

# CVPR'21 Workshop "Vision for All Seasons: Adverse Weather and Lighting Conditions"

[Workshop at CVPR21 <link>](#)

# 1

Keynote - Prof. Wolfram Burgard: Exploiting Knowledge  
from Multiple Modalities for Robust Perception

# 1.A

Idea of Using Thermal Cameras for Transferring Capabilities from  
Day Time to Night Time

# Exploiting Knowledge from Multiple Modalities for Robust Perception

Wolfram Burgard

---

Joint work with: Jannik Zürn, Johan Vertens, Kshitij Sirohi, Rohit  
Mohan, Abhinav Valada ...



UNI  
FREIBURG

AIS Autonomous Intelligent Systems

# Motivation

- We want to minimize labeling efforts
  - New tasks
  - Domain transfer
- In this talk: **using multi-modal setups**

# Robust Semantic Segmentation by Domain Adaption

## HeatNet: Bridging the Day-Night Domain Gap in Semantic Segmentation with Thermal Images

Johan Vertens\*, Jannik Zürn\*, and Wolfram Burgard

**Abstract**— The majority of learning-based semantic segmentation methods are optimized for daytime scenarios and favorable environmental conditions. Real-world driving scenarios, however, entail adverse environmental conditions such as low ambient illumination or glare which remains a challenge for existing semantic segmentation models. To address this challenge, we propose a multimodal semantic segmentation model that can be applied during daytime and nighttime. To this end, besides RGB images, we leverage thermal images to capture the scene under varying illumination levels. We avoid the expensive annotation of nighttime images by leveraging an existing daytime dataset and propose a knowledge transfer framework that transfers the daytime model's knowledge to the nighttime domain. We further employ a novel two-stage training scheme to simultaneously train across the domains and propose a novel two-stage training scheme. Furthermore, due to a lack of thermal data for semantic segmentation, we propose a large-scale dataset containing over 20,000 time-synchronized and aligned RGB-thermal image pairs. In this work, we present a novel target-loss function that allows for instantaneous robust extrinsic and intrinsic thermal camera calibrations. Among others, we empirically show new state-of-the-art results for nighttime semantic segmentation.

I. INTRODUCTION

Robust and accurate semantic segmentation of urban scenes is one of the enabling technologies for autonomous driving in complex and cluttered driving scenarios. Recent research has shown great promise in RGB datasets for autonomous driving [36], [5], which were predominantly demonstrated in favorable daytime illumination conditions. While the reported results demonstrate high accuracies on these datasets [36], [5], [18], these datasets generalize poorly to adverse weather conditions and low illumination levels present at nighttime. This constraint becomes especially apparent in rural areas where artificial lighting is often absent. In addition, the lack of semantic segmentation and situation awareness, robust perception in these conditions is a vital prerequisite.

Transfer learning and domain adaptation approaches aim at narrowing the domain gap between a source domain, where supervised learning from labelled data is possible, to a target domain, where labelled data is either sparse or not available. Such approaches, as demonstrated in [28] or [35], allow for learning representations that are invariant to the source domain. These approaches, however, do not leverage a complementary modality such as thermal infrared images that can contain more relevant information to solve a given task.

\*These authors contributed equally. All authors are with the University of Freiburg Institute for Computer Science, with the Toyota Research Institute, Los Altos, USA. Corresponding author: vertens@informatik.uni-freiburg.de

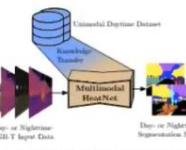
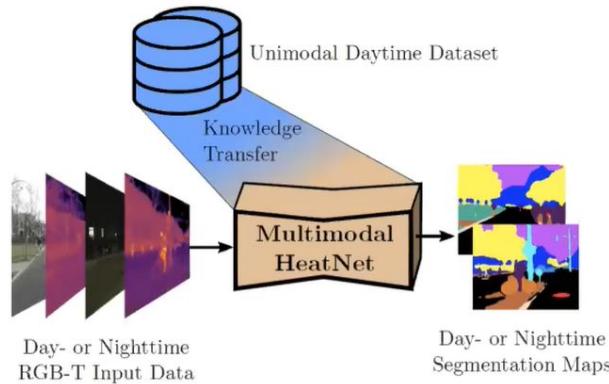


Fig. 1. Our multimodal segmentation network leverages both daytime and nighttime images. We transfer relevant knowledge from a large-scale unimodal daytime dataset to a multimodal segmentation model with a teacher model to simultaneously adapt our model to the nighttime domain by unsupervised domain adaptation.

**HeatNet: Bridging the Day-Night Domain Gap in Semantic Segmentation with Thermal Images. Johan Vertens, Jannik Zürn and Wolfram Burgard, IROS 2020**



# INSIGHTS

- IRIS Published.
- Focused on - how we can actually leverage - the thermal cameras - for better semantic segmentation, in particular, - for situations - where normal RGB cameras face difficulty - like the low light vision.
- Labelling nighttime images is very difficult for semantic segmentation task.

# Motivation



Semantic segmentation prediction during nighttime using a conventional CNN trained on publicly available datasets

