000
001
002
003
004
005
006
007
008
009
010
011
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
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

# Multi-view Spectral Polarization Propagation for Video Glass Segmentation (Supplementary Material)

Anonymous ICCV submission

Paper ID 6969

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)**

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

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 $G$:

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

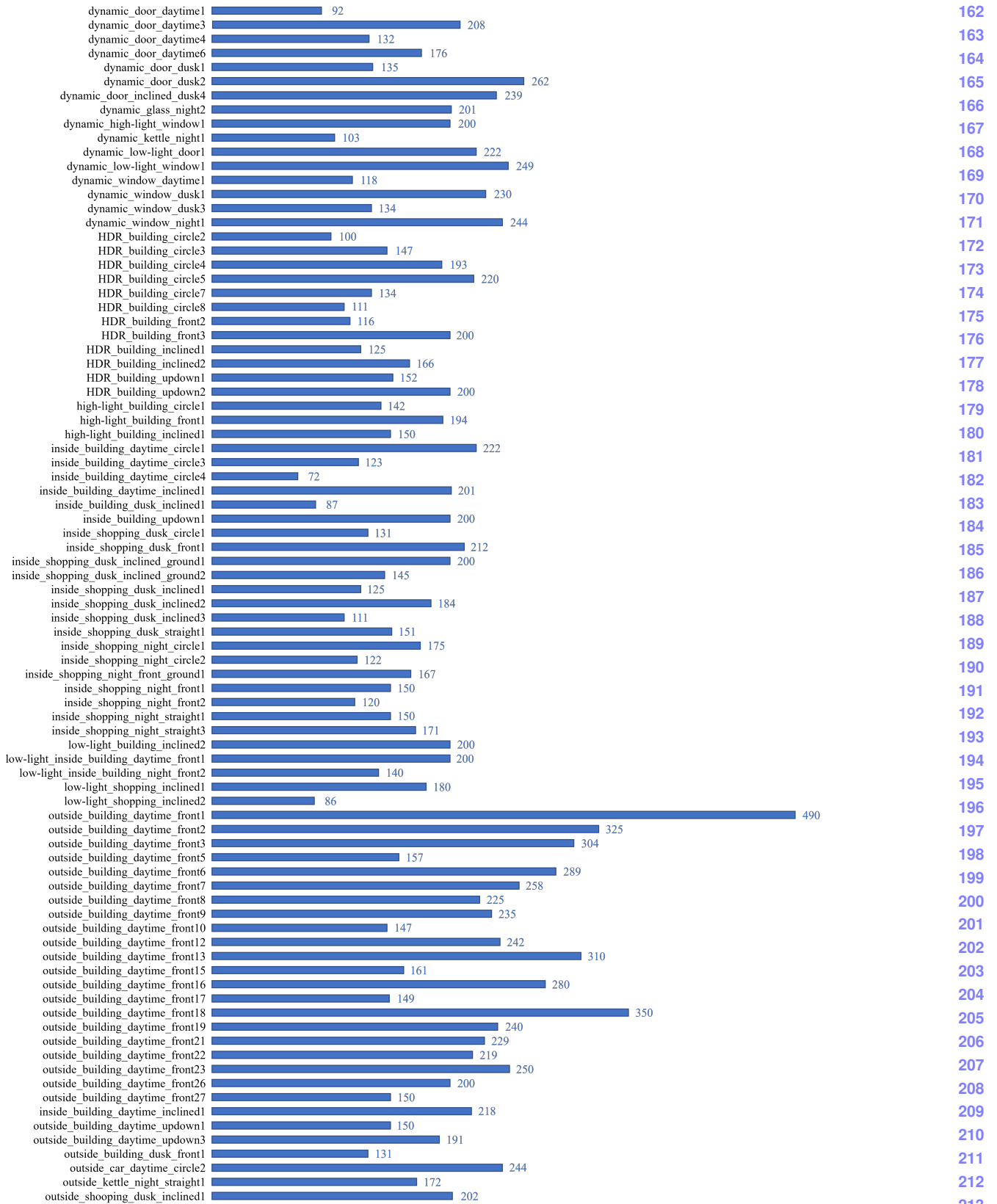where $P(i,j)$ indicates the predicted probability score at location $(i,j)$.

