# GDIP: Gated Differentiable Image Processing for Object-Detection in Adverse Conditions

Sanket Kalwar*[1], Dhruv Patel*[1], Aakash Aanegola[1], Krishna Reddy Konda[3], Sourav Garg[2], K Madhava Krishna[1]
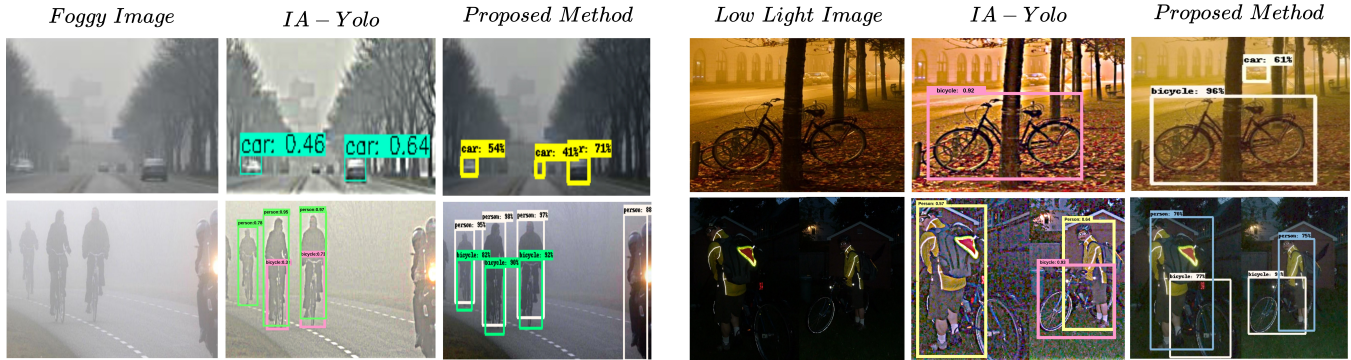
Fig. 1: **Overview**: Object detection is challenging in adverse weather conditions because objects in the scene are only partially visible, resulting in missed detections. We compare our proposed GDIP-Yolo with the current state-of-the-art (SOTA) IA-Yolo [1] qualitatively. GDIP-Yolo detects more objects than IA-Yolo (the middle car (row 1, column 2) in the top foggy image, and the car in the background (row 1, column 5) in the top low-light image), and is more confident about its prediction (best viewed in 4x zoom). The keynote is GDIP-Yolo's informed weighting of differentiable Image Processing blocks that act concurrently on the input image, leading to a vastly superior detection performance than IA-Yolo's sequential image processing framework. Note: The confidence boxes for IA-Yolo are given as a fraction between [0,1], whereas ours is in percentage.

*Abstract*— **Detecting objects under adverse weather and lighting conditions is crucial for the safe and continuous operation of an autonomous vehicle, and remains an unsolved problem. We present a Gated Differentiable Image Processing (GDIP) block, a domain-agnostic network architecture, which can be plugged into existing object detection networks (e.g., Yolo) and trained end-to-end with adverse condition images such as those captured under fog and low lighting. Our proposed GDIP block learns to enhance images directly through the downstream object detection loss. This is achieved by learning parameters of multiple image pre-processing (IP) techniques that operate *concurrently*, with their outputs combined using weights learned through a novel *gating mechanism*. We further improve GDIP through a multi-stage guidance procedure for progressive image enhancement. Finally, trading off accuracy for speed, we propose a variant of GDIP that can be used as a *regularizer* for training Yolo, which eliminates the need for GDIP-based image enhancement during inference, resulting in higher throughput and plausible real-world deployment. We demonstrate significant improvement in detection performance over several state-of-the-art methods through quantitative and qualitative studies on synthetic datasets such as PascalVOC, and real-world foggy (RTTS) and low-lighting (ExDark) datasets.**

\* denotes equal contribution

[1] are with RRC, IIIT Hyderabad, India {sankethkalwar, dhruv.r.patel14, aakash.aanegola}@gmail.com, mkrishna@iiit.ac.in

[2] is with the QUT Centre for Robotics at the Queensland University of Technology (QUT), Brisbane, Australia. s.garg@qut.edu.au

[3] is with ZF TCI, Hyderabad, India krishna.konda@zf.com

† Project page: https://gatedip.github.io

† Code: https://github.com/Gatedip/GDIP-Yolo

## I. INTRODUCTION

Autonomous mobile agents need a high-level understanding of their environment to plan their trajectories and function effectively. This requires robust perception which stems from object detection and semantic segmentation like tasks for recognizing and localizing safety-critical objects such as pedestrians and vehicles. Most object detection methods are designed for and trained with images captured under ideal environmental conditions, and often do not generalize to adverse settings (like foggy and low-light conditions). Recent attempts to achieve such robustness include domain classification based invariant detection [2]–[9], prior-knowledge based feature adaptation [10], [11], adversarially-trained image alignment [12], [13], map-specific domain adaptation [14], and physics-prior based zero-shot learning [15], [16].

