

# FLOWBENCH: A ROBUSTNESS BENCHMARK FOR OPTICAL FLOW ESTIMATION

Anonymous authors

Paper under double-blind review

## ABSTRACT

Optical flow estimation is a crucial computer vision task often applied to safety-critical real-world scenarios like autonomous driving and medical imaging. While optical flow estimation accuracy has greatly benefited from the emergence of deep learning, learning-based methods are also known for their lack of generalization and reliability. However, reliability is paramount when optical flow methods are employed in the real world, where safety is essential. Furthermore, a deeper understanding of the robustness and reliability of learning-based optical flow estimation methods is still lacking, hindering the research community from building methods safe for real-world deployment. Thus we propose FLOWBENCH, a robustness benchmark and evaluation tool for learning-based optical flow methods. FLOWBENCH facilitates streamlined research into the reliability of optical flow methods by benchmarking their robustness to adversarial attacks and out-of-distribution samples. With FLOWBENCH, we benchmark 91 methods across 3 different datasets under 7 diverse adversarial attacks and 23 established common corruptions, making it the most comprehensive robustness analysis of optical flow methods to date. Across this wide range of methods, we consistently find that methods with state-of-the-art performance on established standard benchmarks lack reliability and generalization ability. Moreover, we find interesting correlations between performance, reliability, and generalization ability of optical flow estimation methods, under various lenses such as design choices used, number of parameters, etc. After acceptance, FLOWBENCH will be open-source and publicly available, including the weights of all tested models.

## 1 INTRODUCTION

NEW

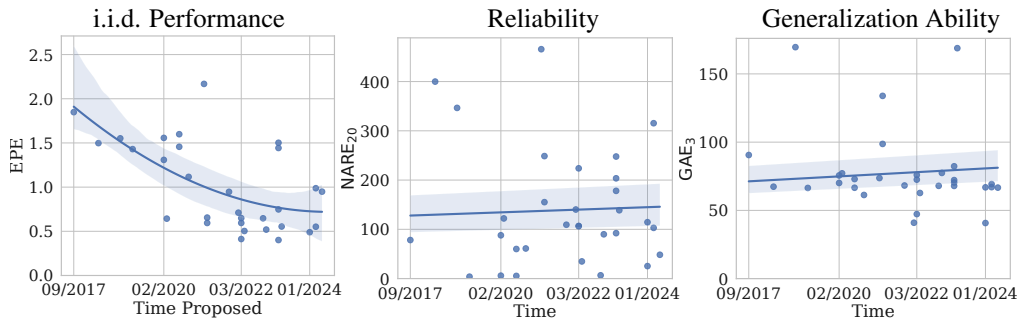


Figure 1: Optical flow estimation methods proposed over time and their reliability and generalization ability. In all three plots, the y-axis represents error, i.e., lower is better. The error of optical flow estimation methods on independent and identically distributed data samples (i.i.d.) has decreased over time, however, their reliability and generalization ability are stagnant if not deteriorating.

The recent growth of Deep Learning (DL) has greatly benefited computer vision, in particular when considering complex tasks such as the estimation of optical flow fields. In optical flow estimation, a method is supposed to estimate the movement of every pixel between at least two consecutive image frames in a subpixel-accurate manner. This task was earlier performed using model-driven

054 approaches such as Horn & Schunck (1981) and Lucas & Kanade (1981). However, these methods  
055 have severe limitations leading to suboptimal estimations and, consequently, to the predominant use  
056 of DL to perform the estimations (Dosovitskiy et al., 2015; Ilg et al., 2017; Jahedi et al., 2024b).  
057 The performance of learning-based optical flow estimation methods has improved over the years on  
058 independent and identically distributed data samples (i.i.d.), leading to lower errors on evaluation as  
059 shown by Fig. 1 (left). At the same time, DL-based methods are known to be unreliable (Geirhos  
060 et al., 2018; Prasad, 2022), they tend to learn shortcuts rather than meaningful feature representa-  
061 tions (Geirhos et al., 2020), and can be easily deteriorated even by small corruptions. This can  
062 become a practical threat, as optical flow estimation is highly relevant in safety-critical applications  
063 such as autonomous driving (Capito et al., 2020; Wang et al., 2021), robotic surgery (Rosa et al.,  
064 2019) and others. Thus, before deploying DL-based optical flow estimation methods, assessing their  
065 vulnerability and generalization ability is of paramount importance to gauge their readiness. We ob-  
066 serve in Fig. 1 that over the years, despite improvement in the performance of learning-based optical  
067 flow estimation methods, their reliability and generalization ability are almost unchanged. Had re-  
068 cent research been focused on these factors, the newly proposed methods could have been more  
069 reliable and ready for practical use. Our proposed FLOWBENCH facilitates this study, streamlining  
070 it for future research to utilize.

071 Many works have highlighted the importance of such a study by reducing model vulnerability (Xu  
072 et al., 2021b; Croce et al., 2023; Agnihotri et al., 2023; Schrodi et al., 2022; Tran et al., 2022;  
073 Grabinski et al., 2022), showing that robustness does follow from high accuracy (Tsipras et al.,  
074 2019; Schmidt et al., 2018; Schmalfluss et al., 2022b) or improving generalization (Hendrycks et al.,  
075 2020; Hoffmann et al., 2021) for various downstream tasks such as image classification, semantic  
076 segmentation, image restoration and others. To facilitate this research, robustness benchmarking  
077 tools and benchmarks like Croce et al. (2021); Jung et al. (2023); Tang et al. (2021) have been  
078 proposed for image classification models. They look into the adversarial and Out-of-Distribution  
079 (OOD) robustness of DL models. However, these works are limited to image classification. A  
080 similar benchmarking tool and comprehensive benchmark for optical flow is amiss.

081 To bridge this gap, we propose FLOWBENCH that facilitates robustness evaluations of optical flow  
082 models against adversarial attacks and image corruptions for OOD data and provides a unified eval-  
083 uation scheme and streamlined code. Using FLOWBENCH, we benchmark 91 model checkpoints  
084 over 3 commonly used optical flow estimation datasets. These model checkpoints include SotA opti-  
085 cal flow estimation methods and evaluation methods including SotA adversarial attacks and image  
086 corruption methods. FLOWBENCH is easy to use and new methods, when proposed, can be easily  
087 integrated to benchmark their performance. This will help researchers build better models that are  
088 not limited to improved performance on identical and independently distributed (i.i.d.) samples and  
089 are less vulnerable to adversarial attacks while generalizing better to image corruptions.

