

000  
001  
002  
003  
004 **Multi-view Spectral Polarization Propagation for Video Glass Segmentation**  
005 **(Supplementary Material)**  
006  
007  
008  
009       Anonymous ICCV submission  
010  
011       Paper ID 6969  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065

This supplementary document provides more details of the proposed PGV-117 dataset (§ 1), the formal definitions of the four quantitative metrics (§ 2), and the detailed calculation of key affinity (§ 3). Five processed video sequences are provided along with this document. The teaser and the visual results shown in section 5.2 of this submission are from these five videos. We also offer additional video results through [Google Drive](#).

## 1. PGV-117 Dataset

The ground truth masks of the proposed PGV-117 dataset are annotated by annotation professionals, resulting in 144,686 glass masks. Each ground truth mask is manually checked to ensure the quality of the annotations.

The proposed dataset consists of 117 sequences and 21,485 frames. The training set offers 85 sequences, 15,838 frames, and the testing set provides 32 sequences, 5,647 frames. [Figure 1](#) and [Figure 2](#) show the number of frames for each sequence in the training and testing set, respectively.

## 2. Formal Definition of Evaluation Metrics

We adopt the four metrics used by Mei *et al.* [5] for evaluating all competing approaches, which are intersection over union (IoU), weighted F-measure ( $F_\beta$ ) [4], mean absolute error (MAE), and balance error rate (BER) [6]. Here, we provide the formal definitions of these four metrics.

## Intersection over union (IoU)

$$IoU = \frac{\sum_{i=1}^H \sum_{j=1}^W (G(i, j) * P_b(i, j))}{\sum_{i=1}^H \sum_{j=1}^W (G(i, j) + P_b(i, j) - G(i, j) * P_b(i, j))}, \quad (1)$$

where  $G$  is the ground truth mask in which the values of the glass region are 1 while those of the non-glass region are 0;  $P_b$  is the predicted mask binarized with a threshold of 0.5; and  $H$  and  $W$  are the height and width of the ground truth mask, respectively.

**Weighted F-measure ( $F_\beta$ )** takes a prediction map's precision and recall into account, which is a common metric used in salient object detection tasks. Based on recent studies [2, 3], the weighted F-measure [4] is more reliable than the traditional  $F_\beta$  [5], and it is used in our evaluation.

**Mean Absolute Error (MAE)** calculates the element-wise distance between a prediction map  $P$  and the corresponding ground truth mask  $G$ :

$$MAE = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W |P(i, j) - G(i, j)| \quad (2)$$

where  $P(i, j)$  indicates the predicted probability score at location  $(i, j)$ .