More recently, learnable image pre-processing methods have emerged as a superior alternative [1], [15]–[23]. However, these methods are either limited to a single preprocessing module [15], [16], [21]–[23], require domain-specific architectural variations [1], [18], [19] or use multiple modules in an arbitrary sequential order [1], [20]. In this work, we address all these limitations and present a novel approach that learns to enhance images for object detection under adverse conditions in an end-to-end manner. This is achieved through a learnable gating-based weighted combination of concurrent image processing operations, dubbed GDIP (Gated Differen-

tiable Image Processing). Our proposed GDIP method integrated with Yolo significantly outperforms the current SOTA, Image Adaptive Yolo (IA-Yolo [1]), which relies on an arbitrary sequential image preprocessing. On real-word fog (RTTS [24]) and low-light (Ex-Dark [25]) datasets, GDIP-Yolo leads IA-Yolo in mAP by 5.76 and 15.89 respectively. The key contributions of this paper are listed below:

1) a novel *gating mechanism* that enables *concurrent* relative weighting of multiple differentiable image processing modules to enhance images for object detection under adverse environmental conditions;

2) a *multi-level* version of GDIP where an image is progressively enhanced through multiple GDIP blocks, each guided by a different layer of the image encoder; and

3) an adaptation of GDIP as a *training regularizer* which directly improves object detection training for adverse conditions, eliminating the need of GDIP during inference, thus saving compute time with a minor drop in performance.

## II. RELATED WORK

Object detection is the problem of localizing and classifying objects in the scene and has seen a recent uptick in popularity due to its applications in autonomous vehicles and more. There are two primary approaches to the object detection problem, two-stage detection, and single-stage detection. Two-stage detectors like FasterRCNN [26] and MaskRCNN [27] utilize a region proposal network (RPN) which generates proposals of plausible regions of interest, sent to the downstream network that performs classification. Two-stage approaches are computationally expensive, reducing their application range. Single-stage detectors like Yolo [28], RetinaNet [29], SSD [30] and FCOS [31] bypass the heavy RPN and directly extract objects with their associated labels. Nevertheless, under adverse conditions, both types of networks fail to detect objects.

**Adverse Conditions:** Typical object detection techniques fail in adverse weather conditions, and transfer learning has proven to be a viable way to employ object detection in adverse weather conditions. Chen et al. [2] approach this problem from a domain adaptation perspective and utilize image and instance-level features to reduce the domain shift. Singadi et al. [10] use weather-specific knowledge and define a prior-adversarial loss with feature recovery to mitigate weather effects on detection. Multiscale Domain Adaptive Yolo (MS-DAYOLO) [32] uses classifiers for each domain at different scales to learn domain invariant features. Zhang et al. [12] employ image-level feature alignment to match local and global features.

**Differentiable Image Pre-processing:** Another popular approach to the problem is to perform image enhancement before object detection. In Exposure [20], a deep Reinforcement Learning model learns a policy to apply a sequence of enhancement operations. AOD-Net [33] dehazes images using a CNN designed on a re-formulated atmospheric scattering model. Dong et al. [34] use an encoder-decoder architecture (U-Net) with the strength-operate-subtract boosting strategy to help dehaze images. GridDehazeNet [18] employs a multi-scale attention mechanism with pre and post-processing modules to generate better inputs and reduce artifacts in the final dehazed image. He et al. [35] use a dark channel prior (one color channel in most pixels will be low) to dehaze images but do not perform object detection. Some models target only a single adverse condition, Guo et al. [19] perform light enhancement by estimating light enhancement curves that are applied iteratively to the input image, boosting face detection performance. Zeng et al. [36] learn multiple look-up tables and use CNN predictions to fuse the look-up tables into one and transform the color and tone of a source image. DSNet [17] uses two subnets for image restoration and object detection to boost performance in adverse weather conditions. IA-Yolo [1] uses a CNN to predict differentiable image processing parameters trained in conjunction with Yolov3 [28] for object detection to enhance images, and perform object detection in an end-to-end fashion.

Unlike existing methods, GDIP is a domain agnostic network architecture that handles multiple image processing operations concurrently. Additionally, it has a unique advantage with its utility as a training regularizer, which eliminates image enhancement overhead during inference resulting in higher throughput.

## III. PROPOSED METHOD

We propose a Gated Differentiable Image Processing (GDIP) framework that learns to enhance input images for object detection in adverse environmental conditions. The GDIP block learns parameters for multiple Image Processing (IP) operations performed concurrently and learns the optimal weights to combine their output. We use the following IP operations (similar to IA-Yolo [1]): tone correction ($T$), contrast balance ($C$), sharpening ($S$), defogging ($DF$), gamma correction ($G$), white balancing ($WB$), and the identity operation ($I$). Unlike IA-Yolo's sequential image enhancement, GDIP enhances images through a weighted combination of concurrent IP operations.