090 The main contributions of this work are as follows:

- 091 • We provide a benchmarking tool FLOWBENCH to evaluate the performance of most DL-  
092 based optical flow estimation methods over different datasets and make 91 checkpoints  
093 over different datasets publicly available for streamlined benchmarking while enabling the  
094 research community to add further checkpoints.
- 095 • We benchmark the aforementioned models against SotA and other commonly used adver-  
096 sarial attacks and common corruptions that can be easily queried using FLOWBENCH.
- 097 • We perform an in-depth analysis using FLOWBENCH and present interesting findings show-  
098 ing that methods that are SotA on i.i.d. are remarkably less reliable and generalize worse  
099 than other non-SotA methods.
- 100 • We analyze correlations between performance, reliability, and generalization abilities of  
101 optical flow estimation methods, under various lenses such as design choices, and the num-  
102 ber of learnable parameters.
- 103 • We show that the optimization of white-box adversarial attacks for optical flow estimation  
104 can be performed even without the availability of ground truth predictions, furthering the  
105 scope of study in their reliability.
- 106 • We show that the optimization of white-box adversarial attacks for optical flow estimation  
107 can be performed even without the availability of ground truth predictions, furthering the  
108 scope of study in their reliability.

## 2 RELATED WORK

FLOWBENCH is the first robustness benchmarking tool and benchmark for optical flow estimation methods that unifies adversarial and OOD robustness, taking inspiration from robustness benchmarks for other vision tasks such as image classification. While several previous works provide benchmarking tools for optical flow estimation, they only facilitate benchmarking of either adversarial or OOD robustness and are less comprehensive than FLOWBENCH. FLOWBENCH leverages the individual strengths of prior benchmarking tools, but casts them into a unified and easy-to-use robustness benchmark. Following, we discuss these related works in detail.

### 2.1 ROBUSTNESS BENCHMARKING FOR IMAGE CLASSIFICATION METHODS

Goodfellow et al. (2015) proposed the Fast Sign Gradient Method (FGSM) attack which gave rise to the domain of adversarial attacks on image classification. Complementing adversarial attacks, Hendrycks & Dietterich (2019) proposed 2D Common Corruptions for image classification tasks on the CIFAR-100 (Krizhevsky et al., 2009) and ImageNet-1k (Russakovsky et al., 2015) datasets and their variants. Since then, most adversarial attacks and OOD Robustness works have focused on image classification tasks, warranting a consolidated benchmarking tool and benchmark for robustness. In the case of image classification, this gap was filled by multiple works such as RobustBench (Croce et al., 2021) and RobustArts (Tang et al., 2021). Both works make multiple image classification model checkpoints publicly available, including checkpoints trained for improved robustness. Moreover, RobustBench is a benchmarking tool that facilitates evaluating both adversarial and OOD robustness of image classification models. Other similar benchmarking tools exist, like DeepFool (Moosavi-Dezfooli et al., 2016), Torchattacks (Kim, 2020), and Foolbox (Rauber et al., 2020). Yet, these are merely benchmarking tools and do not provide a comprehensive benchmark - they only facilitate evaluating adversarial robustness but not the OOD robustness of the method. As of now, no benchmarking tool or benchmark exists for optical flow estimation methods' robustness evaluations. Thus, we propose FLOWBENCH which enables benchmarking adversarial and OOD robustness and makes a multitude of model checkpoints available, providing the research community with the much needed tools.

### 2.2 BENCHMARKING OPTICAL FLOW ESTIMATION METHODS

Optical flow estimation has been a problem attempted to be solved for a long time. Over time multiple works have been proposed to streamline research in this direction by providing benchmarking libraries for i.i.d. performance of proposed methods. Such libraries include *mmflow* (Contributors, 2021), *ptflow* (Morimitsu, 2021), and *Spring* (Mehl et al., 2023). These libraries also provide model checkpoints to facilitate evaluations. *Spring*, also provides a benchmark but the performance evaluations are limited to their proposed Spring dataset. Whereas, both *mmflow* and *ptflow* do not provide a benchmark but enable benchmarking on multiple optical flow datasets such as FlyingThings3D (Mayer et al., 2016), KITTI2015 (Menze & Geiger, 2015) and MPI Sintel (Butler et al., 2012). However, the evaluation abilities of these benchmarking tools are limited to i.i.d. data. Thus, we built FLOWBENCH, using *ptflow* and publicly available model checkpoints to extend method evaluations to adversarial and OOD Robustness consolidating research towards reliability and generalization ability of optical flow estimation methods. Additionally, FLOWBENCH is the first to provide a comprehensive benchmark on existing optical flow estimation methods over 3 datasets and multiple adversarial attacks and image corruptions.

### 2.3 ADVERSARIAL ATTACKS

As discussed in Sec. 1, DL models tend to learn shortcuts to map data samples from input to target distribution (Geirhos et al., 2020), leading to the model learning inefficient feature representations. In their work, Goodfellow et al. (2015) showed that this inefficient learning of feature representations can be easily exploited. Goodfellow et al. (2015) added noise to the input data samples which was optimized to increase loss using model information, such that the model was fooled into making incorrect predictions. This demonstrated the vulnerability and unreliability of model predictions as the perturbed input samples still appeared semantically similar to the human eye. They named this attack the **Fast Sign Gradient Method** (FGSM). This attack led to an increased inter-

est by the research community to better optimize the noise inspiring multiple other works such as Basic Iteration method (BIM) (Kurakin et al., 2018), Projected Gradient Descent (PGD) (Kurakin et al., 2017), Auto-PGD (APGD) (Wong et al., 2020) and CosPGD (Agnihotri et al., 2024) which were direct extensions to FGSM, and other attacks such as Perturbation-Constrained Flow Attack (PCFA) (Schmalfuss et al., 2022b) and Adversarial Weather (Schmalfuss et al., 2023), which are indirect extensions of FGSM.

### 3 FLOWBENCH USAGE

In the following, we describe the benchmarking tool, FLOWBENCH. It is built using `ptlflow` (Morimitsu, 2021), and supports 36 unique architectures (new architectures added to `ptlflow` over time are compatible with FLOWBENCH) and distinct datasets, namely FlyingThings3D (Mayer et al., 2016), KITTI2015 (Menze & Geiger, 2015), MPI Sintel (Butler et al., 2012) (clean and final) and Spring (Mehl et al., 2023) datasets (please refer Appendix C for additional details on the datasets). It enables training and evaluations on all aforementioned datasets including evaluations using SotA adversarial attacks such as CosPGD (Agnihotri et al., 2024) and PCFA (Schmalfuss et al., 2022b), Adversarial weather (Schmalfuss et al., 2023), and other commonly used adversarial attacks like BIM (Kurakin et al., 2018), PGD (Kurakin et al., 2017), FGSM (Goodfellow et al., 2015), under various Lipschitz ( $l_p$ ) norm bounds.

Additionally, it enables evaluations for Out-of-Distribution (OOD) robustness by corrupting the inference samples using 2D Common Corruptions (Hendrycks & Dietterich, 2019) and 3D Common Corruptions (Kar et al., 2022).