**Labeling night-time images is extremely painful!**

# INSIGHTS

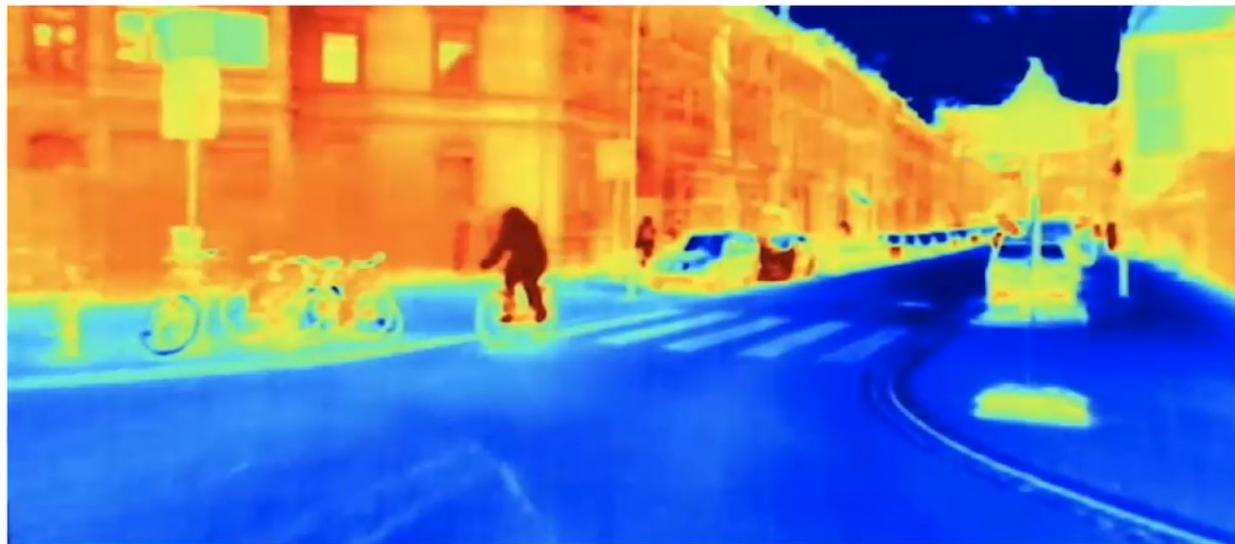
- The idea of this work is - to leverage daytime semantic segmentation AND another modality - to bootstrap semantic segmentation for nighttime.
- NEXT SLIDE - A Typical Night Scene
- AND THEN - Overlaying with Thermal Infrared Images.