### A. *Gated Differentiable Image Processing (GDIP) block:*

The GDIP block (shown in Fig. 2) consists of multiple gated image processing modules, referred to as $Gb_{IP}$, that individually enhance images, which are then combined through the weights predicted by the gates. Each $Gb$ module contains a linear layer, a differentiable image processing operation, a gate (shifted *tanh* function that returns a value between 0 and 1), and a normalization operation. The linear layer (purple linear block in Fig. 2) computes two entities: the parameters required by the differentiable IP block and a scalar value that serves as an input to its corresponding gate. The individual linear layers of every $Gb$ module are passed a common feature embedding as input, obtained from a shared vision encoder (described later). The output of the IP operation (using the predicted parameters) gets multiplied by the scalar output of the gate. The weighted outputs of individual $Gb$
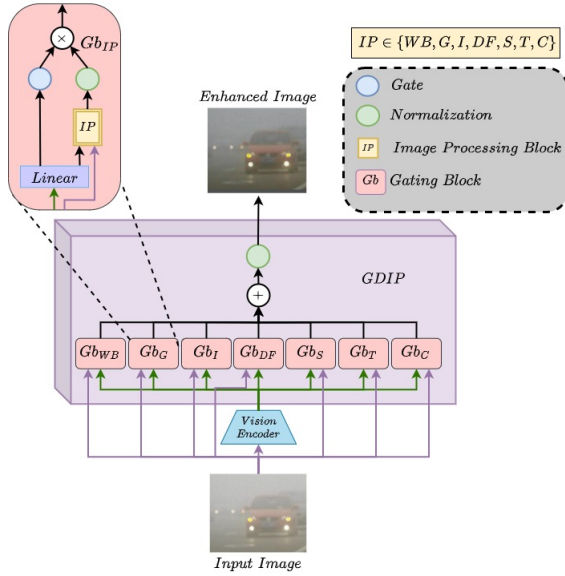
Fig. 2: **GDIP Block:** The peach blocks represent different $Gb_{IP}$ and their structures can be seen in the expanded general $Gb$ block (top left corner). Yolo and the object detection output are not shown due to space constraints.

TABLE I: Image Processing (IP) Operations. The parameters in the second column are computed by the linear layer in GDIP. Further details on the equations can be found in IA-Yolo [1]

| IP Operation | Parameters | Function |
|---|---|---|
| Tone | $t_i$ | $I_{tone} = (L_{t_r}(r_i), L_{t_g}(g_i), L_{t_b}(b_i))$ |
| Contrast | $\alpha$ | $I_{contrast} = \alpha * En(I) + (1 - \alpha) * I$ |
| Sharpening | $\lambda$ | $I_{sharpen} = I + \lambda * (I - Gaussian(I))$ |
| Defogging | $\omega$ | $I = I_{defog} * t(\omega) + A * (1 - t(\omega))$ |
| Gamma | $\gamma$ | $I_{gamma} = I^{\gamma}$ |
| White balance | $W_r, W_g, W_b$ | $I_{wb} = (W_r r_i, W_g g_i, W_b b_i)$ |
| Identity | - | $I_{identity} = I$ |

blocks are finally aggregated to obtain an enhanced image. Expressed mathematically, the output of the GDIP block is:

$$z = N(\sum_i N(f_i(x)) * w_i) \tag{1}$$

where $x$ is the input image captured under adverse environmental conditions, $z$ is the enhanced clear image, $f_i(x)$ represents the $i^{th}$ IP operation (top-right in Fig. 2) weighted by its respective scalar gate output $w_i \in [0, 1]$, and $N$ is the min-max normalization operation. Normalization ensures that the pixel intensity range of the output of all the image processing operations are the same. The IP operations are expressed mathematically in Table I, see [1] for a detailed description.

**Vision Encoder:** Our proposed GDIP block requires latent embeddings to compute image processing parameters and gate values. For this purpose, we employ a vision encoder comprising five convolutional layers (each with a kernel size of three and a stride of one). The number of channels in each layer is double the previous, starting from 64 in the first layer and 1024 in the final layer. Each convolution operation is followed by average pooling (with kernel size
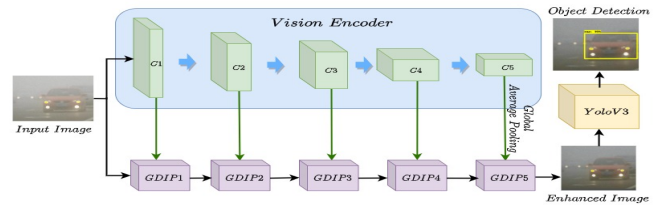


Fig. 3: **MGDIP-Yolo:** The GDIP blocks perform enhancements progressively on the input image starting with $GDIP_1$ obtaining intermediate features from $C1$ (each GDIP block still performs IP operations concurrently). The final enhanced image is passed to Yolo and the whole pipeline is trained end-to-end.

three and stride two), while the last layer is followed by global average pooling, the output of which is a 1x1x1024. This is then projected to a 256-dimensional latent space using a fully connected layer. The GDIP block takes this 256-dimensional embedding from the vision encoder along with the adverse input image and performs image enhancement after computing the necessary parameters.