We follow the nomenclature set by RobustBench (Croce et al., 2021) and use “`threat_model`” to define the kind of evaluation to be performed. When “`threat_model`” is defined to be “None”, the evaluation is performed on unperturbed and unaltered images, if the “`threat_model`” is defined to be an adversarial attack, for example “PGD”, “CosPGD” or “PCFA”, then FLOWBENCH performs an adversarial attack using the user-defined parameters. We elaborate on this in Appendix E.1. Whereas, if “`threat_model`” is defined to be “2DCommonCorruptions” or “3DCommonCorruptions”, the FLOWBENCH performs evaluations after perturbing the images with 2D Common Corruptions and 3D Common Corruptions respectively. We elaborate on this in Appendix E.2. If the queried evaluation already exists in the benchmark provided by this work, then FLOWBENCH simply retrieves the evaluations, thus saving computation.

FLOWBENCH enables the use of all the attacks mentioned in Sec. 2.3 to help users better study the reliability of their optical flow methods. We choose to specifically include these white-box adversarial attacks as they either serve as the common benchmark for adversarial attacks in classification literature (FGSM, BIM, PGD, APGD) or they are unique attacks proposed specifically for pixel-wise prediction tasks (CosPGD) and optical flow estimation (PCFA and Adversarial Weather). These attacks can either be *Non-targeted* which are designed to simply fool the model into making incorrect predictions, irrespective of what the model eventually predicts, or can be *Targeted*, where the model is fooled to make a certain prediction. Most attacks can be, designed to be either Targeted or Non-targeted, these include, FGSM, BIM, PGD, APGD, CosPGD, and Adversarial Weather. However, by design, some attacks are limited to being only one of the two, for example, PCFA which is a targeted attack.

Following we show the basic commands to use FLOWBENCH. We describe each attack and common corruption supported by FLOWBENCH in detail in Appendix E. Please refer to Appendix G for details on the arguments and function calls.

#### 3.1 MODEL ZOO

It is a challenge to find all checkpoints, while training them is a time and compute exhaustive process. Thus we gather available model checkpoints from various sources such as `ptlflow` (Morimitsu, 2021) and `mmflow` (Contributors, 2021). The trained checkpoints for all models available in FLOWBENCH can be obtained using the following lines of code:

```
from flowbench.evals import load_model
model = load_model(model_name='RAFT', dataset='KITTI2015')
```

Each model checkpoint can be retrieved with the pair of ‘model\_name’, the name of the model, and ‘dataset’, the dataset for which the checkpoint was last finetuned. In Appendix F we provide a complete overview of all the 91 available pairs of model checkpoints and datasets.

### 3.2 ADVERSARIAL ATTACKS

FLOWBENCH can be used to evaluate models on the discussed adversarial attacks using the following lines of code (please refer Appendix G.1 for details regarding the arguments):

```
from flowbench.evals import evaluate
model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
                          threat_model='CosPGD', iterations=20, alpha=0.01,
                          epsilon=8/255, lp_norm='Linf', targeted=True,
                          optim_wrt='ground_truth', retrieve_existing=True)
```

### 3.3 OOD ROBUSTNESS

NEW

FLOWBENCH can be used to evaluate models on the 2D and 3D Common Corruptions using the following lines of code, following is an example for the latter (please refer Appendix G.3 (2D Common Corruptions) and Appendix G.4 (3D Common Corruption) for details regarding the arguments):

```
from flowbench.evals import evaluate
model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
                          threat_model='3DCommonCorruption',
                          severity=3, retrieve_existing=True)
```

## 4 METRICS FOR ANALYSIS AT SCALE

Analysis of optical flow estimation methods at the same scale as this work, especially under the lens of reliability and generalization ability has not been attempted before. The most commonly (Schrodi et al., 2022; Schmalfluss et al., 2022a; Agnihotri et al., 2024; Dosovitskiy et al., 2015) used metric for evaluating the performance of a method is calculating the mean End-Point-Error (EPE) between the predicted optical flow and the ground truth for all pairs of frames in a given dataset. However, this does not reflect the reliability and generalization ability of the method. Moreover, this work has performed over 4500 experiments in total, and analyzing the EPE from each experiment would not lead to a fruitful finding. Thus, we attempt to simplify this with our proposed metrics, the [Reliability Error](#) and [Generalization Ability Error](#).

The objective of any optical flow estimation method is to obtain an EPE of zero or as low as possible. The larger the EPE, the worse the performance of the method. Most works (Dosovitskiy et al., 2015; Teed & Deng, 2020; Ilg et al., 2017; Huang et al., 2022) report the mean EPE value over a dataset as a measure of the method’s performance. For reliability and generalization, we look at the maximum possible value of mean EPE across attacks over multiple datasets. That is, we ask the question “What is the worst possible performance of a given method?”. An answer to this question tells us about the reliability and generalization ability of a method. In the following, we describe the measures for different scenarios in detail.

### 4.1 GENERALIZATION ABILITY ERROR

NEW

Inspired by multiple works (Croce et al., 2021; Hendrycks et al., 2020; Hoffmann et al., 2021) that use OOD Robustness of methods for evaluating the generalization ability of the method, even evaluate over every common corruptions, that is 2D Common Corruptions and 3D Common Corruptions combined. Then, we find the maximum of the mean EPE w.r.t. the ground truth for a given method, across all corruptions at a given severity and report this as [Generalization Ability Error](#) denoted by  $GAE_{severity\ level}$ . For example, for severity 3, the measure would be denoted by  $GAE_3$ . The less the GAE value, the better the generalization ability of the given optical flow estimation method. These corruptions perturb the images to cause distributions and domain shifts, such shifts often confuse the methods into making incorrect predictions.

For calculating GAE, we use all 15 2D Common Corruptions: ‘Gaussian Noise’, ‘Shot Noise’, ‘Impulse Noise’, ‘Defocus Blur’, ‘Frosted Glass Blur’, ‘Motion Blur’, ‘Zoom Blur’, ‘Snow’, ‘Frost’, ‘Fog’, ‘Brightness’, ‘Contrast’, ‘Elastic Transform’, ‘Pixelate’, ‘JPEG Compression’, and eight 3D Common Corruptions: ‘Color Quantization’, ‘Far Focus’, ‘Fog 3D’, ‘ISO Noise’, ‘Low Light’, ‘Near Focus’, ‘XY Motion Blur’, and ‘Z Motion Blur’. All the common corruptions are at severity 3. Kar et al. (2022) offers more 3D Common Corruptions, however computing them is resource intensive. Thus, given our limited resources and an overlap in the corruptions between 2D Common Corruptions and 3D Common Corruptions, we focus on generating 3D Common Corruptions that might be unique from their 2D counterpart, require fewer sources to generate, and are interesting from an optical flow estimation perspective. In Appendix A we show that these synthetic common corruptions can indeed be used as a proxy for possible corruptions when in the wild in the real world.

