Based Retrieval with Contextual Embeddings. In *ACCV*, volume 12626, 460–479.
- Baraldi, L.; Grana, C.; and Cucchiara, R. 2015. A Deep Siamese Network for Scene Detection in Broadcast Videos. In *ACM MM*, 1199–1202.
- BBC. 2006. Planet earth. <https://www.bbc.co.uk/programmes/b006mywy>. Accessed: 2024-07-02.
- Chen, S.; Liu, C.; Hao, X.; Nie, X.; Arap, M.; and Hamid, R. 2023. Movies2Scenes: Using Movie Metadata to Learn Scene Representation. In *CVPR*, 6535–6544.
- Chen, S.; Nie, X.; Fan, D.; Zhang, D.; Bhat, V.; and Hamid, R. 2021. Shot Contrastive Self-Supervised Learning for Scene Boundary Detection. In *CVPR*, 9796–9805.
- Ding, X.; Zhang, Y.; Ge, Y.; Zhao, S.; Song, L.; Yue, X.; and Shan, Y. 2024. UniRepLKNet: A Universal Perception Large-Kernel ConvNet for Audio Video Point Cloud Time-Series and Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 5513–5524.
- Graves, A.; and Schmidhuber, J. 2005. Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Networks*, 18(5-6): 602–610.
- Gu, A.; Goel, K.; and Ré, C. 2022. Efficiently Modeling Long Sequences with Structured State Spaces. In *ICLR*.
- Hamilton, W. L.; Ying, Z.; and Leskovec, J. 2017. Inductive Representation Learning on Large Graphs. In *NeurIPS*, 1024–1034.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2015. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. In *ICCV*, 1026–1034.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep Residual Learning for Image Recognition. In *CVPR*, 770–778.
- He, T.; Zhang, Z.; Zhang, H.; Zhang, Z.; Xie, J.; and Li, M. 2019. Bag of Tricks for Image Classification with Convolutional Neural Networks. In *CVPR*, 558–567.
- He, X.; Deng, K.; Wang, X.; Li, Y.; Zhang, Y.; and Wang, M. 2020. LightGCN: Simplifying and Powering Graph Convolution Network for Recommendation. In *ACM SIGIR*, 639–648.
- Huang, Q.; Xiong, Y.; Rao, A.; Wang, J.; and Lin, D. 2020. MovieNet: A Holistic Dataset for Movie Understanding. In *ECCV*, volume 12349, 709–727.
- Islam, M. M.; Hasan, M.; Athrey, K. S.; Braskich, T.; and Bertasius, G. 2023. Efficient Movie Scene Detection using State-Space Transformers. In *CVPR*, 18749–18758.
- Kingma, D. P.; and Ba, J. 2015. Adam: A Method for Stochastic Optimization. In *ICLR*.
- Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
- Mun, J.; Shin, M.; Han, G.; Lee, S.; Ha, S.; Lee, J.; and Kim, E. 2022. BaSSL: Boundary-aware Self-Supervised Learning for Video Scene Segmentation. In *ACCV*, volume 13844, 485–501.
- Ngo, C.-W.; Pong, T.-C.; and Zhang, H. 2002. Motion-Based Video Representation for Scene Change Detection. *IJCV*, 50(2): 127–142.
- Odobez, J.; Gatica-Perez, D.; and Guillemot, M. 2003. Spectral Structuring of Home Videos. In *CIVR*, volume 2728, 310–320.
- Rao, A.; Xu, L.; Xiong, Y.; Xu, G.; Huang, Q.; Zhou, B.; and Lin, D. 2020. A Local-to-Global Approach to Multi-Modal Movie Scene Segmentation. In *CVPR*, 10143–10152.
- Rasheed, Z.; and Shah, M. 2005. Detection and representation of scenes in videos. *IEEE TMM*, 7(6): 1097–1105.
- Rotman, D.; Porat, D.; and Ashour, G. 2016. Robust and Efficient Video Scene Detection Using Optimal Sequential Grouping. In *ISM*, 275–280.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M. S.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 115(3): 211–252.
- Simonyan, K.; and Zisserman, A. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *ICLR*.
- Tan, J.; Wang, H.; Li, J.; Ou, Z.; and Qian, Z. 2024a. Neighbor Relations Matter in Video Scene Detection. In *CVPR*, 18473–18482.
- Tan, J.; Wang, H.; and Yuan, J. 2024. Characters Link Shots: Character Attention Network for Movie Scene Segmentation. *ACM TOMM*, 20(4): 94:1–94:23.
- Tan, J.; Wang, Y.; Wu, G.; and Wang, L. 2023. Temporal Perceiver: A General Architecture for Arbitrary Boundary Detection. *IEEE TPAMI*, 45(10): 12506–12520.
- Tan, J.; Yang, P.; Chen, L.; and Wang, H. 2024b. Temporal Scene Montage for Self-Supervised Video Scene Boundary Detection. *ACM TOMM*, 20(7): 215:1–215:19.
- Velickovic, P.; Cucurull, G.; Casanova, A.; Romero, A.; Liò, P.; and Bengio, Y. 2018. Graph Attention Networks. In *ICLR*.
- Vicol, P.; Tapaswi, M.; Castrejón, L.; and Fidler, S. 2018. MovieGraphs: Towards Understanding Human-Centric Situations From Videos. In *CVPR*, 8581–8590.
- Wei, X.; Shi, Z.; Zhang, T.; Yu, X.; and Xiao, L. 2023. Multimodal High-order Relation Transformer for Scene Boundary Detection. In *ICCV*, 22081–22090.
- Wu, H.; Chen, K.; Luo, Y.; Qiao, R.; Ren, B.; Liu, H.; Xie, W.; and Shen, L. 2022. Scene Consistency Representation Learning for Video Scene Segmentation. In *CVPR*, 14001–14010.