**GDIP-Yolo:** To integrate GDIP with Yolo, we use the vision encoder with GDIP to perform image enhancements (depicted in Fig. 2), and use the enhanced image as input to Yolo. Integrating GDIP with Yolo in this fashion ensures that our architecture doesn't require any additional loss formulation and uses Yolo's standard object detection loss [37] (referred to as $L_{obj}$) to train the network for object detection end-to-end.

### B. *Multi-Level GDIP (MGDIP):*

GDIP-Yolo contains a single GDIP block, which is fed with latent embeddings obtained from the vision encoder. Since, we only use the last layer of the vision encoder for this purpose, it limits the extent of information available for GDIP to learn parameters for image processing modules. Thus, we propose multi-level progressive image enhancement, achieved by integrating a GDIP block with every layer of the vision encoder, dubbed MGDIP-Yolo. Note that the individual image processing modules within a single GDIP block still operate concurrently with their corresponding gates providing relative weightings. As shown in Fig. 3, MGDIP progressively enhances images by feeding the output from one GDIP block as input to the next, where individual GDIP blocks are guided by the features extracted from different layers of the vision encoder. The final enhanced output from MGDIP is passed to Yolo for object detection. MGDIP-Yolo is trained in an end-to-end manner using the standard object detection loss $L_{obj}$, similar to GDIP-Yolo.

We hypothesize that utilizing embeddings from different layers provides GDIP access to multiple feature scales, each of which can have a varied relevance for different image processing operations. This is based on the understanding that earlier layers in CNNs capture lower level information (local information like edges) and later layers capture high-level (global) information. Thus, MGDIP gains the ability to use the local/global feature properties to selectively apply image processing operations.
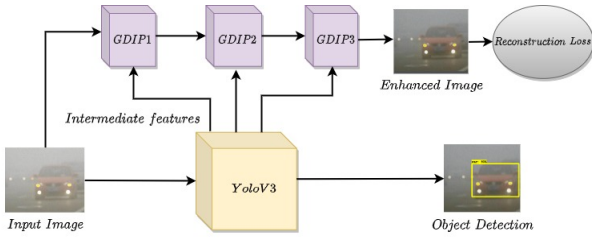
Fig. 4: **GDIP as a regularizer:** The purple GDIP blocks and the reconstruction loss help regularize Yolo's features. During inference, everything apart from Yolo is removed.

### C. *GDIP block as a regularizer:*

In this section, we demonstrate how GDIP can also be employed as a feature regularization technique to improve Yolo's performance while maintaining its throughput.

The original GDIP block used a vision encoder to obtain feature embeddings. Alternatively, multiple GDIP blocks can be connected to intermediate layers of Yolo, bypassing the need for a vision encoder and directly using Yolo's embeddings to construct an enhanced output, as shown in Fig. 4. Note that this enhanced output is not the input to Yolo but rather a byproduct that we use for training regularization. The reconstruction loss (Eq. 2) is calculated between this output and the clear version of the input image as a combination of $L_1$ norm and Mean Square Error loss $L_{MSE}$. The overall loss function used is shown in Eq. 3, where $\alpha$ is the weight of the reconstruction loss and is empirically set to $1 \times 10^{-4}$.

$$L_{Reg} = L_1 + L_{MSE} \quad (2) \qquad L_{total} = L_{obj} + \alpha L_{Reg} \quad (3)$$

Inclusion of the reconstruction loss in the formulation helps Yolo learn features that are invariant to adverse conditions, resulting in better performance when compared to standalone Yolo. Since the GDIP blocks exist solely to refine Yolo's features, it is only required during training and can be removed during inference. This results in an unchanged network architecture (Yolo) that performs better in adverse weather conditions along with higher throughput.

### IV. EXPERIMENTAL SETUP

#### A. *Datasets*

**Foggy Conditions:** We use the RTTS dataset [24], a collection of 4322 natural foggy images with five annotated classes - person, car, bus, bicycle and motorcycle - primarily for testing. The PascalVOC train/val datasets (2007 and 2012) [38], [39] have 22136 clear images that we use as a base to create a synthetic training set. We select images from PascalVOC having objects belonging to the five classes from RTTS and create two datasets, one with clear images (VOCNormal) and one with augmented foggy images generated using the atmospheric scattering model (ASM) [40]. We employ the ASM to generate 10 different levels of fog to include variance in our synthetic training set. We subsample and prepare a synthetic testing set in a similar fashion from the PascalVOC 2007 test set with 4952 images (referred to as V_F_Ts). We employ a *hybrid* strategy where we use a

mix of foggy and clear images (in a 2:1 ratio) to help our model learn fog-invariant features.

**Low-lighting Conditions:** The ExDark dataset [25] is a collection of 7363 real-world images with 10 object classes in low lighting conditions that we use to evaluate our models. Similar to preparing the foggy dataset, we select images from PascalVOC having objects from the 10 classes of ExDark and apply a gamma filter to emulate a low-lighting condition. Mathematically, $I_{dark} = I^\gamma$, where $\gamma$ is sampled uniformly from the range of 1.5 to 5, $I$ is the normalized clear image, and $I_{dark}$ is the synthetic dark image. Using the same selection and image processing methods, we generate a synthetic low-light test set (V_D_Ts) from the PascalVOC test set. During training, we employ a *hybrid* strategy (similar to the foggy setting) by using a mix of dark and clear images.