## 4.2 RELIABILITY ERROR

NEW

An adversarial attack is a perturbation made on the input images to fool a method into changing its predictions while the input image looks semantically similar to a human observer. Most works that focus on the reliability of optical flow estimation methods perform adversarial attacks, however, these works either focus on targeted attacks or on non-targeted attacks, not both at the same time. The objective of targeted attacks is to optimally perturb the input image such that the method predictions are changed towards a specifically desired target, for example, a target can be a  $\vec{0}$  flow i.e. attacking so that the flow prediction at all pixels should become zero. Conversely, non-targeted adversarial attacks do not intend to shift the method’s predictions to a specific target, they simply intend to fool the method into making any incorrect predictions. To streamline research into the reliability of these methods, we perform both targeted and non-targeted attacks.

**Non-Targeted Attacks.** For non-targeted attacks, we measure the EPE w.r.t. the ground truth, in this case, the higher the EPE value, the worse the performance of the optical flow estimation method. The notation for this metric is,  $\text{NARE}_{\text{attack iterations}}$ , where NARE stands for Non-targeted Attack Reliability Error, and the subscript informs the number of attack iterations used for optimizing the attack. For example, when 20 attack iterations were used to optimize the attack then the metric would be  $\text{NARE}_{20}$ . The higher the NARE value, the worse the reliability of the optical flow method.

**Targeted Attacks.** For targeted attacks, we measure the EPE w.r.t. the target flow, however, to standardize notations, we report the negative EPE in this case, thus, the higher the value, the worse the performance of the optical flow estimation method. The notation for this metric is,  $\text{TARE}_{\text{attack iterations}}^{\text{target}}$ , where TARE stands for Targeted Attack Reliability Error and the superscript informs about the target used (zero vector or negative of the initial flow prediction) and the subscript informs about the number of attack iterations used for optimizing the attack. For example, when the target is  $\vec{0}$  and 20 attack iterations were used to optimize the attack then the metric would be  $\text{TARE}_{20}^{\vec{0}}$ . The higher the TARE value, the worse is the reliability of the optical flow method.

For calculating TARE and NARE values we used BIM, PGD, and CosPGD attack with step size  $\alpha=0.01$ , perturbation budget  $\epsilon = \frac{8}{255}$  under the  $\ell_\infty$ -norm bound, as targeted and non-targeted attacks respectively. We use  $\ell_\infty$ -norm bound as we observe in Appendix H that there is a high correlation between the performance of optical flow estimation methods when attacked using  $\ell_\infty$ -norm bounded attacks and  $\ell_2$ -norm bounded attacks. We use 20 attack iterations for calculating TARE and NARE as we observe in Appendix H, that at a lower number of iterations, the gap in performance of different optical flow estimation methods is small, thus an in-depth analysis would be difficult, and we do not go beyond 20 attack iterations as computing each attack step for an adversarial attack is very expensive, and as shown by Agnihotri et al. (2024) and Schmalfluss et al. (2022b), 20 iterations are enough to optimize an attack to truly understand the performance of the attacked method.

## 5 ANALYSIS AND INTERESTING FINDINGS

To demonstrate the potential of FLOWBENCH, we use it to perform multiple analyses which provide us with a better understanding of many optical flow estimation methods, including novel findings.



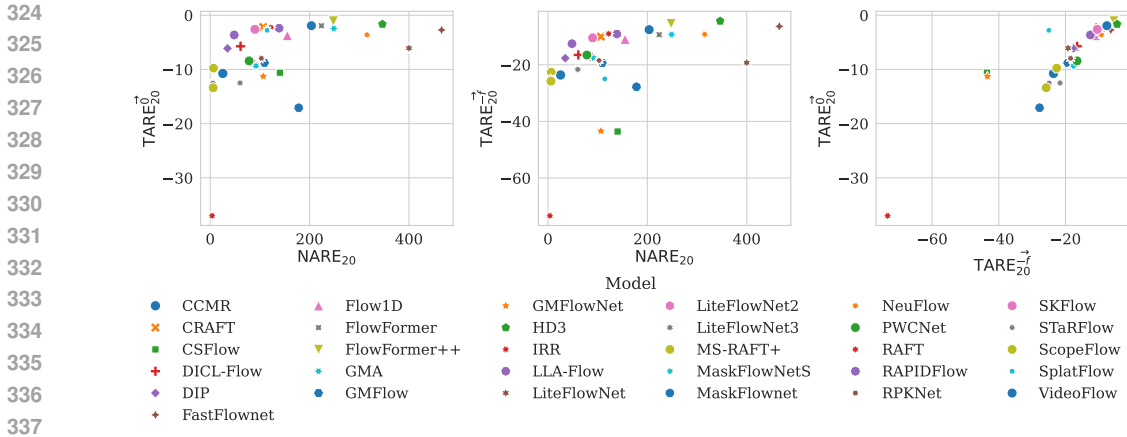


Figure 2: Analysing correlations between Targeted and Non-targeted adversarial attacks. A model is more reliable if it has a low NARM value and a high TARM value.

Following, we discuss the observations made in the comprehensive robustness benchmark created using FLOWBENCH. Please refer to Appendix C for details on the dataset, Appendix D for additional implementation details, and Appendix H for additional results from the benchmarking.

### 5.1 TARGETED V/S NON-TARGETED ADVERSARIAL ATTACKS

We benchmark the performance of all prominent DL-based optical flow estimation methods across three datasets, namely KITTI2015, MPI Sintel (clean), and MPI Sintel (final) against SotA and commonly used adversarial attacks such as BIM, PGD, and CosPGD. Then, we compare the NARE and TARE values (introduced in Sec. 4.2) and find correlations in their performance. These are reliability metrics, higher NARE and TARE values indicates low reliability and vice versa. Please refer to Appendix D for more implementation details. We observe in Fig. 2 that there is a very high correlation between the  $TARE_{20}^0$  and  $TARE_{20}^1$  values of every optical flow estimation method. This shows that evaluating either one of the values can serve as a reliable proxy for the other. We use this finding in the later analysis. Additionally, in Fig. 2 we observe that most optical flow estimation methods like ScopeFlow (Bar-Haim & Wolf, 2020), MS-RAFT+ (Jahedi et al., 2024b) and StarFlow (Godet et al., 2021) are relatively more susceptible to targeted attacks than they are to non-targeted attacks. On the other hand, some methods are highly susceptible to both and thus very unreliable, these include SKFlow (Sun et al., 2022), FastFlowNet (Kong et al., 2021), HD3 (Yin et al., 2019) and some SotA methods like FlowFormer (Huang et al., 2022) and FlowFormer++ (Shi et al., 2023b). Interestingly, IRR (Hur & Roth, 2019) stands out as the most reliable optical flow estimation method as it is robust to both targeted and non-targeted adversarial attacks. While ScopeFlow (Bar-Haim & Wolf, 2020), GMFlowNet (Zhao et al., 2022) and MaskFlowNet (Zhao et al., 2020) are less reliable than IRR but more reliable than the other methods.