- Wu, Z.; Pan, S.; Chen, F.; Long, G.; Zhang, C.; and Yu, P. S. 2021. A Comprehensive Survey on Graph Neural Networks. *IEEE TNNLS*, 32(1): 4–24.
- Yan, S.; Xiong, Y.; and Lin, D. 2018. Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. In *AAAI*, 7444–7452.
- Yang, Y.; Huang, Y.; Guo, W.; Xu, B.; and Xia, D. 2023. Towards Global Video Scene Segmentation with Context-Aware Transformer. In *AAAI*, 3206–3213.
- Yao, L.; Mao, C.; and Luo, Y. 2019. Graph Convolutional Networks for Text Classification. In *AAAI*, 7370–7377.
- Zhou, B.; Lapedriza, À.; Khosla, A.; Oliva, A.; and Torralba, A. 2018. Places: A 10 Million Image Database for Scene Recognition. *IEEE TPAMI*, 40(6): 1452–1464.

Supplementary Material for Modality-Aware Shot Relating and Comparing for Video Scene Detection

The following results are obtained in a fully supervised manner on the MovieNet (Huang et al. 2020) dataset, unless otherwise specified.

6 Additional Experimental Results

6.1 Ablation Studies on Place Continuity Graph

Detail-Wide Shot Affiliation. The composition of detail-wide shot affiliation, as defined in Eq. (4), is ablated in Table 6. The results indicate that combining shot similarity and temporal proximity yields optimal performance. This underscores the importance of considering both shot similarity and temporal proximity when determining the affiliation between wide shots and detail shots.

Local-Global Shot Affiliation	AP	mIoU	F1
Shot Similarity S_{ij}^P	65.2	64.7	63.3
Temporal Proximity D_{ij}	63.9	60.1	59.9
Both	67.4	65.8	63.8

Table 6: Comparison of variants of local-global shots affiliation as defined in Eq. (4).