#### B. *Training Setup*

Training for both foggy and low-lighting setting is done by resizing images to $448 \times 448 \times 3$ pixels and with a batch size of 6 for 80 epochs. We use a cosine learning rate scheduler with learning rates ranging from $1 \times 10^{-6}$ to $1 \times 10^{-4}$ and an SGD optimizer with a weight decay of $5 \times 10^{-4}$.

### V. RESULTS AND ANALYSES

We provide qualitative and quantitative results that establish our proposed method's superiority and evaluate design choices through ablation studies. We also show that GDIP variants provide flexibility to prioritize speed or accuracy based on application requirements.

#### A. *Qualitative Analysis*

We compare the results of GDIP-Yolo with the current SOTA IA-Yolo, shown in Fig. 1 on real-world data. Unlike IA-Yolo, our method clears fog and improves lighting conditions without changing underlying color distributions. Our method is able to detect far-off objects such as cars and bicycles (see the third column in Fig. 1), which are generally missed in extreme foggy conditions by SOTA IA-Yolo (see the second column). The last column (low-light conditions) clearly indicates GDIP-Yolo enhances lighting without artifacts, helping the model detect all objects in the scene (the car in the background, for example). We quantify the improvement extended by GDIP in the next section.

#### B. *Quantitative Analysis*

We compare our proposed method with other SOTA works using the standard object detection evaluation metric - mean average precision (mAP). All mAP values are calculated at an IoU (Intersection over Union) of 0.5.

*a) Foggy Conditions:* In Table II, we compare our proposed variants of the GDIP with other competing methods on VOCNormal Test set (V_N_Ts), synthetic VOCFoggy Test set (V_F_Ts), and the real-world foggy dataset RTTS. The second column in the table shows the training data used by each of the methods, where "Hybrid" implies the use of both clear and foggy data. We set YoloV3 as the baseline, which is trained on a mix of foggy and clear

TABLE II: Quantitative results for *foggy* conditions on the V_N_Ts (VOCNormal Test set), V_F_Ts (VOCFoggy Test set) and real-world RTTS dataset. Best and second best mAP scores are bold and italicized, respectively.

| | Methods | Train Data | V_N_Ts (mAP) | V_F_Ts (mAP) | RTTS (mAP) |
|---|---|---|---|---|---|
| Baseline | Yolov3 [28] | Hybrid | 64.13 | 63.40 | 30.80 |
| Defog | MSBDN [34] | VOC_Norm | - | 57.38 | 30.20 |
| | GridDehaze [18] | VOC_Norm | - | 58.23 | 31.42 |
| Domain Adaptation | DAYolo [32] | Hybrid | 56.51 | 55.11 | 29.93 |
| Multi-task | DSNet [17] | Hybrid | 53.29 | 67.40 | 28.91 |
| Image Adaptive | IA-Yolo [1] | Hybrid | 73.23 | 72.03 | 37.08 |
| Proposed Method | GDIP-Yolo | Hybrid | *73.70* | 71.92 | *42.42* |
| | MGDIP-Yolo | Hybrid | **75.36** | **73.37** | **42.84** |
| | GDIP Regularizer | Hybrid | 73.17 | 72.77 | 39.52 |

TABLE III: Quantitative results for *low-lighting* conditions on the V_N_Ts (VOCNormal Test set), V_D_Ts (VOCDark Test set) and real-world ExDark dataset. Best and second best mAP scores are bold and italicized, respectively.

| | Methods | Train Data | V_N_Ts (mAP) | V_D_Ts (mAP) | ExDark (mAP) |
|---|---|---|---|---|---|
| Baseline | Yolov3 [28] | Hybrid | 62.73 | 52.28 | 37.03 |
| Enhance | ZeroDCE [19] | VOC_Norm | - | 33.57 | 34.41 |
| Domain Adaptation | DAYolo [32] | Hybrid | 41.68 | 21.53 | 18.15 |
| Multi-task | DSNet [17] | Hybrid | **64.08** | 43.75 | 36.97 |
| Image Adaptive | IA-Yolo [1] | Hybrid | 56.01 | 48.44 | 26.67 |
| Proposed Method | GDIP-Yolo | Hybrid | *63.23* | *57.85* | **42.56** |
| | MGDIP-Yolo | Hybrid | 62.86 | **57.91** | *40.96* |
| | GDIP Regularizer | Hybrid | 62.30 | 57.67 | 40.72 |

images to validate if we can improve performance by using data augmentation. We also compare our results against a diverse range of methods based on domain adaption (DA-Yolo [32]), multi-task learning (DSNet [17]), defogging as pre-processing (MSBDN [34], GridDehaze [18]), and adaptive image enhancement (IA-Yolo [1]). It can be observed in Table II that our proposed variants of GDIP establish a new SOTA across different fog datasets.