### 5.2 RELIABILITY V/S GENERALIZATION

Following we analyze if there is a correlation between the reliability and generalization ability of optical flow estimation methods. We observe in Fig. 3, that most methods that have a good performance also generalize better, however, methods like FlowFormer++, while having good i.i.d. performance have a relatively poor generalization ability. As observed in Sec. 5.1, HD3 Yin et al. (2019) stands out as having poor performance and poor generalization ability. Interestingly, as shown by Fig. 3, there is a correlation between the generalization ability ( $GAE_3$  values, introduced in Sec. 4.1, higher GAE value indicates lower generalization ability) and reliability when measured using non-targeted adversarial attacks ( $NARE_{20}$  values). Additionally, most methods identified in Sec. 5.1 to be reliable, for example, CSFlow, MaskFlowNet also have considerable generalization ability compared to the other methods. However, IRR which stood out as the most reliable method has low generalization abilities. It is interesting to note that CCMR (Jahedi et al., 2024a) offers a good trade-off as it has reasonably good performance, reliability, and generalization abilities.

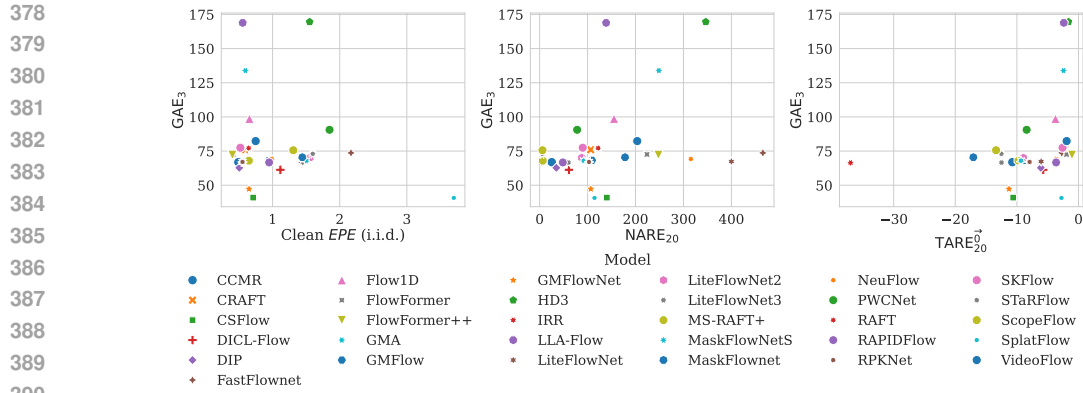


Figure 3: Analysing correlations between reliability and generalization ability of optical flow estimation methods.

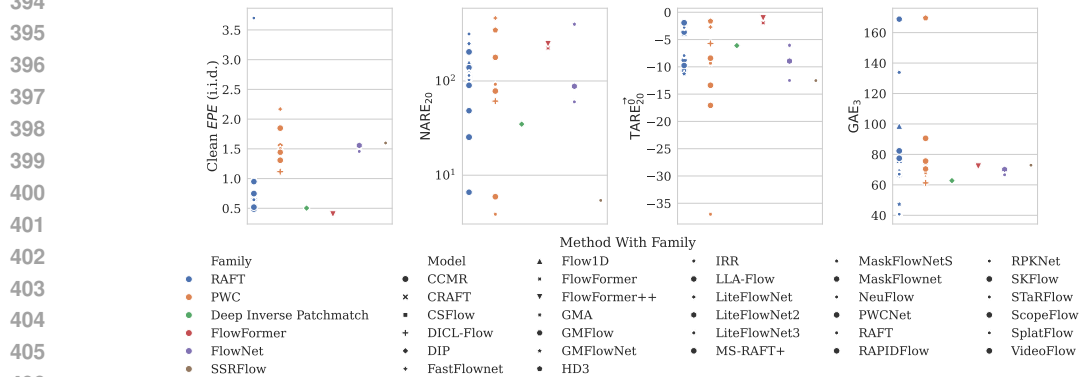


Figure 4: Analysing correlations between the method family to which the optical flow estimation method belongs and its corresponding performance, reliability, and generalization ability.

### 5.3 ANALYZING METHOD FAMILIES

NEW

Optical flow estimation methods proposed over the years use different training strategies and architecture designs. However, there exist many architectural similarities between the methods, and based on these most methods can be broadly classified into four method families: FlowNet-family, PWC-family, RAFT-family, and FlowFormer-family (please refer to Appendix B for detailed justifications). In Fig. 4 we observe that methods belonging to the FlowFormer-family and RAFT-family and DIP Zheng et al. (2022) have the best i.i.d. performance, however, given their relatively higher  $NARE_{20}$  and  $TARE_{20}$  values, some exceeding 100, they appear to not be reliable. Here we observe that IRR (Hur & Roth, 2019) stands out as one of the most reliable methods under adversarial attacks. Given that the primary differences between IRR-PWC and other methods from the PWC family are the classical energy minimization-inspired approach and the use of residual networks to propose an iterative residual refinement, it makes an interesting finding.

When considering generalization ability under common corruptions, we observe all methods to have poor performance. Methods such as LLA-Flow Xu et al. (2023b) and HD3 Yin et al. (2019) from the RAFT-family and PWC-family respectively have  $GAE_3$  values over 160! Here, SplatFlow Wang et al. (2024) stands out, given that the primary difference between SplatFlow and other RAFT-family methods is the use of splatting for feature matching by SplatFlow, it is an interesting finding.

Additionally, we observe in Fig. 4 that compared to other method families, FlowFormer-family is very susceptible to targeted adversarial attacks. Given that the FlowFormer family comprises only transformer-based architectures for optical flow estimation, this is very interesting, as this contradicts the observations made for transformer-based methods for image classification (Paul & Chen, 2022;



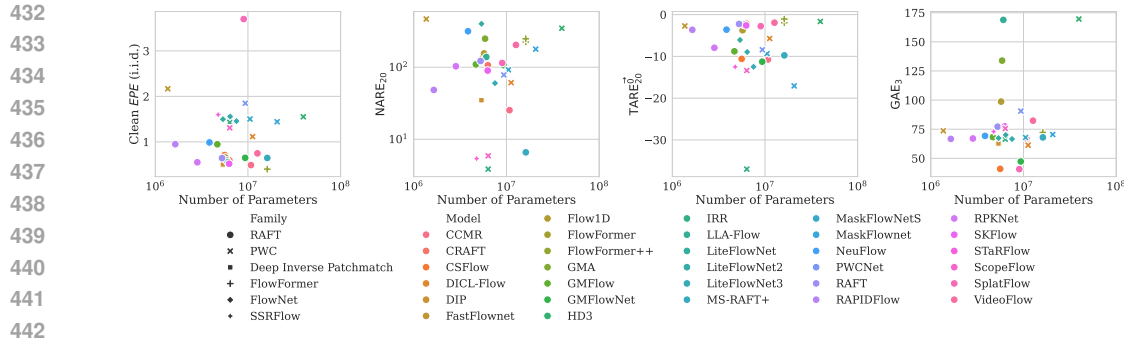


Figure 5: Analysing correlation between the number of learnable parameters in a DL-based optical flow estimation method and its performance, reliability, and generalization ability. Colors show the different optical flow methods while marker styles show the method family to which they belong.