Necessary of Using Both Detail-to-Wide (D2W) and Wide-to-Detail (W2D) Graphs. Table 7 presents the results of assessing the contributions of different components of our proposed PCG to overall performance. Combining D2W and W2D achieves significantly better performance compared to using only D2W or only W2D.

D2W	W2D	AP
✓	-	63.6
-	✓	65.3
✓	✓	67.4

Table 7: Ablation study on PCG graph in different inclusions of D2W and W2D graphs. ✓ signifies “included”, while - “excluded”.

Visualization of Edge Connections in EJG and PCG. In Fig. 10, we present a sample result of edge connections of EJG and PCG, built upon Sec. 3.2. As expected, EJG connects similar but distant shots, while PCG links consecutive shots depicting the same location. Consequently, EJG is advantageous for modeling long-range shot relations, whereas PCG is effective for modeling short-range shot relations.

6.2 Comparison with MHRT under Various Modality Combinations

Fig. 11 compares the performance of our MASRC with MHRT¹ (Wei et al. 2023), a state-of-the-art multi-modal

¹The code for MHRT is not available; hence, we reproduced it for comparisons.

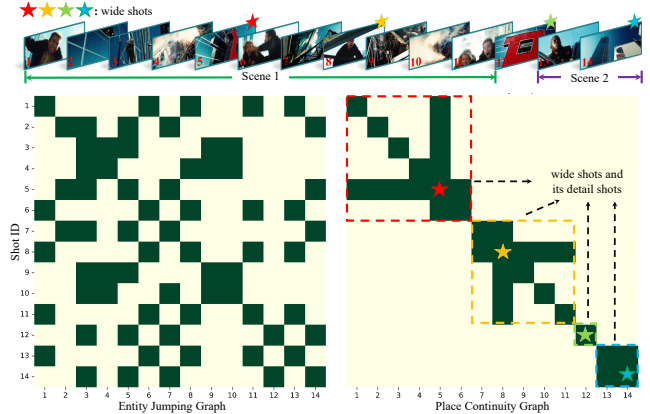


Figure 10: Visualization of edges connections (in dark color) of EJG and PCG. Stars indicate different wide shots discovered by the place modality, and the others are detail shots.

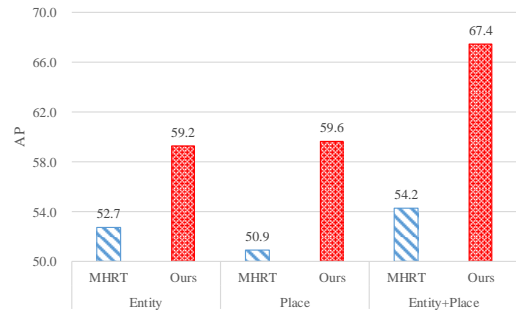


Figure 11: Comparison between proposed MASRC and MHRT (Wei et al. 2023) under various combinations of multi-modality semantic inputs.

video scene detection method, under the same entity and place feature inputs. Our model beats MHRT regardless of the combination of different visual semantic modalities. This superiority of ours is mainly thanks to diverse shot relationship modeling and context comparison convolution for ending shot detection.

6.3 Computation Cost Analysis

Under identical hardware conditions, Table 8 presents the number of training parameters, the training/inference throughput measured in samples per second (Sam/s), and GPU memory costs per sample (GB/Sam.) for state-of-the-art methods. A sample corresponds to a time window with a length of 14 shots, equating to 14 nodes in our graph. As can be seen, our model requires fewer resources in terms of training params, training/inference throughput, and GPU memory costs compared to others utilizing the same backbone.



Figure 12: Qualitative comparison of the proposed method with the previous method MHRT (Wei et al. 2023). GT denotes ground-truth scene boundaries for reference. The colored vertical lines represent the video scene boundaries.



Figure 13: Qualitative comparison of the proposed method with the previous method MHRT (Wei et al. 2023). GT denotes ground-truth scene boundaries for reference. The red vertical lines represent the video scene boundaries.