ICCV
#6969

ICCV
#6969

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

| Sequence | Count |
|---|---|
| dynamic_door_daytime1 | 92 |
| dynamic_door_daytime3 | 208 |
| dynamic_door_daytime4 | 132 |
| dynamic_door_daytime6 | 176 |
| dynamic_door_dusk1 | 135 |
| dynamic_door_dusk2 | 262 |
| dynamic_door_inclined_dusk4 | 239 |
| dynamic_glass_night2 | 201 |
| dynamic_high-light_window1 | 200 |
| dynamic_kettle_night1 | 103 |
| dynamic_low-light_door1 | 222 |
| dynamic_low-light_window1 | 249 |
| dynamic_window_daytime1 | 118 |
| dynamic_window_dusk1 | 230 |
| dynamic_window_dusk3 | 134 |
| dynamic_window_night1 | 244 |
| HDR_building_circle2 | 100 |
| HDR_building_circle3 | 147 |
| HDR_building_circle4 | 193 |
| HDR_building_circle5 | 220 |
| HDR_building_circle7 | 134 |
| HDR_building_circle8 | 111 |
| HDR_building_front2 | 116 |
| HDR_building_front3 | 200 |
| HDR_building_inclined1 | 125 |
| HDR_building_inclined2 | 166 |
| HDR_building_updown1 | 152 |
| HDR_building_updown2 | 200 |
| high-light_building_circle1 | 142 |
| high-light_building_front1 | 194 |
| high-light_building_inclined1 | 150 |
| inside_building_daytime_circle1 | 222 |
| inside_building_daytime_circle3 | 123 |
| inside_building_daytime_circle4 | 72 |
| inside_building_daytime_inclined1 | 201 |
| inside_building_dusk_inclined1 | 87 |
| inside_building_updown1 | 200 |
| inside_shopping_dusk_circle1 | 131 |
| inside_shopping_dusk_front1 | 212 |
| inside_shopping_dusk_inclined_ground1 | 200 |
| inside_shopping_dusk_inclined_ground2 | 145 |
| inside_shopping_dusk_inclined1 | 125 |
| inside_shopping_dusk_inclined2 | 184 |
| inside_shopping_dusk_inclined3 | 111 |
| inside_shopping_dusk_straight1 | 151 |
| inside_shopping_night_circle1 | 175 |
| inside_shopping_night_circle2 | 122 |
| inside_shopping_night_front_ground1 | 167 |
| inside_shopping_night_front1 | 150 |
| inside_shopping_night_front2 | 120 |
| inside_shopping_night_straight1 | 150 |
| inside_shopping_night_straight3 | 171 |
| low-light_building_inclined2 | 200 |
| low-light_inside_building_daytime_front1 | 200 |
| low-light_inside_building_night_front2 | 140 |
| low-light_shopping_inclined1 | 180 |
| low-light_shopping_inclined2 | 86 |
| outside_building_daytime_front1 | 490 |
| outside_building_daytime_front2 | 325 |
| outside_building_daytime_front3 | 304 |
| outside_building_daytime_front5 | 157 |
| outside_building_daytime_front6 | 289 |
| outside_building_daytime_front7 | 258 |
| outside_building_daytime_front8 | 225 |
| outside_building_daytime_front9 | 235 |
| outside_building_daytime_front10 | 147 |
| outside_building_daytime_front12 | 242 |
| outside_building_daytime_front13 | 310 |
| outside_building_daytime_front15 | 161 |
| outside_building_daytime_front16 | 280 |
| outside_building_daytime_front17 | 149 |
| outside_building_daytime_front18 | 350 |
| outside_building_daytime_front19 | 240 |
| outside_building_daytime_front21 | 229 |
| outside_building_daytime_front22 | 219 |
| outside_building_daytime_front23 | 250 |
| outside_building_daytime_front26 | 200 |
| outside_building_daytime_front27 | 150 |
| inside_building_daytime_inclined1 | 218 |
| outside_building_daytime_updown1 | 150 |
| outside_building_daytime_updown3 | 191 |
| outside_building_dusk_front1 | 131 |
| outside_car_daytime_circle2 | 244 |
| outside_kettle_night_straight1 | 172 |
| outside_shooping_dusk_inclined1 | 202 |