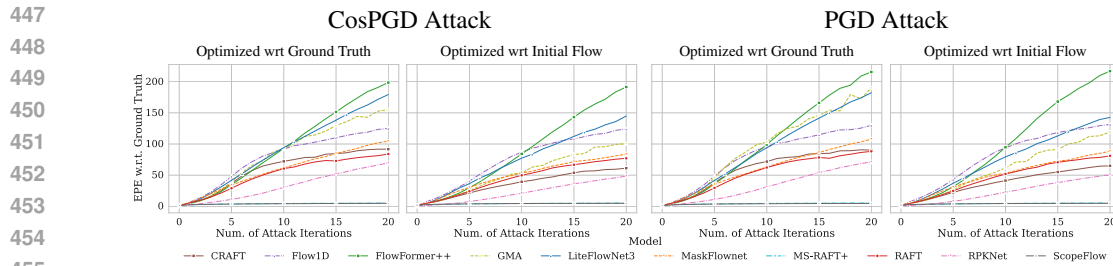


Figure 6: Performance of interesting optical flow estimation methods under different non-targeted adversarial attacks optimized using initial flow predictions on the KITTI2015 dataset.

Hoyer et al., 2022) and semantic segmentation (Xie et al., 2021), where they appear to be more robust. However, this lack of generalization ability can also be attributed to the use of dynamic positional cost queries by both FlowFormer and FlowFormer++, thus, we require more models to be proposed in the FlowFormer family to be certain.

#### 5.4 IMPACT OF THE NUMBER OF LEARNABLE PARAMETERS

Many works for classification have shown that Deep Neural Networks with more parameters and less vulnerable to adversarial attacks and generalize better to common corruptions (Liu et al., 2022; Ding et al., 2022; Hoffmann et al., 2021). It would be interesting to see if the same holds true for optical flow estimation methods. Thus, we analyze this in Fig. 5 and observe that while the number of learnable parameters has an impact on the performance of the methods to some extent (other than the exceptions of MaskFlowNet and HD3), the same does not hold for reliability and generalization ability. Methods such as FlowFormer, FlowFormer++ (FlowFormer-family), and VideoFlow (RAFT-family) have relatively more parameters than other methods however they are less reliable and have a poor generalization ability. On the other hand, methods like CSFlow and SplatFlow (both RAFT-family) have significantly fewer parameters but are more reliable and generalize better than the other methods.

#### 5.5 OPTIMIZING TARGETED ATTACKS USING INITIAL FLOW PREDICTIONS

Based on the observation in Sec. 5, we identify several interesting methods whose performance warrants additional analysis and discussion. Following, we discuss our observations in detail.

One of the major limitations of white-box adversarial attacks is that they require access to the ground truth to optimize the attack (Agnihotri et al., 2024). However, access to the ground truth is not guaranteed in every scenario. Additionally as discussed by Schmalfluss et al. (2022b), robustness is a measure of the difference in a model’s prediction on perturbed input w.r.t. the model’s prediction on clean input samples. Thus, the goal of an attack should be to fool the method into changing its initial

486 predictions (predictions when the method is not attacked), independent of the ground truth. Thus, we  
 487 attempt to optimize the adversarial attack w.r.t. to the initial flow prediction on the unperturbed input  
 488 sample before any attacks, as access to this is almost guaranteed. This helps us ascertain if initial  
 489 flow predictions can be used as a proxy to ground truth while optimizing attacks. Thus, in Eq. (4),  
 490 Eq. (8), Eq. (9) and there places where applicable  $Y=X^{\text{clean}}$  (please refer Appendix E.1). However,  
 491 this optimization is only possible for attacks that introduce certain randomness in the initial input  
 492 sample, as shown by Eq. (7). This allows for there to exist a non-zero loss between the predictions  
 493 on the clean input samples and the perturbed input samples allowing for optimization. We report the  
 494 evaluations for CosPGD and PGD attack using the KITTI2015 dataset for 10 interesting methods in  
 495 Fig. 6. We choose the optical flow estimation methods on the basis of their performance in Sec. 5  
 496 and their performance on i.i.d. samples. For additional evaluation using more models please refer  
 497 to Appendix H. We observe in Fig. 6 that there appears a high correlation in the performance of  
 498 all considered methods under attack when optimized using the ground truth flow and the initial flow  
 499 prediction. Thus, initial flow predictions from methods do serve as a strong proxy to the ground truth  
 500 for optimizing attacks. This new finding over a big sample, helps advance study in the reliability of  
 501 optical flow methods, even when ground truth predictions are not available.

## 502 6 CONCLUSION

NEW

505 FLOWBENCH is the first robustness benchmarking tool and a novel benchmark for optical flow esti-  
 506 mation methods. It currently supports 91 model checkpoints, over distinct datasets, and all relevant  
 507 robustness evaluation methods including SotA adversarial attacks and image corruptions. We dis-  
 508 cuss the unique features of FLOWBENCH in detail and demonstrate that the library is user-friendly.  
 509 Adding new evaluation methods or optical flow estimation methods to FLOWBENCH is easy and  
 510 intuitive. In Sec. 5.1, we find that there is a high correlation in the performance of optical flow  
 511 estimation methods against targeted attacks using different targets, thus saving compute for future  
 512 works as they need to evaluate only against one target. In Sec. 5.2, we observe the methods known  
 513 to be SotA on i.i.d. samples are not reliable, and do not generalize well to image corruptions,  
 514 demonstrating the gap in current research when considering real-world applications. Additionally,  
 515 we observe here that there is no apparent correlation between generalization abilities and the reli-  
 516 ability of optical flow estimation methods. In Sec. 5.3, we show that methods from the FlowFormer  
 517 family have good i.i.d. performance but are the most unreliable under targeted attacks, also that IRR  
 518 stands out to have marginally better reliability. generalization abilities. In Sec. 5.4, we show that,  
 519 unlike image classification, increasing the number of learnable parameters does not help increase the  
 520 robustness of optical flow estimation methods, however, a couple of RAFT variants have marginally  
 521 better generalization abilities even with fewer parameters. These observation helps us conclude that  
 522 based on current works, different approaches might be required to attain reliability under attacks and  
 523 generalization ability to image corruptions. Lastly, we show that white-box adversarial attacks on  
 524 optical flow estimation methods can be independent of the availability of ground truth information,  
 525 and can harness the information in the initial flow predictions to optimize attacks, thus overcoming a  
 526 huge limitation in the field. Such an in-depth understanding of reliability and generalization abilities  
 527 to optical flow estimation methods can only be obtained using our proposed FLOWBENCH. We are  
 528 certain that FLOWBENCH will be immensely helpful in gathering more such interesting findings and  
 529 its comprehensive and consolidated nature would make things easier for the research community.