All GDIP variants perform significantly better than SOTA methods on RTTS, which tests the generalizability of our method to real-world conditions. Our basic GDIP-Yolo variant outperforms the SOTA method IA-Yolo by 5.34 mAP. This can be attributed to the concurrently weighted IP operations unlike IA-Yolo's fixed sequential pipeline. MGDIP-Yolo further improves upon the GDIP-Yolo by 0.42 mAP and does so consistently across all datasets. It emerges superior to all other methods and GDIP variants, as it benefits from multi-scale information. Our regularizer variant outperforms all SOTA methods, while being as fast as vanilla Yolo (see Subsection V-E), emerging as an alternative to GDIP-Yolo and MGDIP-Yolo with an accuracy-speed trade-off.

*b) Low-lighting conditions:* We compare our proposed variants with other SOTA methods on the real-world ExDark dataset, synthetic low-lighting VOCDark test set (V_D_Ts), and VOCNormal test set (V_N_Ts), as shown in Table III. Here once again, we set YoloV3 as the baseline trained on hybrid data of a mix of dark and clear images. In addition, we also compare against a diverse range of methods based on light enhancement as pre-processing (ZeroDCE [19]), domain adaptation (DA-Yolo [32]), multi-task learning (DSNet [17]) and adaptive image enhancement (IA-Yolo [1]).

Our proposed GDIP-Yolo outperforms all the existing methods on the real-world ExDark dataset and achieves an absolute increase of 16 mAP over the previous SOTA IA-Yolo. Additionally, MGDIP-Yolo and GDIP as a regularizer variants also perform superior to other methods in this setting. In the synthetic low-lighting setting, MGDIP-

Yolo approach emerges superior, while the other GDIP variants also significantly improve the performance over other methods. On the VOCNormal test set, our proposed method performs comparable to DSNet [17]. To conclude, our proposed variants are significantly better than existing approaches for low-light settings (synthetic and real-world).

*C. Detection Statistics*

We present True and False Positives (TP, FP) and False Negatives (FN) of the number of object detections as an interpretable statistical measure as mAP does not convey the actual detections. The TP (Fig. 5 left) and FN (Fig. 5 middle) plots show substantial improvement at high object detection confidence thresholds both for RTTS and Ex-Dark datasets for GDIP-Yolo vis a vis SOTA IA-Yolo. For Autonomous Driving applications, the TP and FN statistics are critical, as not detecting an object when present can be catastrophic, and on these vital statistics, the significantly superior performance of GDIP is evident. GDIP-Yolo evaluates to comparable FP metrics vis a vis SOTA at high confidence
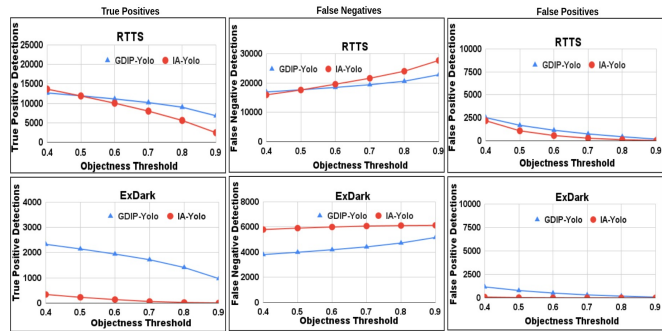


Fig. 5: **Alternative Metrics:** We indicate the True Positives (TP), False Negatives (FN) and False Positives (FP) on the RTTS and ExDark datasets for GDIP-Yolo (blue) and IA-Yolo (red). The x and y axes in the plots represent the objectness (detection confidence) threshold and the absolute number of TP/FN/FP. Detection confidence reflects the likelihood of box containing the object.
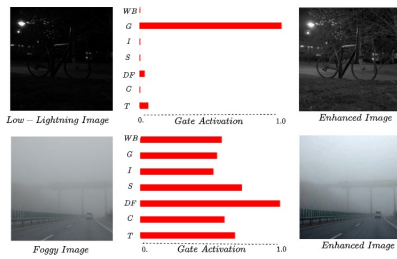
Fig. 6: **GDIP gate firing pattern:** The bar charts show the gates firing for low lighting and foggy images. GDIP learns the optimum enhancements to apply based on the input image features.

thresholds on which it shows vastly superior performance on the TP, and FN metrics.

### D. Ablation Studies

In this section, we demonstrate experimental support for using the gating mechanism and normalization layer in our proposed GDIP block. We validate that incorporating gating and normalization in the GDIP block provides a stable improvement in the downstream detection task for both foggy and low-light conditions (as shown in Table IV).