# Motivation



Thermal infrared images exhibit small domain gap between day- and nighttime

# Motivation

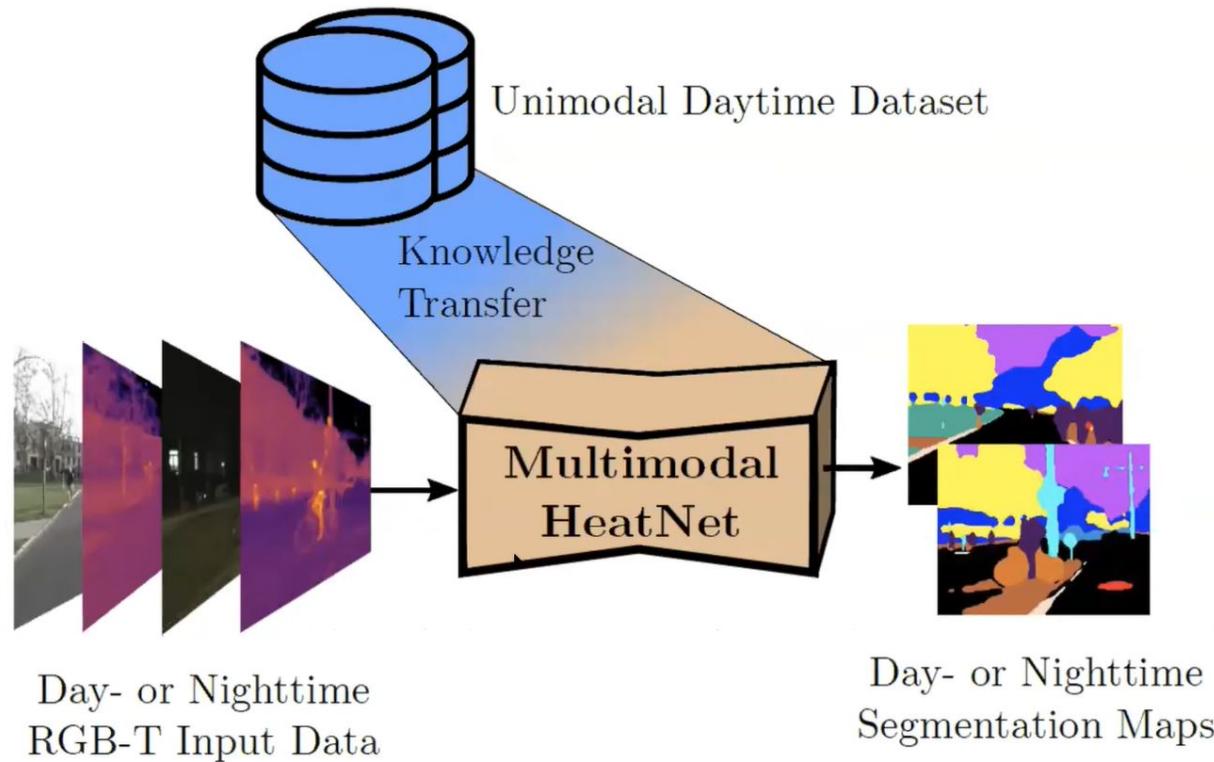


Thermal infrared images exhibit small domain gap between day- and nighttime

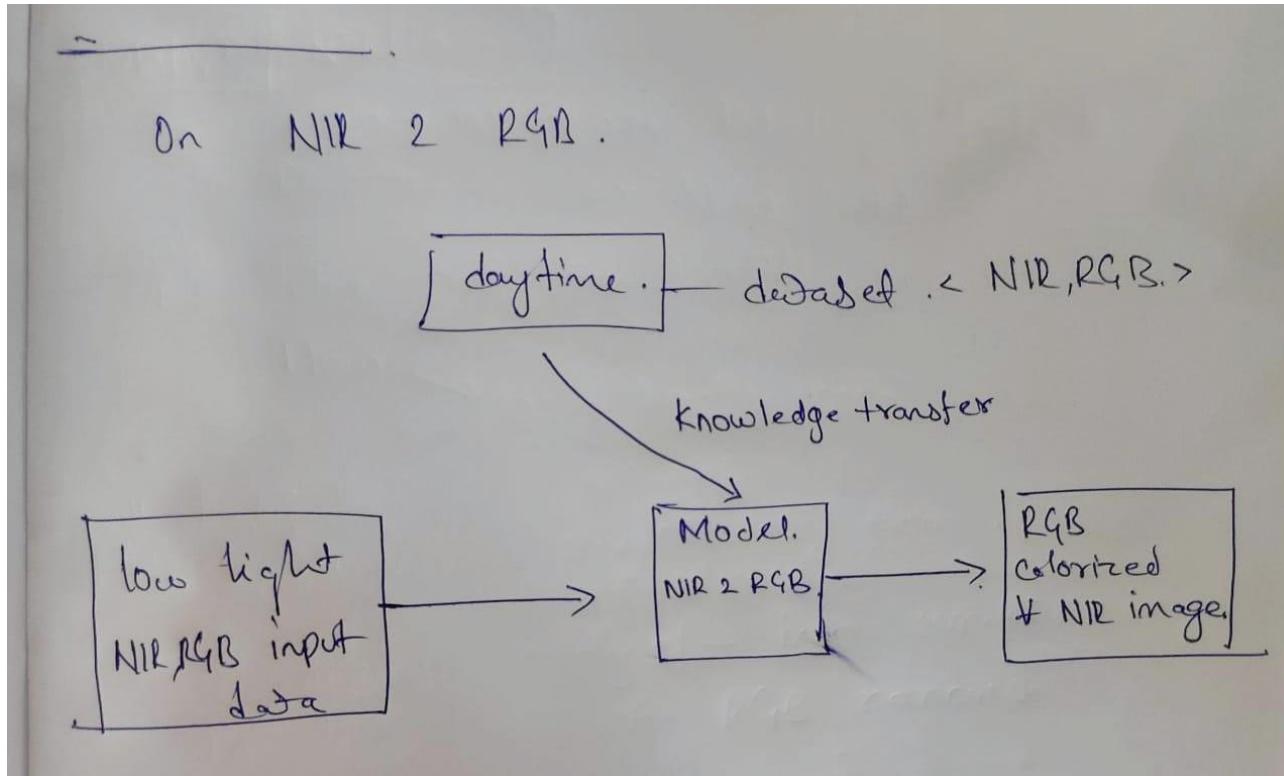
# INSIGHTS

- THERMAL INFRARED IMAGES exhibit small domain gap between day and night time.
- **NEXT SLIDE :** Shows the Approach - To Train a Multimodal HeatMap - That takes RGB - Thermal Images - AND - Creates Semantic Segmentation Map - For Day & Night time Images.
- At the same time, - have to leverage all the knowledge we have - for Daytime Semantic Segmentation.

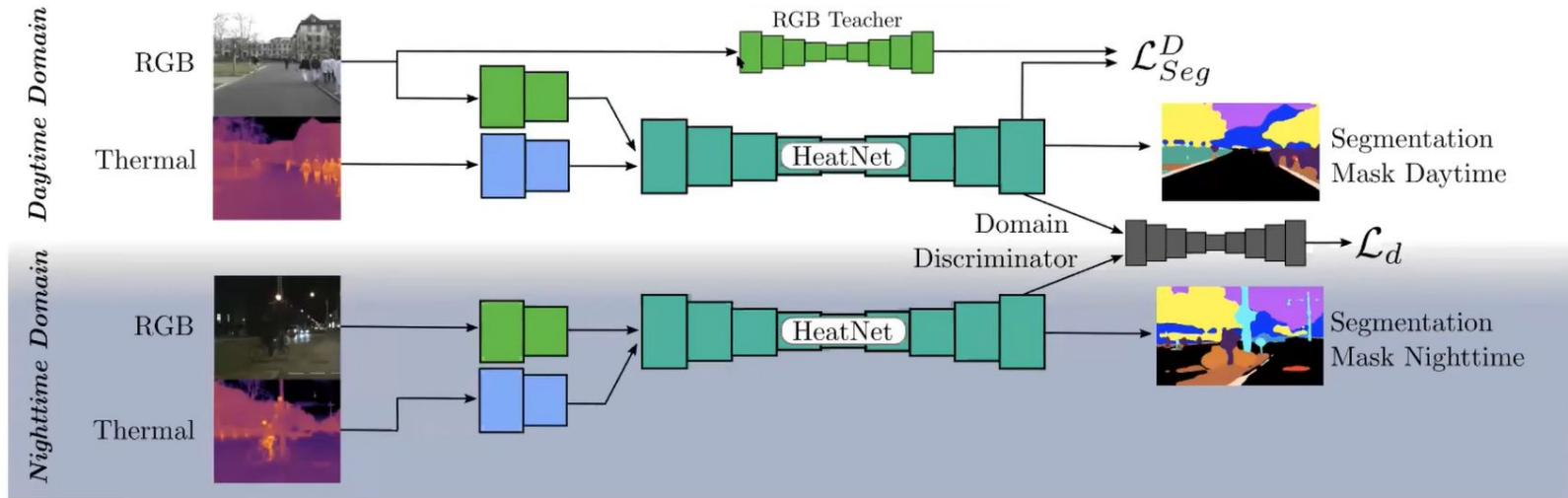
# Approach



# INSIGHTS



# Approach



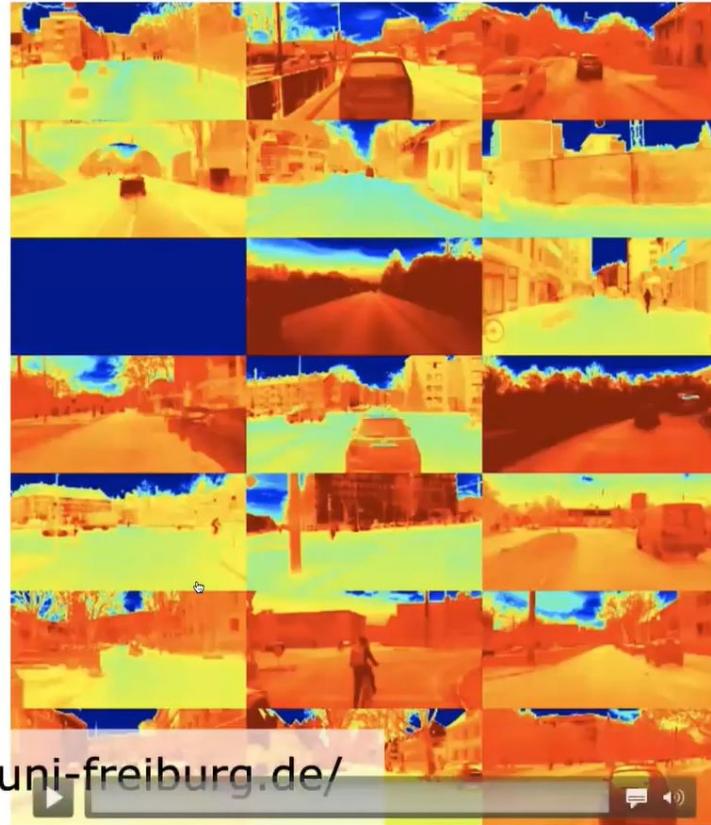
$$\mathcal{L}_{p_1} = \mathcal{L}_s^D + \lambda[0 - C(S_N)]^2, \quad \mathcal{L}_{p_2} = \frac{1}{HW} \sum_{h,w} \begin{cases} [0 - C(S_X)]^2, & \text{if } X = D \\ [1 - C(S_X)]^2, & \text{if } X = N \end{cases}$$