529 **Future Work.** For optical flow estimation, patch attacks are also interesting and widely stud-  
 530 ied (Ranjan et al., 2019; Schrodi et al., 2022; Scheurer et al., 2024). We plan to add such patch  
 531 attacks to FLOWBENCH in future iterations. Schmalfluss et al. (2022b) proposed optimizing adver-  
 532 sarial noise jointly for the consecutive image frames and also over the entire evaluation set. Only  
 533 PCFA supports such optimization regimes in FLOWBENCH, so it would be interesting to extend such  
 534 optimization to other adversarial attacks as well. Croce et al. (2021) show that the training methods  
 535 used significantly impact the robustness of image classification methods. The same might be true  
 536 for optical flow estimation methods, thus robustness evaluations under the lens of different training  
 537 setups used would make an interesting extension to the analysis in this work. Lastly, traditional non-  
 538 DL-based optical flow estimation methods might be more robust to adversarial attacks than current  
 539 DL-based methods. Thus, it would be interesting to study their robustness and hopefully adapt them  
 to increase the reliability of current methods.

## REPRODUCIBILITY STATEMENT

Every experiment in this work is reproducible and is part of an effort toward open-source work. FLOWBENCH will be open source and publicly available, including all evaluation logs and model checkpoint weights. This work intends to help the research community build more reliable and generalizable optical flow estimation methods such that they are ready for deployment in the real world even under safety-critical applications. FLOWBENCH is built upon ptflow and thus any new model added with ptflow would most likely be supported by FLOWBENCH as well.

There always exists stochasticity when evaluating adversarial attacks, due to the randomness these attacks exploit, and when evaluating common corruptions due to different seeds and calculation approximations made by different python libraries. Therefore, for transparency and reproducibility, we evaluate different runs on the same seed and different runs of different seeds. We report these evaluations in Appendix K, using Tab. 2 for adversarial attacks and Tab. 3 for common corruptions, and observe that the variance is extremely low and the analysis performed in this work still stands.

## REFERENCES

- Shashank Agnihotri, Julia Grabinski, and Margret Keuper. Improving stability during upsampling – on the importance of spatial context, 2023.
- Shashank Agnihotri, Steffen Jung, and Margret Keuper. CosPGD: an efficient white-box adversarial attack for pixel-wise prediction tasks. In *Proc. International Conference on Machine Learning (ICML)*, 2024.
- Anonymous. SemSegBench: Robustness-Aware Benchmarking Of Semantic Segmentation.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEit: BERT pre-training of image transformers. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=p-BhZSz59o4>.
- Aviram Bar-Haim and Lior Wolf. Scopeflow: Dynamic scene scoping for optical flow. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7998–8007, 2020.
- Daniel J. Butler, Jonas Wulff, Garrett B. Stanley, and Michael J. Black. A naturalistic open source movie for optical flow evaluation. In *Proc. European Conference on Computer Vision (ECCV)*, LNCS 7577, pp. 611–625, 2012.
- Linda Capito, Umit Ozguner, and Keith Redmill. Optical flow based visual potential field for autonomous driving. In *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 885–891. IEEE, 2020.
- Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Advances in neural information processing systems*, 34:9355–9366, 2021.
- MMFlow Contributors. MMFlow: Openmmlab optical flow toolbox and benchmark. <https://github.com/open-mmlab/mmlflow>, 2021.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding, 2016.
- Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo DeBenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. RobustBench: a standardized adversarial robustness benchmark. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Francesco Croce, Naman D Singh, and Matthias Hein. Robust semantic segmentation: Strong adversarial attacks and fast training of robust models. *arXiv preprint arXiv:2306.12941*, 2023.

- 594 Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-  
595 efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*,  
596 35:16344–16359, 2022.
- 597
- 598 Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to  
599 31x31: Revisiting large kernel design in cnns. In *Proc. of the IEEE/CVF conference on computer  
600 vision and pattern recognition*, pp. 11963–11975, 2022.
- 601
- 602 Qiaole Dong, Chenjie Cao, and Yanwei Fu. Rethinking optical flow from geometric matching con-  
603 sistent perspective. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition  
604 (CVPR)*, pp. 1337–1347, 2023.
- 605
- 606 Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov,  
607 Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with  
608 convolutional networks. In *Proc. of the IEEE international conference on computer vision*, pp.  
2758–2766, 2015.
- 609
- 610 Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and  
611 Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias im-  
612 proves accuracy and robustness. In *International Conference on Learning Representations*, 2018.
- 613
- 614 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,  
615 Matthias Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature  
616 Machine Intelligence*, 2(11):665–673, nov 2020.
- 617
- 618 Pierre Godet, Alexandre Boulch, Aurélien Plyer, and Guy Le Besnerais. STaRFlow: A spatiotem-  
619 poral recurrent cell for lightweight multi-frame optical flow estimation. In *Proc. IEEE International  
620 Conference on Pattern Recognition (ICPR)*, pp. 2462–2469, 2021.
- 621
- 622 Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial  
623 examples. In *Proc. International Conference on Learning Representations (ICLR)*, 2015.
- 624
- 625 Julia Grabinski, Paul Gavrikov, Janis Keuper, and Margret Keuper. Robust models are less over-  
626 confident. *NeurIPS*, 2022.
- 627
- 628 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common  
629 corruptions and perturbations. In *Proc. International Conference on Learning Representations  
630 (ICLR)*, 2019.
- 631
- 632 Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-  
633 narayanan. AugMix: A simple data processing method to improve robustness and uncertainty.  
634 *Proc. of the International Conference on Learning Representations (ICLR)*, 2020.
- 635
- 636 J Hoffmann, S Agnihotri, Tonmoy Saikia, and Thomas Brox. Towards improving robustness of  
637 compressed cnns. In *ICML Workshop on Uncertainty and Robustness in Deep Learning (UDL)*,  
638 2021.
- 639
- 640 Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):  
641 185–203, 1981.
- 642
- 643 Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and  
644 training strategies for domain-adaptive semantic segmentation. In *Proceedings of the IEEE/CVF  
645 conference on computer vision and pattern recognition*, pp. 9924–9935, 2022.
- 646
- 647 Zhaoyang Huang, Xiaoyu Shi, Chao Zhang, Qiang Wang, Ka Chun Cheung, Hongwei Qin, Jifeng  
Dai, and Hongsheng Li. FlowFormer: A transformer architecture for optical flow. In *Proc.  
European Conference on Computer Vision (ECCV)*, LNCS 13677, pp. 668–685, 2022.
- Tak-Wai Hui and Chen Change Loy. LiteFlowNet3: Resolving correspondence ambiguity for more  
accurate optical flow estimation. In *Proc. European Conference on Computer Vision (ECCV)*,  
LNCS 12365, pp. 169–184, 2020.