## ICCV 2023 Submission #6969. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

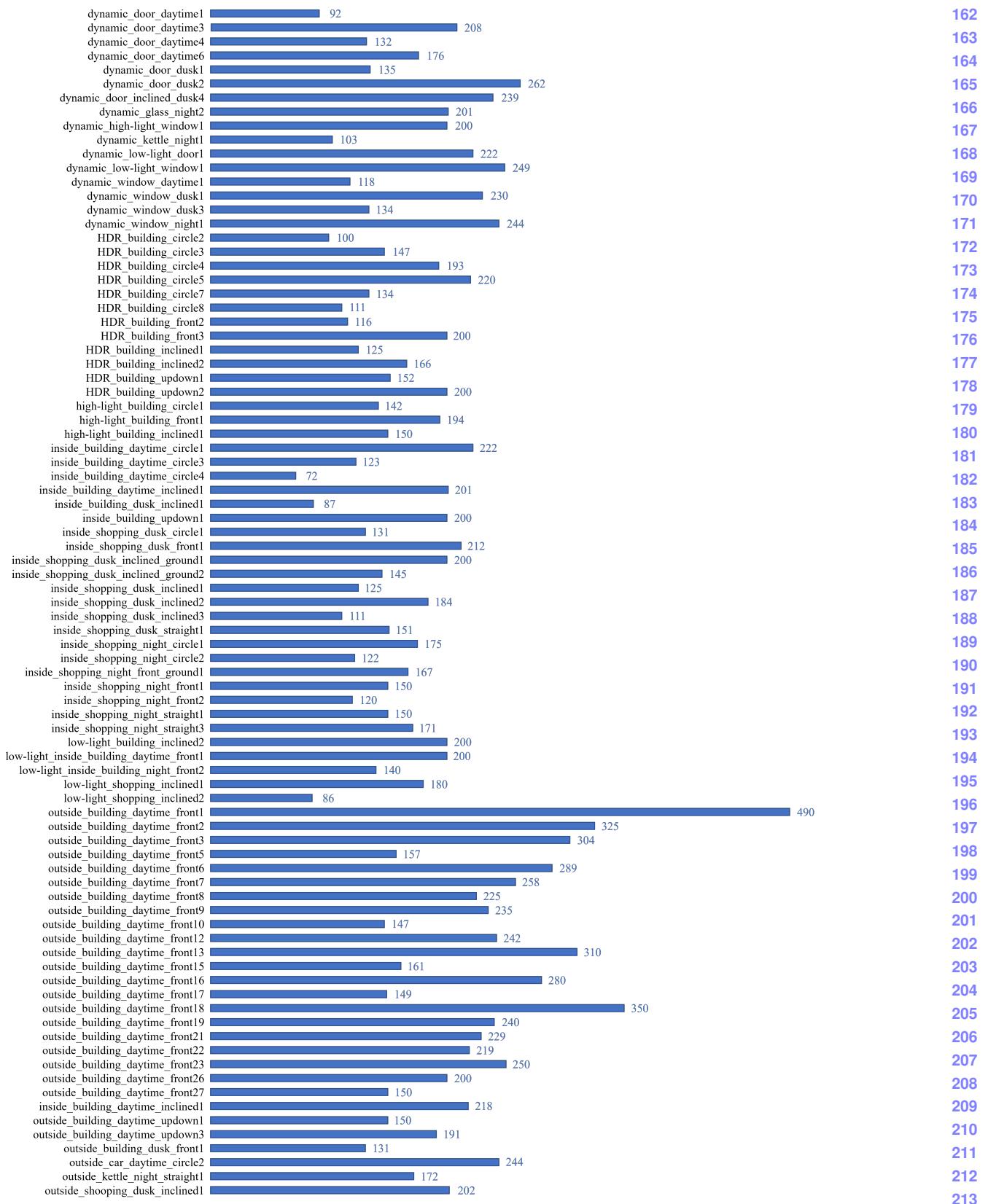


Figure 1. Summary for training set distribution of PGV-117, which includes 85 video sequences in total.