# Dataset

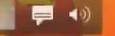


- Freiburg Thermal Dataset
- Five day- and three nighttime collections
- Multiple seasons
- 12,000 daytime images
- 8,000 nighttime images
- GPS and IMU data
- LiDAR point clouds
- 64 evaluation images

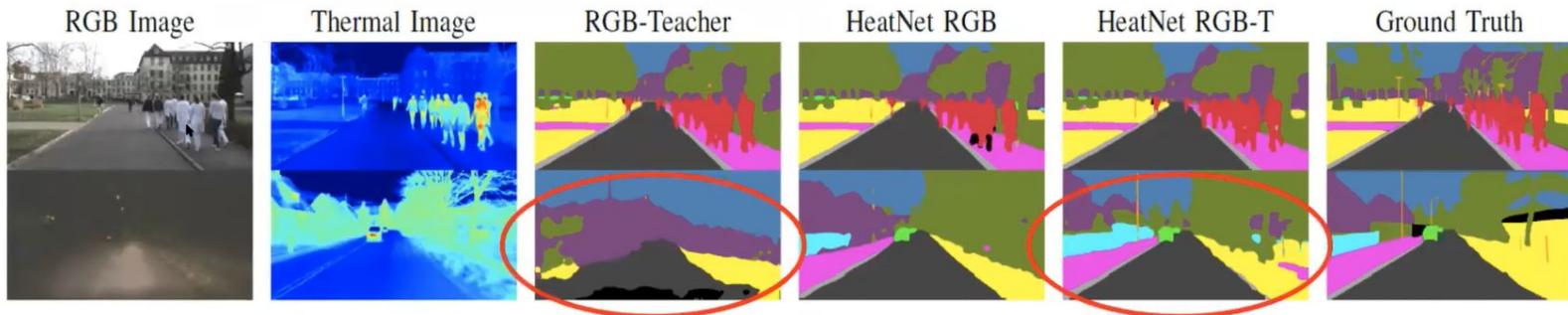
# Dataset



<http://thermal.cs.uni-freiburg.de/>



# Experimental Results



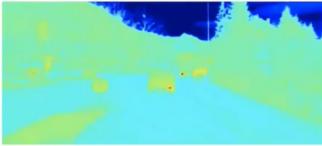
# DAYTIME AND NIGHT TIME

## Experimental Results

RGB



Thermal



HeatNet

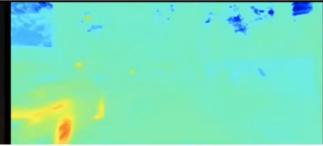


## Experimental Results

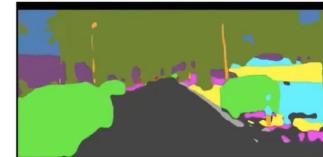
RGB



Thermal



HeatNet



# INSIGHTS

Instead of Daytime and Nighttime, we can do for without rain and in the rain.

# Quantitative Evaluation

Train On	Test On	Model	RGB	T	Road	Sidewalk	Building	Curb	Fence	Pole	Vegetation	Terrain	Sky	Person	Car	Bicycle	Mean
MF	MF	MFNet [9] RTFNet-50 [25] HeatNet	✓ ✓ ✓	✓ ✓ ✓	- - -	- - -	- - -	- - -	- - -	- - -	- - -	- - -	58.9 <b>67.8</b> 56.4	65.9 <b>86.3</b> 68.8	42.9 <b>58.2</b> 33.9	55.9 <b>70.7</b> 53.0	
	FR-T Day/Night	MFNet [9] RTFNet-50 [25] HeatNet	✓ ✓ ✓	✓ ✓ ✓	- - 86.7	57.5 67.7	46.4 41.5	41.5 43.8	57.9 44.1	43.8 63.7	57.9 63.1	42.8 63.2 <b>85.6</b>	27.0 61.5 63.1	24.5 51.3 <b>58.2</b>	31.4 58.6 59.7		
	MF	HeatNet	✓	✓	-	-	-	-	-	-	-	-	51.6	61.8	30.2	47.9	
(Vistas) FR-T	FR-T Day	RGB Teacher HeatNet	✓ ✓	✗ ✓	<b>89.7</b> 89.4	<b>67.0</b> 65.6	73.8 <b>74.8</b>	56.9 <b>59.7</b>	48.8 <b>52.9</b>	53.8 <b>54.3</b>	73.8 <b>74.1</b>	62.8 <b>65.1</b>	84.3 <b>84.5</b>	72.0 <b>74.0</b> <b>91.2</b>	90.1 <b>64.1</b>	60.4 <b>70.8</b>	
FR-T (Vistas) FR-T	FR-T Night	Thermal Teacher RGB Teacher HeatNet	✗ ✓ ✓	✓ ✗ ✓	84.9 76.3 <b>86.4</b>	60.5 22.6 <b>60.9</b>	<b>65.5</b> 53.4 <b>45.5</b>	43.1 10.8 <b>45.5</b>	31.8 14.1 <b>35.5</b>	38.1 31.6 <b>42.0</b>	51.8 10.4 <b>52.5</b>	40.1 13.5 <b>52.3</b>	72.6 47.7 <b>73.9</b>	49.6 28.0 <b>54.9</b>	87.1 74.3 85.7	<b>56.9</b> 45.2 53.3	57.0 35.7 <b>59.0</b>
FR-T FR-T	FR-T Day/Night	HeatNet HeatNet RGB-only	✓ ✓	✓ ✗	<b>87.9</b> 82.7	<b>63.3</b> 56.0	<b>70.1</b> 66.0	<b>52.6</b> 45.3	<b>44.2</b> 34.0	<b>48.2</b> 37.8	<b>63.3</b> 58.4	<b>58.9</b> 49.5	<b>79.2</b> 71.0	<b>64.5</b> 54.4	<b>88.5</b> 84.2	<b>58.7</b> 57.4	<b>64.9</b> 58.0
(Vistas) FR-T	BDD Night [34]	RGB Teacher HeatNet RGB-only	✓ ✓	✗ ✗	68.8 <b>87.1</b>	21.5 <b>40.0</b>	32.9 <b>50.2</b>	- -	0.0 <b>25.9</b>	12.3 <b>22.9</b>	11.5 <b>12.8</b>	6.6 <b>8.5</b>	<b>27.2</b> 25.0	24.5 <b>27.4</b>	40.4 <b>68.3</b>	-	24.6 <b>36.8</b>