Methods	Modalities	Train Par.	Train		Inference	
			Sam./s	GB/sam.	Sam./s	GB/sam.
LGSS (Rao et al. 2020)	P	13.7 M	32	0.37	68	0.17
Ours	P	12.9 M	3788	0.0013	5632	0.0005
MHRT (Wei et al. 2023)	E+P	27.8 M	2509	0.0020	3584	0.0014
Ours	E+P	21.4 M	3093	0.0019	5171	0.0009

Table 8: Number of training parameters (Train Par.), training/inference throughput measured in samples per second (Sam./s), and GPU memory costs per sample (GB/Sam.) for state-of-the-art methods. E is short for entity modality and P is for place modality.

6.4 Visualization of Short Scene Detection.

In Fig. 12, we present a sample results of video scene detection on the MovieNet dataset. Ground-truth annotations are provided for reference, and we include the previous state-of-the-art method, MHRT (Wei et al. 2023), for comparison. The results demonstrate that MHRT is susceptible to camera shot changes, including switching, zooming, and movement, leading to over-segmented video scenes. In contrast, our proposed method effectively handles these challenging cases without experiencing over-segmentation. This is attributed to our approach on establishing the diverse relations between shots and carefully comparing the relations between shots.

6.5 Visualization of Long Scene Detection.

Fig. 13 depicts a result of long scene consisting of 76 shots. It portrays a taxi driver picking up two separate groups of passengers during his night shift before finishing his work. The sequence involves numerous rapid shot transitions. Previous SOTA (Wei et al. 2023) incorrectly segments the scene into 12 scenes, while our method reduces false positive detection and segments the scene into 5 scenes. This is attributed to our approach on establishing modality-aware shot relations and comparing the relations between shots.

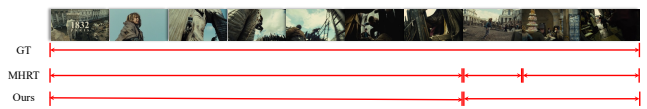


Figure 14: Qualitative comparison of the proposed method with the previous method MHRT (Wei et al. 2023). GT denotes ground-truth scene boundaries for reference. The red vertical lines represent the video scene boundaries.

6.6 Limitation.

We find some commonalities in cases where our method and previous SOTA method (Wei et al. 2023) may fail. Fig. 14 presents a short scene from "Les Misérables". Both our method and previous SOTA (Wei et al. 2023) segment this scene into two parts. The inconsistency between shots arises from frequent motion transitions, which weaken the coherence of entities and places, leading to incorrect segmentation. Addressing this challenge may involve incorporating background audio, such as music, which can provide another type of contextual consistency. Exploring the integration of audio into our model is a future research direction.

6.7 Impact of MCD on Detecting Video Scenes of Different Lengths

Fig. 15 provides a deeper analysis of the impact of the proposed MCD on video scene detection. We let length 12 to distinguish between short scenes (each of which has ≤ 12 shots) and long scenes (each of which has > 12 shots) based on the fact that the average shot number in MovieNet scenes is 12.6. As shown in Fig. 15, compared with the model without MCD, the model with MCD improves both short and long scene detection. Notably, in long scene detection, the inclusion of MCD (72.6% AP) significantly enhances the

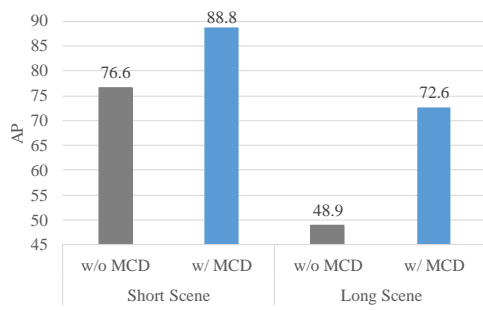


Figure 15: Influence of w/o and w/ MCD on detecting video scenes with varied lengths.

performance compared to the exclusion of MCD.