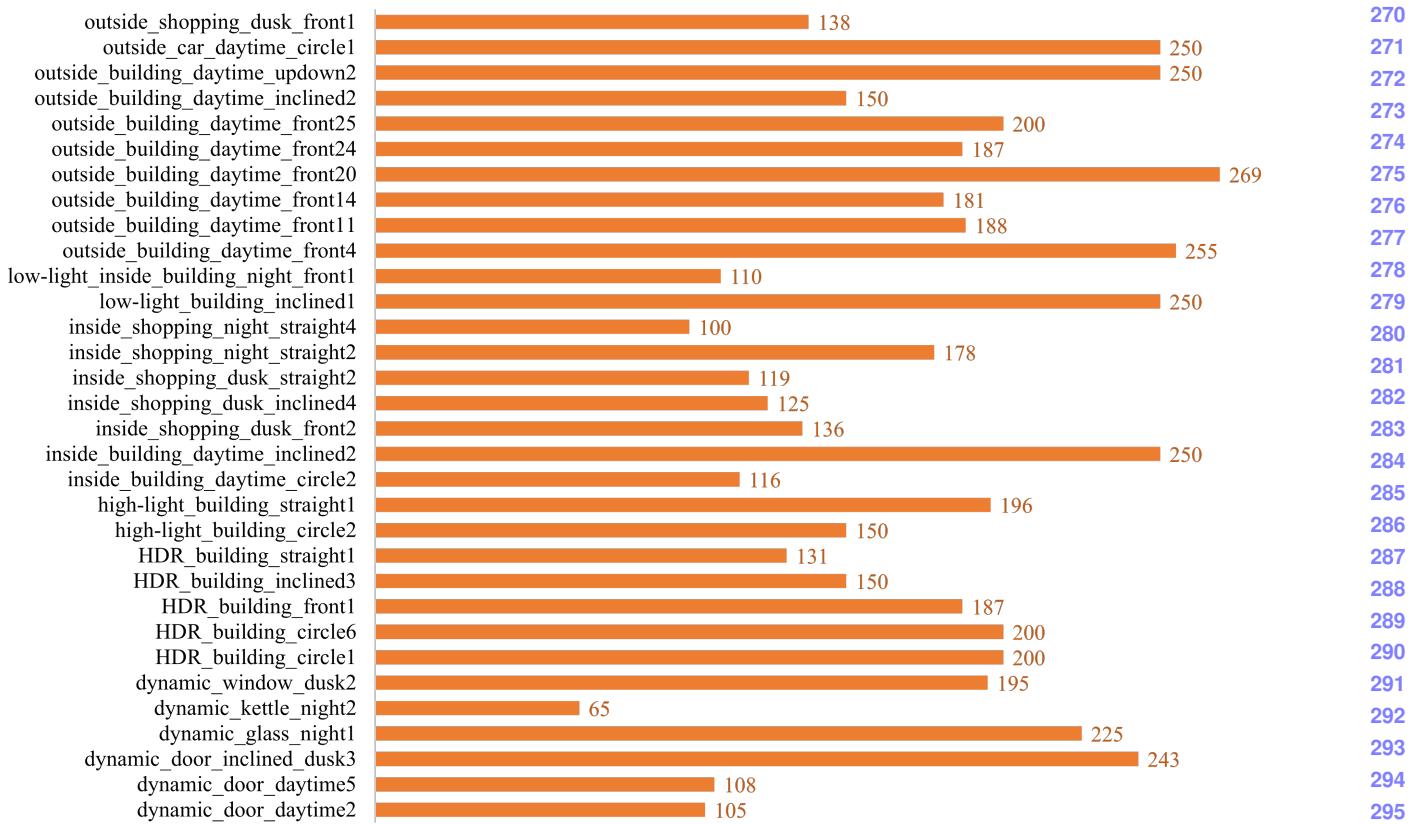


Figure 2. Summary for testing set distribution of PGV-117, which includes 32 video sequences in total. In order not to lose generality, all lighting conditions, camera motion patterns, and dynamics are also included in the testing set.

**Balance error rate (BER)** is a common metric used in shadow detection tasks. Formally, it is defined as:

$$BER = \left(1 - \frac{1}{2} \left( \frac{TP}{N_p} + \frac{TN}{N_n} \right) \right) \times 100 \quad (3)$$

where  $TP$ ,  $TN$ ,  $N_p$ , and  $N_n$  represent the numbers of true positive pixels, true negative pixels, glass pixels, and non-glass pixels, respectively.

### 3. Computing Affinity

For the query frame  $t$ , we relate multi-view RGB-P information by exploring the relationship between the PGI key of  $t$  with the keys in the memory (0 to  $t - 1$ ). After generating the query key  $k^Q$  and memory keys  $k^M$ , we refer to [7, 1] to calculate the affinity between  $k^Q$  and  $k^M$ :

$$\begin{aligned} a &= \xi[(k^M)^2], \\ b &= 2 * [(k^M)^T \circledast k^Q], \\ A &= (-a + b) / \sqrt{CK}, \\ A &= \frac{\exp(A_{ij})}{\sum_n(\exp(A_{nj}))}. \end{aligned} \quad (4)$$

where  $\xi[\cdot]$  represents the summation and unsqueeze operation,  $\circledast$  means matrix multiplication.  $CK$  is the number of channels for the key features, and  $i$  denotes the affinity value at the  $i$ -th position. The affinity considers both RGB similarity and multi-view spectral consistency between the query and memory frames. The value encoder output  $v^M$  corresponding to  $k^M$  contains features from the tripartite memory values, and the multiplication of affinity and  $v^M$  correlates the query frame and historical information to obtain  $v^Q$  to participate in the readout of the memory bank.

324	<b>References</b>	378
325		379
326	[1] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient	380
327	video object segmentation. In <i>NeurIPS</i> , 2021. 3	381
328	[2] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In	382
329	<i>ICCV</i> , 2017. 1	383
330	[3] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary fore-	384
331	ground map evaluation. In <i>IJCAI</i> , 2018. 1	385
332	[4] Ran Margolin, Lih Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps? In <i>CVPR</i> , 2014. 1	386
333	[5] Haiyang Mei, Bo Dong, Wen Dong, Jiaxi Yang, Seung-Hwan Baek, Felix Heide, Pieter Peers, Xiaopeng Wei, and Xin Yang. Glass	387
334	segmentation using intensity and spectral polarization cues. In <i>CVPR</i> , 2022. 1	388
335	[6] Vu Nguyen, Tomas F Yago Vicente, Maozheng Zhao, Minh Hoai, and Dimitris Samaras. Shadow detection with conditional generative	389
336	adversarial networks. In <i>ICCV</i> , 2017. 1	390
337	[7] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In	391
338	<i>ICCV</i> , 2019. 3	392
339		393
340		394
341		395
342		396
343		397
344		398
345		399
346		400
347		401
348		402
349		403
350		404
351		405
352		406
353		407
354		408
355		409
356		410
357		411
358		412
359		413
360		414
361		415
362		416
363		417
364		418
365		419
366		420
367		421
368		422
369		423
370		424
371		425
372		426
373		427
374		428
375		429
376		430
377		431