# 2

## NIGHT IMAGES: WHAT INFORMATION WE CAN EXTRACT & ENHANCE?

# INSIGHTS

- Night Images Problem - LOW LIGHT
- **WHAT IF** we multiply the pixels with some constant values?
  - NOISY IMAGES
  - INTERESTINGLY – NOISE CAN BE DIFFERENTIATED FROM STRUCTURES

# Night Image Problem: Low Light



# INSIGHTS

- Night Images Problem: Glow/ Glare/ Floodlight.

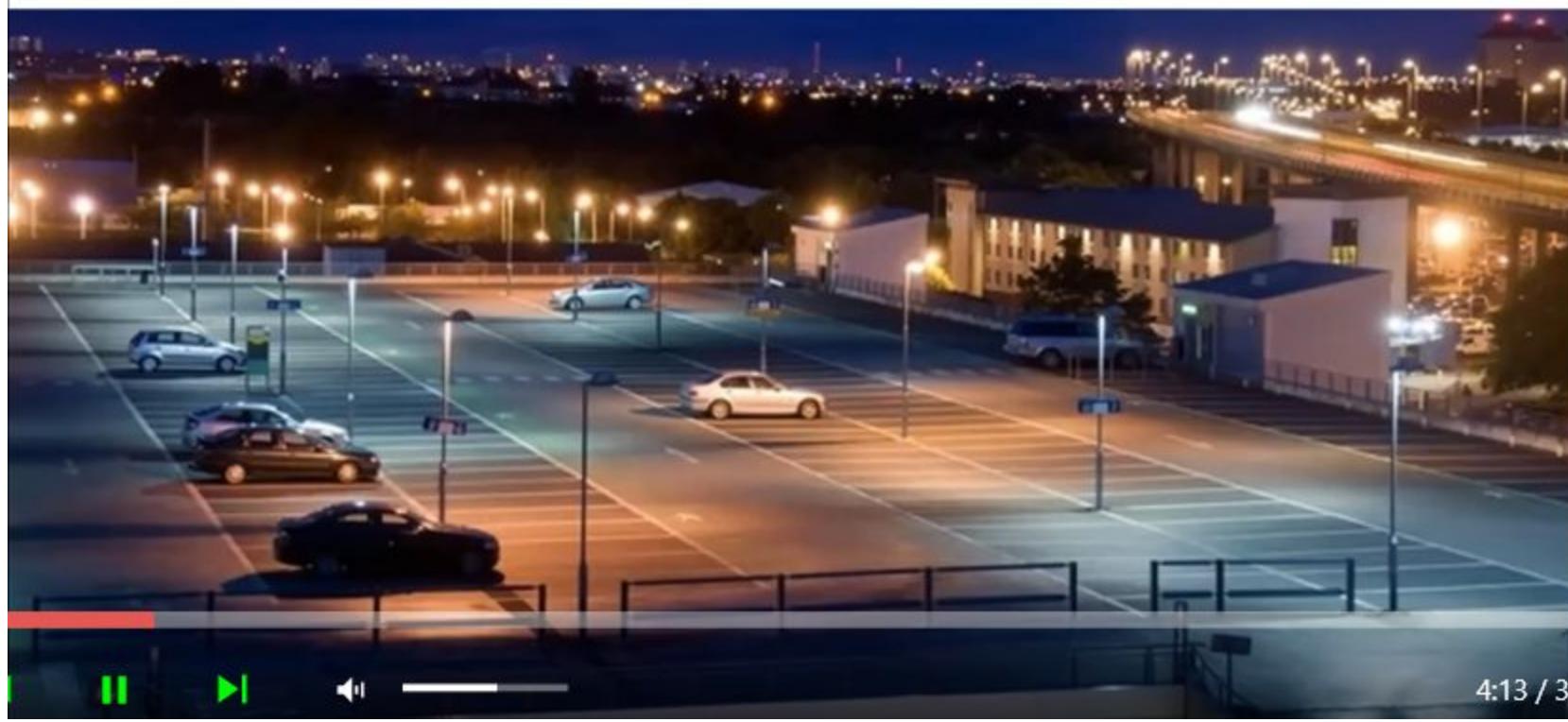
## **Night Image Problem: Glow/Glare/Floodlight**