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

ICCV
#6969

ICCV
#6969

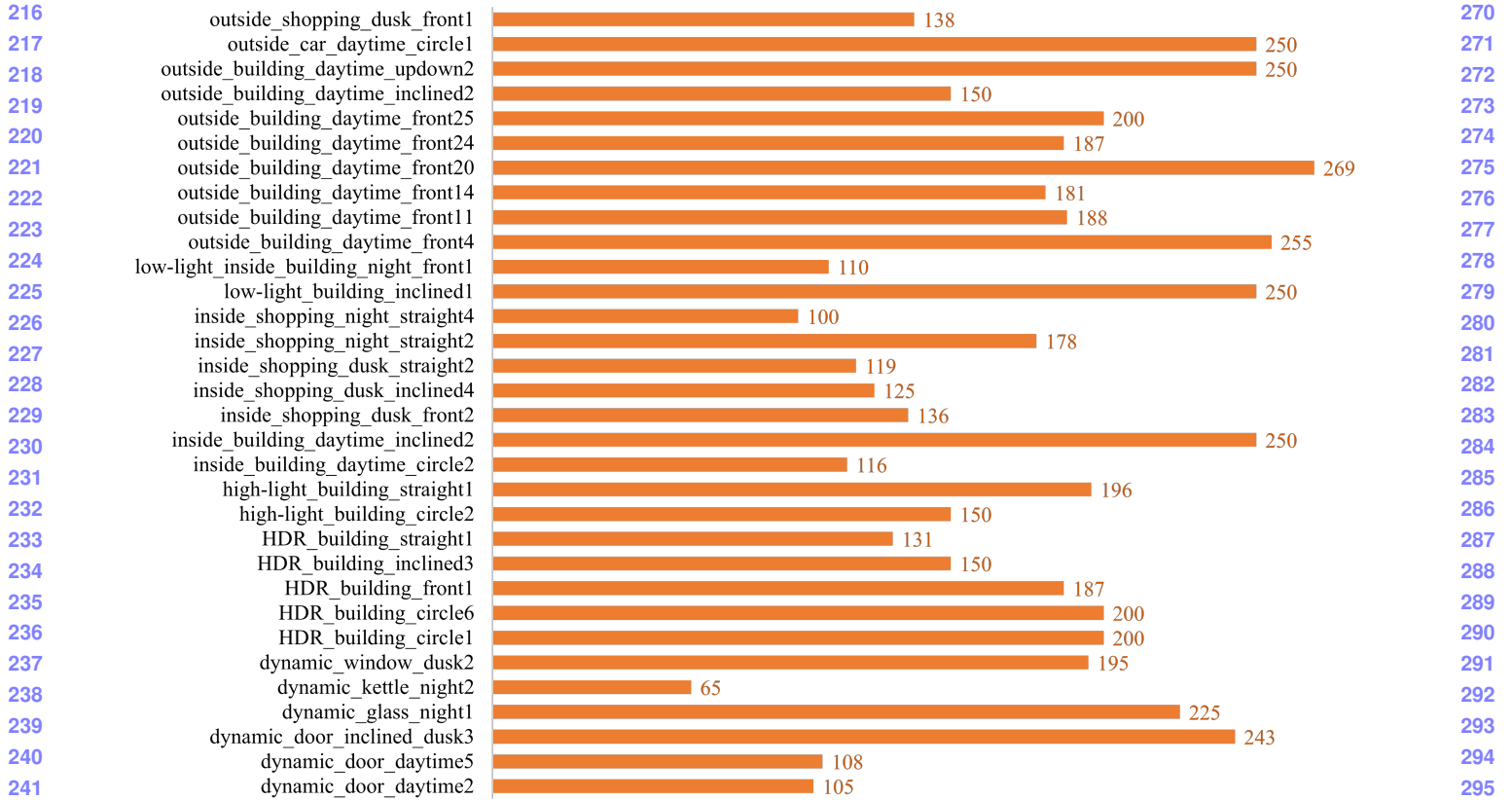ICCV 2023 Submission #6969. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

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:

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

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 $k^M$:

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

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.

ICCV
#6969

ICCV
#6969

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

# References

[1] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. In *NeurIPS*, 2021. 3

[2] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *ICCV*, 2017. 1

[3] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. In *IJCAI*, 2018. 1

[4] Ran Margolin, Lihi Zelnik-Manor, and Ayellet Tal. How to evaluate foreground maps? In *CVPR*, 2014. 1

[5] Haiyang Mei, Bo Dong, Wen Dong, Jiaxi Yang, Seung-Hwan Baek, Felix Heide, Pieter Peers, Xiaopeng Wei, and Xin Yang. Glass segmentation using intensity and spectral polarization cues. In *CVPR*, 2022. 1

[6] Vu Nguyen, Tomas F Yago Vicente, Maozheng Zhao, Minh Hoai, and Dimitris Samaras. Shadow detection with conditional generative adversarial networks. In *ICCV*, 2017. 1

[7] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *ICCV*, 2019. 3