**Single Best vs Weighted Combination:** Our proposed gating mechanism helps combine image processing operations through relative weighting (see Eq. 1). To illustrate its effect, we remove this mechanism and use a single best image processing operation based on the highest gate value (referred to as GDIP-max in Table IV), expressed as $z = N(f_{i^*}(x))$ where $i^* = \arg\max_i w_i$. Without the proposed gating, performance reduces by 10.4 mAP for the RTTS dataset (comparing row 1 and 2 in Table IV). For the ExDark dataset the performance increases by a negligible amount - 0.15 mAP, insubstantial compared to the drop observed in the foggy setting. This study indicates that incorporating enhancements from multiple image processing operations is necessary as no single operation is sufficient for dealing with adverse conditions.

**Uniform vs Predicted Weighting:** In this experiment, we compare our proposed relative weighting of image processing operations with a uniform weighting across all operations (referred to as GDIP w/o gates in Table IV), expressed mathematically as $z = N(\sum_i N(f_i(x)))$. We observe that mAP reduces by 0.77 and 0.27 for RTTS and ExDark, respectively. This performance drop can be attributed to the fact that depending on the environmental conditions of images, image processing operations need to be weighted differently, achieved through gating. This is evident from Fig. 6, where activation of the gamma gate (G) is prominent for low light conditions, and in the case of foggy conditions, most of the gates remain activated with defogging (DF) being the highest. The gate activation patterns also help improve the interpretability of GDIP by indicating which IP operations are performed by what proportions based on the input image.

**With and Without Normalization:** We also verify the necessity for the normalization layer after each image processing operation by removing them (GDIP-unnormalized in

TABLE IV: Ablation Study

| Method | RTTS (mAP) | ExDark (mAP) |
|---|---|---|
| GDIP | **42.42** | 42.56 |
| GDIP-max | 31.99 | **42.71** |
| GDIP-unnormalized | 40.61 | 40.2 |
| GDIP w/o gates | 41.65 | 42.29 |

TABLE V: Real-time performance on GeForce GTX 1080

| Methods | FPS |
|---|---|
| Yolo V3 | 68.39 ± 1.5 |
| IA-Yolo | 22.84 ± 0.0 |
| GDIP-Yolo | 29.78 ± 0.1 |
| MGDIP-Yolo (top-down) | 11.38 ± 0.02 |
| MGDIP-Yolo (bottom-up) | 11.25 ± 0.03 |
| GDIP as regularizer | **68.39 ± 1.5** |

Table IV), expressed as $z = N(\sum_i f_i(x) * w_i)$. This leads to a considerable performance drop of around 1.8 and 2.35 mAP in RTTS and ExDark, respectively.

Overall, these ablation studies indicate that GDIP with normalization and the gating mechanism leads to the best overall performance irrespective of environmental conditions and is a promising solution for the object detection task.

### E. Real-time performance

In Table V, we compare the real-time performance of our proposed GDIP variants with other techniques. GDIP peforms the fastest as a regularizer at around 68 fps on a Nvidia GTX 1080Ti, which is the same as YoloV3. Our basic variant - GDIP-Yolo operates at 7 fps higher than IA-Yolo, while achieving SOTA mAP on real-world fog and night datasets.

## VI. CONCLUSION

We presented GDIP and MGDIP as domain-agnostic network architectures for object detection in adverse weather conditions, which can be used with existing object detection networks and trained under different adverse conditions, as we demonstrated for fog and low lighting. We also presented a training regularizer variant of GDIP, which improves the baseline Yolo performance under adverse conditions while maintaining its original throughput. All our GDIP variants result in a new state-of-the-art on challenging real-world datasets both under foggy and low-lighting conditions, while only having trained on synthetic adverse condition data, thus exhibiting significant generalization capability.

In future, this work can be extended to other adverse condition types (e.g., haze, rain, snow, etc.) along with additional relevant image pre-processing operations which are easy to integrate given their concurrent processing and relative weighting within GDIP. For long-term autonomy and highly-safe operations, dealing with adverse conditions is crucial for autonomous vehicles, and this work pushes the boundaries of robust perception, getting a step closer to the ubiquity of autonomous vehicles.

## References

[1] W. Liu, G. Ren, R. Yu, S. Guo, J. Zhu, and L. Zhang, "Image-adaptive yolo for object detection in adverse weather conditions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022. 1, 2, 3, 5

[2] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool, "Domain adaptive faster r-cnn for object detection in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3339–3348. 1, 2

[3] H.-K. Hsu, C.-H. Yao, Y.-H. Tsai, W.-C. Hung, H.-Y. Tseng, M. Singh, and M.-H. Yang, "Progressive domain adaptation for object detection," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, March 2020. 1

[4] K. Tian, C. Zhang, Y. Wang, S. Xiang, and C. Pan, "Knowledge mining and transferring for domain adaptive object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 9133–9142. 1

[5] A. Wu, R. Liu, Y. Han, L. Zhu, and Y. Yang, "Vector-decomposed disentanglement for domain-invariant object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9342–9351. 1

[6] A. Wu, Y. Han, L. Zhu, and Y. Yang, "Instance-invariant domain adaptive object detection via progressive disentanglement," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 1

[7] C. Lin, Z. Yuan, S. Zhao, P. Sun, C. Wang, and J. Cai, "Domain-invariant disentangled network for generalizable object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 8771–8780. 1