# Night Image Problem: Uneven Light Distribution



# Night Image Problem: Light Colors



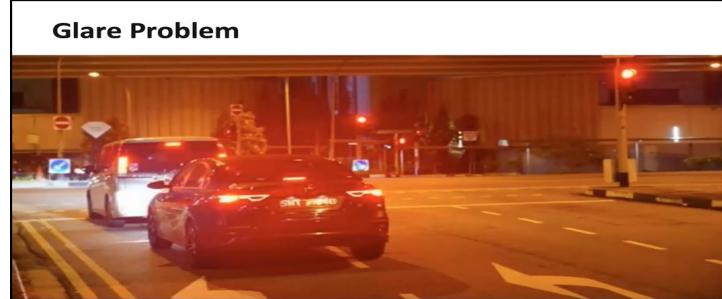


# INSIGHTS - NIGHT IMG ENHANCEMENT

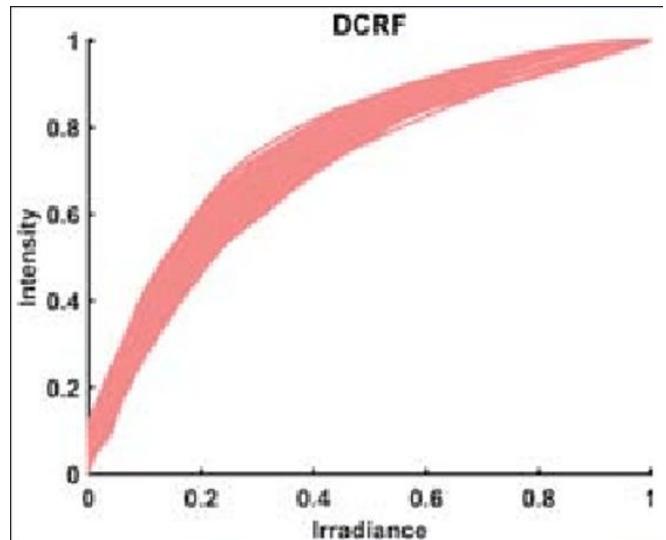
- Camera response function (CRF) relates scene irradiance to image intensities.
  - Estimation of CRF is a fundamental and necessary step in many computer vision applications such as the generation of high dynamic image, bidirectional reflectance distribution function (BRDF) estimation etc.
  - CRF is non-linear.
- Irradiance vs Radiance
  - In terms of explanation, it can be said that Radiation is the number of photons that are being emitted by a single source. Irradiation, on the other hand, is one where the radiation is falling on the surface is being calculated.
- **IDEA : BOOST the intensity, BUT have to keep the color => ESTIMATE THE CRF.**

## INSIGHTS - NIGHT IMG ENHANCEMENT

- NOT JUST BOOST THE IMAGE -
- HAVE TO SUPPRESS LIGHT IN IMGS LIKE THIS
- HAVE TO REMOVE GLARE HERE



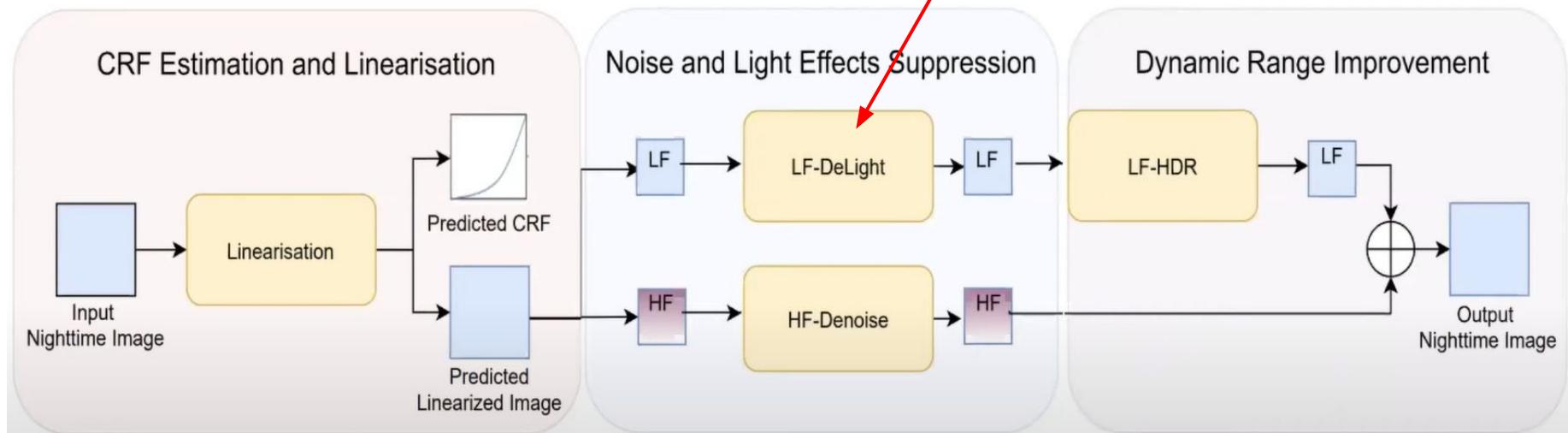
CRF of Camera



**Problem:** Boosting low light regions, at the same time suppress glare and noise

# Pipeline

De Glare



# HDR

- HDR stands for high dynamic range.
- Put simply, *it's the range of light and dark tones in your photos.*
- The human eye has a very high dynamic range — which is why we can see details in both shadows and highlights.

# Supervised Training: CRF

- CRF Loss:

$$\mathcal{L}_{\text{mse}} = \|\hat{\mathbf{g}} - \mathbf{g}^{gt}\|_2 \quad \hat{\mathbf{g}} = \mathbf{g}_0 + \sum_{i=1}^{11} \mathbf{h}_i \mathbf{c}_i$$

- Image-Linearization Loss:

$$\mathcal{L}_{\text{lin}} = \|\hat{\mathbf{g}}(\mathbf{Z}') - \mathbf{g}^{gt}(\mathbf{Z}')\|_1$$

# Supervised Training: HDR

- For decomposing the linearized image to low and high frequency layers, we employ: ‘Fast end-to-end trainable guided filter’, CVPR’18.
- HDR Loss:

$$\mathcal{L}_{\text{HDR}} = \left\| \frac{\log(1 + \mu \hat{\mathbf{X}})}{\log(1 + \mu)} - \frac{\log(1 + \mu \mathbf{X}^{gt})}{\log(1 + \mu)} \right\|_1$$

$$\hat{\mathbf{X}} = \frac{1}{K} \sum_{k=1}^K (\mathbf{HrLF}_{\mathbf{x}_k} + \mathbf{DnHF}_{\mathbf{x}_k})$$

Loss function is based on  $\mu$ -law (used for tone mapping)  
 $\mathbf{X}$  is an HDR image;  $K$  is the number of filters

# Unsupervised Test-Time Training: CRF

- CRF-Monotonicity Loss:

$$\mathcal{L}_{\text{mon}} = \sum_{t=0}^1 H \left( -\frac{\partial \hat{\mathbf{g}}(t)}{\partial t} \right)$$

- CRF-Pixel-Linearization Loss:

$$\mathcal{L}_{es} = \sum_{i=1}^S \left( \frac{|\mathbf{n}_{\mathbf{Y}_{es}}^{\min} - \mathbf{n}_{\mathbf{Y}_{es}}^{\max}| \times |\mathbf{n}_{\mathbf{Y}_{es}}^{\min} - \mathbf{n}_{\mathbf{Y}_{es}}^i|}{|\mathbf{n}_{\mathbf{Y}_{es}}^{\min} - \mathbf{n}_{\mathbf{Y}_{es}}^{\max}|} \right),$$

$$\mathcal{L}_{\text{distlin}} = \sum_{e=1}^E \left( \sum_{s=1}^S (\mathcal{L}_{es}) \right),$$

E = #patches; Each patch has a size of S x S

# Light Effect Suppression:



(a) Input



(b) LF Map



(d) LF Map (w/o l. e.)



(e) Output (w/o l. e.)

# Results:



Input



Our Method



Ground-Truth



HDRCNN [2]



SingleHDR [10]



DrTMO [3]

# Results:



Input

Our Method

LIME [6]



ZeroDCE [5]

EnlightenGAN [7]

SingleHDR [10]

# Results:

GLARE



Input Image



Our Method



LIME [6]



ZeroDCE [5]



EnlightenGAN [7]



SingleHDR [10]



## Publications:

- Nighttime Haze Removal with Glow and Multiple Light Colors ICCV'15
- Nighttime Defogging Using High-Low Frequency Decomposition and Grayscale-Color Networks, ECCV'20
- Single-Image Camera Response Function Using Prediction Consistency and Gradual Refinement, ACCV'20
- Nighttime Stereo Depth Estimation using Joint Translation-Stereo Learning: Light Effects and Uninformative Regions, 3DV'20
- Nighttime Visibility Enhancement by Increasing the Dynamic Range and Suppression of Light Effects, CVPR'21.

# NIGHTTIME HAZE REMOVAL WITH GLOW AND MULTIPLE LIGHT COLORS

- 2015 IEEE International Conference on Computer Vision.
- This paper focuses on dehazing nighttime images.

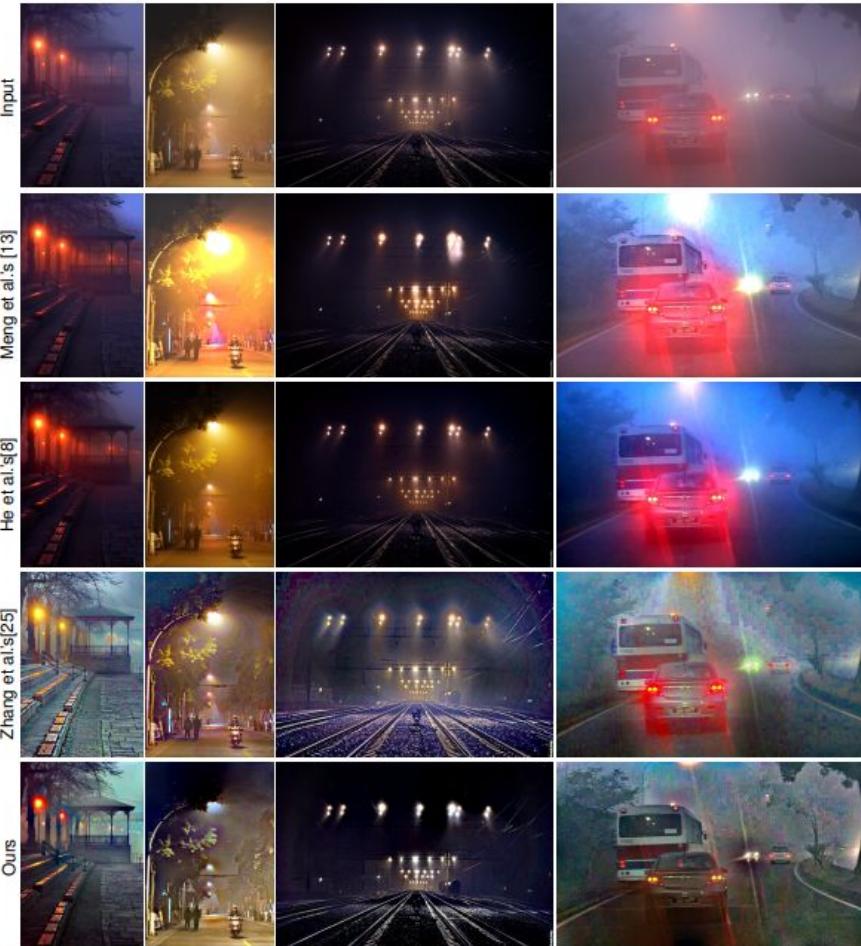
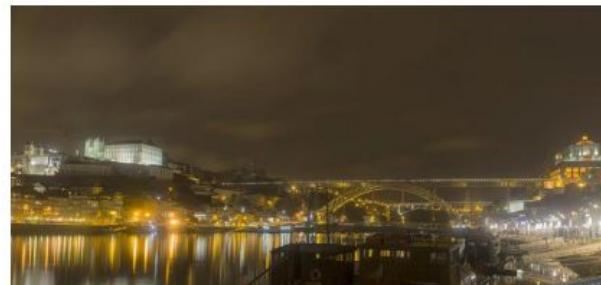


Figure 10. The qualitative comparisons of Meng et al.'s method [13], He et al.'s method [8], Zhang et al.'s method [25], and ours using various nighttime images.

# NIGHTTIME DEFOGGING USING HIGH-LOW FREQUENCY DECOMPOSITION AND GRayscale-COLOR NETWORKS

- We address the problem of nighttime defogging from a single image by introducing a framework consisting of two modules: grayscale and color modules.