[8] Q. Gu, Q. Zhou, M. Xu, Z. Feng, G. Cheng, X. Lu, J. Shi, and L. Ma, "Pit: Position-invariant transform for cross-fov domain adaptation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 8761–8770. 1

[9] Z. He and L. Zhang, "Domain adaptive object detection via asymmetric tri-way faster-rcnn," in *European conference on computer vision*. Springer, 2020, pp. 309–324. 1

[10] V. A. Sindagi, P. Oza, R. Yasarla, and V. M. Patel, "Prior-based domain adaptive object detection for hazy and rainy conditions," in *European Conference on Computer Vision*. Springer, 2020, pp. 763–780. 1, 2

[11] Z. Zhang, L. Zhao, Y. Liu, S. Zhang, and J. Yang, "Unified density-aware image dehazing and object detection in real-world hazy scenes," in *Proceedings of the Asian Conference on Computer Vision*, 2020. 1

[12] S. Zhang, H. Tuo, J. Hu, and Z. Jing, "Domain adaptive yolo for one-stage cross-domain detection," in *Asian Conference on Machine Learning*. PMLR, 2021, pp. 785–797. 1, 2

[13] C. Zhuang, X. Han, W. Huang, and M. Scott, "ifan: Image-instance full alignment networks for adaptive object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 122–13 129. 1

[14] C. Sakaridis, D. Dai, and L. Van Gool, "Map-guided curriculum domain adaptation and uncertainty-aware evaluation for semantic nighttime image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 1

[15] A. Lengyel, S. Garg, M. Milford, and J. C. van Gemert, "Zero-shot day-night domain adaptation with a physics prior," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4399–4409. 1

[16] S. Zheng and G. Gupta, "Semantic-guided zero-shot learning for low-light image/video enhancement," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*, January 2022, pp. 581–590. 1

[17] S.-C. Huang, T.-H. Le, and D.-W. Jaw, "Dsnet: Joint semantic learning for object detection in inclement weather conditions," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 8, pp. 2623–2633, 2020. 1, 2, 5

[18] X. Liu, Y. Ma, Z. Shi, and J. Chen, "Griddehazenet: Attention-based multi-scale network for image dehazing," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 7314–7323. 1, 2, 5

[19] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-reference deep curve estimation for low-light image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1780–1789. 1, 2, 5

[20] Y. Hu, H. He, C. Xu, B. Wang, and S. Lin, "Exposure: A white-box photo post-processing framework," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 2, pp. 1–17, 2018. 1, 2

[21] W. Yang, S. Wang, Y. Fang, Y. Wang, and J. Liu, "From fidelity to perceptual quality: A semi-supervised approach for low-light image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1

[22] G. Li, Y. Yang, X. Qu, D. Cao, and K. Li, "A deep learning based image enhancement approach for autonomous driving at night," *Knowledge-Based Systems*, vol. 213, p. 106617, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0950705120307462 1

[23] Z. Zhang, L. Zhao, Y. Liu, S. Zhang, and J. Yang, "Unified density-aware image dehazing and object detection in real-world hazy scenes," in *Proceedings of the Asian Conference on Computer Vision (ACCV)*, November 2020. 1

[24] B. Li, W. Ren, D. Fu, D. Tao, D. Feng, W. Zeng, and Z. Wang, "Benchmarking single image dehazing and beyond," *IEEE Transactions on Image Processing*, vol. PP, pp. 1–1, 08 2018. 2, 4

[25] Y. P. Loh and C. S. Chan, "Getting to know low-light images with the exclusively dark dataset," *Computer Vision and Image Understanding*, vol. 178, pp. 30–42, 2019. 2, 4

[26] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015. 2

[27] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969. 2

[28] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018. 2, 5

[29] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988. 2

[30] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37. 2

[31] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9627–9636. 2

[32] M. Hnewa and H. Radha, "Multiscale domain adaptive yolo for cross-domain object detection," in *2021 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2021, pp. 3323–3327. 2, 5

[33] B. Li, X. Peng, Z. Wang, J. Xu, and D. Feng, "Aod-net: All-in-one dehazing network," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4770–4778. 2

[34] H. Dong, J. Pan, L. Xiang, Z. Hu, X. Zhang, F. Wang, and M.-H. Yang, "Multi-scale boosted dehazing network with dense feature fusion," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2157–2167. 2, 5

[35] K. He, J. Sun, and X. Tang, "Single image haze removal using dark channel prior," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 12, pp. 2341–2353, 2010. 2

[36] H. Zeng, J. Cai, L. Li, Z. Cao, and L. Zhang, "Learning image-adaptive 3d lookup tables for high performance photo enhancement in real-time," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2

[37] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788. 3

[38] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results," http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html. 4

[39] ——, "The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results," http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html. 4

[40] S. Narasimhan and S. Nayar, "Vision and the atmosphere," *International Journal of Computer Vision*, vol. 48, pp. 233–254, 07 2002. 4