European Conference on  
Computer Vision  
ECV 2020



Nighttime foggy Image



Our Result

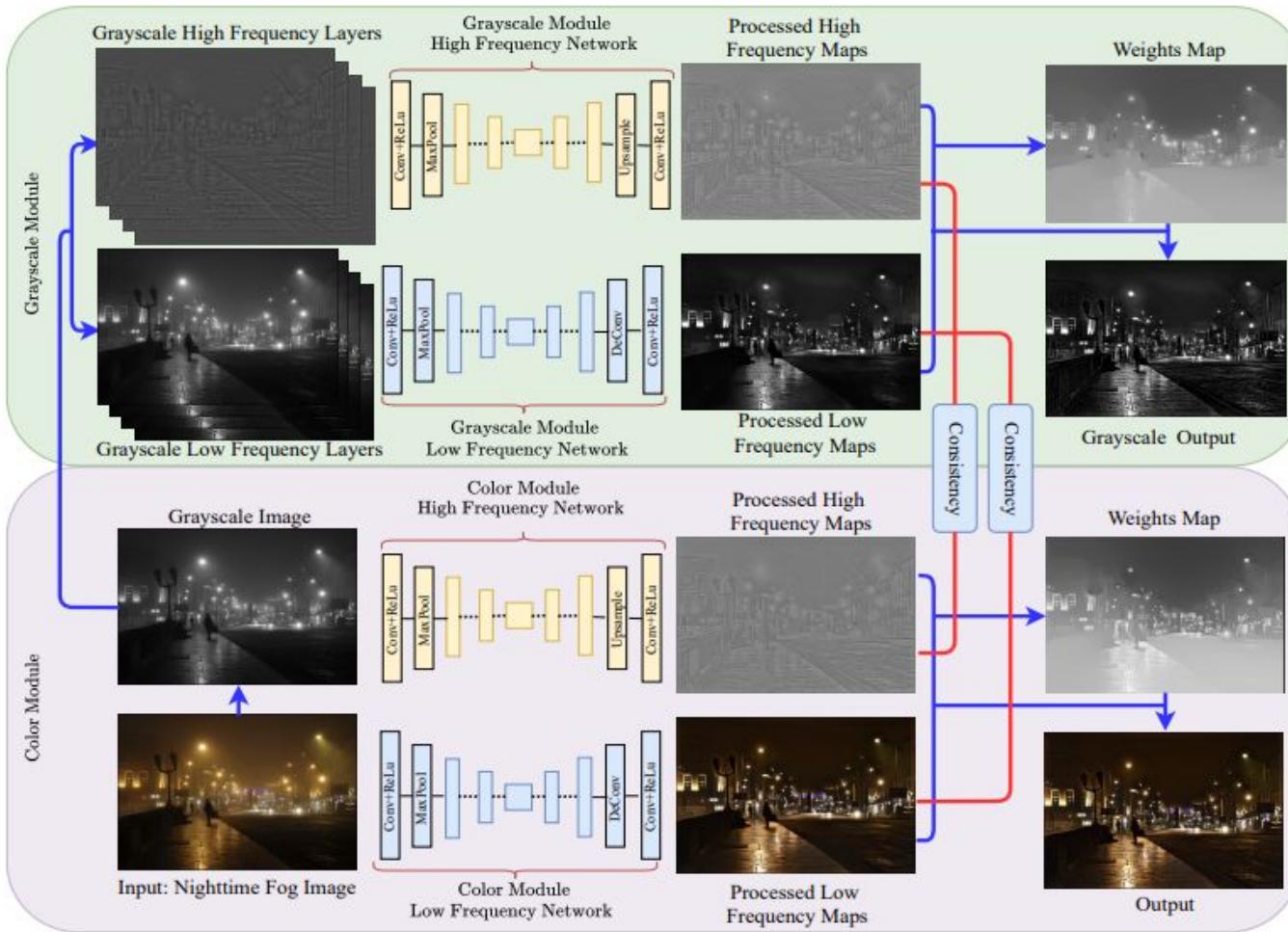


Li et al. [17]



Zhang et al. [26]

# MODEL PIPELINE





Input Image



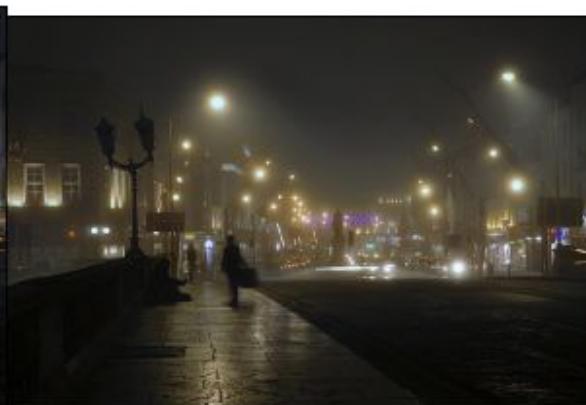
Our Result



Li et al. [17]



Zhang et al. [26]



Ancuti et al. [1]



EPDN [21]



Input Image



Our Result



Li et al. [17]



Zhang et al. [26]



Ancuti et al. [1]



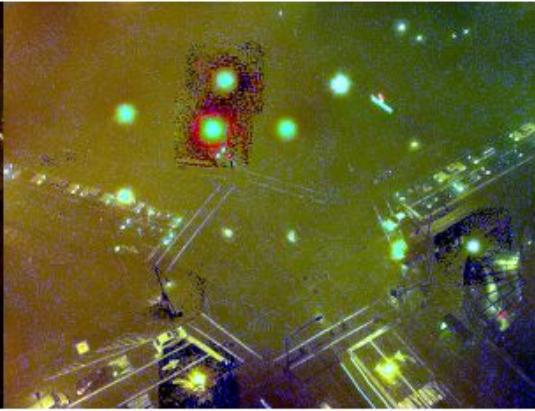
EPDN [21]



Input Image



Our Result



Li et al. [17]



Zhang et al. [26]



Ancuti et al. [1]



EPDN [21]

## 6 Conclusion

We have introduced a learning-based nighttime defogging method. To our knowledge, this is the first time, a deep learning-based method is dedicated to handle nighttime defogging problem. To achieve our goal, we design grayscale and color modules, which rely mainly on the high/low frequency layers to enhance textures and at the same time suppress glow, fog and noise. Due to the lack of paired real ground-truths, our training process employs both paired synthetic data and unpaired real data. For this, we introduce new consistency losses between the outputs of the grayscale and color modules. Experimental results and evaluations, both quantitative and qualitative, show the effectiveness of our method.