- 648 Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. LiteFlowNet: A lightweight convolutional neural  
649 network for optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and*  
650 *Pattern Recognition (CVPR)*, pp. 8981–8989, 2018.
- 651 Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow CNN—revisiting  
652 data fidelity and regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*  
653 *(TPAMI)*, 43(8):2555–2569, 2020.
- 654 Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion  
655 estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
656 pp. 5754–5763, 2019.
- 657 Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox.  
658 FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *Proc. IEEE/CVF Con-*  
659 *ference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2462–2470, 2017.
- 660 Azin Jahedi, Maximilian Luz, Marc Rivinius, and Andrés Bruhn. CCMR: High resolution optical  
661 flow estimation via coarse-to-fine context-guided motion reasoning. In *Proc. IEEE/CVF Winter*  
662 *Conference on Applications of Computer Vision (WACV)*, pp. 6899–6908, 2024a.
- 663 Azin Jahedi, Maximilian Luz, Marc Rivinius, Lukas Mehl, and Andrés Bruhn. MS-RAFT+: high  
664 resolution multi-scale raft. *International Journal of Computer Vision (IJCV)*, 132(5):1835–1856,  
665 2024b.
- 666 Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate  
667 hidden motions with global motion aggregation. In *Proc. IEEE/CVF International Conference on*  
668 *Computer Vision (ICCV)*, pp. 9772–9781, 2021a.
- 669 Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning optical flow from a few matches.  
670 In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16592–  
671 16600, 2021b.
- 672 Steffen Jung, Jovita Lukasik, and Margret Keuper. Neural architecture design and robustness: A  
673 dataset. In *ICLR. OpenReview. net*, 2023.
- 674 Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and  
675 data augmentation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*  
676 *(CVPR)*, pp. 18963–18974, 2022.
- 677 Sarra Khairi, Etienne Meunier, Renaud Fraise, and Patrick Bouthemy. Efficient local correlation  
678 volume for unsupervised optical flow estimation on small moving objects in large satellite images.  
679 In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 440–  
680 448, 2024.
- 681 Hoki Kim. Torchattacks: A pytorch repository for adversarial attacks. *arXiv preprint*  
682 *arXiv:2010.01950*, 2020.
- 683 Daniel Kondermann, Rahul Nair, Stephan Meister, Wolfgang Mischler, Burkhard Güssefeld, Sabine  
684 Hofmann, Claus Brenner, and Bernd Jähne. Stereo ground truth with error bars. In *Asian Confer-*  
685 *ence on Computer Vision, ACCV 2014*, 2014.
- 686 Lingtong Kong, Chunhua Shen, and Jie Yang. FastFlowNet: A lightweight network for fast optical  
687 flow estimation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pp.  
688 10310–10316, 2021.
- 689 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 and cifar-100 datasets, learning multi-  
690 ple layers of features from tiny images. URL: <https://www.cs.toronto.edu/kriz/cifar.html>, 6(1):  
691 1, 2009.
- 692 Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In  
693 *Proc. International Conference on Learning Representations (ICLR)*, 2017.
- 694 Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world.  
695 In *Artificial Intelligence Safety and Security*, pp. 99–112. Chapman and Hall/CRC, 2018.

- 702 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.  
703 A convnet for the 2020s. In *Proc. of the IEEE/CVF conference on computer vision and pattern*  
704 *recognition*, pp. 11976–11986, 2022.
- 705  
706 Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to  
707 stereo vision. In *IJCAI’81: 7th international joint conference on Artificial intelligence*, volume 2,  
708 pp. 674–679, 1981.
- 709 Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy,  
710 and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow,  
711 and scene flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern*  
712 *Recognition (CVPR)*, pp. 4040–4048, 2016.
- 713 Lukas Mehl, Jenny Schmalfluss, Azin Jahedi, Yaroslava Nalivayko, and Andrés Bruhn. Spring: A  
714 high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo. In  
715 *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4981–  
716 4991, 2023.
- 717  
718 Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proc. IEEE/CVF*  
719 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3061–3070, 2015.
- 720  
721 Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and  
722 accurate method to fool deep neural networks. In *Proc. of the IEEE conference on computer*  
723 *vision and pattern recognition*, pp. 2574–2582, 2016.
- 724  
725 Henrique Morimitsu. Pytorch lightning optical flow. <https://github.com/hmorimitsu/ptlflow>, 2021.
- 726  
727 Henrique Morimitsu, Xiaobin Zhu, Roberto M Cesar, Xiangyang Ji, and Xu-Cheng Yin. RAPID-  
728 Flow: Recurrent adaptable pyramids with iterative decoding for efficient optical flow estimation.  
729 In *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2946–2952, 2024a.
- 730  
731 Henrique Morimitsu, Xiaobin Zhu, Xiangyang Ji, and Xu-Cheng Yin. Recurrent partial kernel net-  
732 work for efficient optical flow estimation. In *AAAI Conference on Artificial Intelligence (AAAI)*,  
733 2024b.
- 734  
735 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
736 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward  
737 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
738 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance  
739 deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035.  
740 Curran Associates, Inc., 2019.
- 741  
742 Sayak Paul and Pin-Yu Chen. Vision transformers are robust learners. In *Proceedings of the AAAI*  
743 *conference on Artificial Intelligence*, volume 36, pp. 2071–2081, 2022.
- 744  
745 Adarsh Prasad. *Towards Robust and Resilient Machine Learning*. PhD thesis, Carnegie Mellon  
746 University, 2022.
- 747  
748 Anurag Ranjan, Joel Janai, Andreas Geiger, and Michael J Black. Attacking optical flow. In *Proc.*  
749 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2404–2413, 2019.
- 750  
751 Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel. Foolbox native: Fast  
752 adversarial attacks to benchmark the robustness of machine learning models in pytorch, tensor-  
753 flow, and jax. *Journal of Open Source Software*, 5(53):2607, 2020. doi: 10.21105/joss.02607.  
754 URL <https://doi.org/10.21105/joss.02607>.
- 755  
756 Stephan R Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *Proceedings*  
757 *of the IEEE international conference on computer vision*, pp. 2213–2222, 2017.
- 758  
759 Benoît Rosa, Valentin Bordoux, and Florent Nageotte. Combining differential kinematics and op-  
760 tical flow for automatic labeling of continuum robots in minimally invasive surgery. *Frontiers in*  
761 *Robotics and AI*, 6:86, 2019.



- 756 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
757 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.  
758 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*  
759 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 760 Christos Sakaridis, Dengxin Dai, and Luc Van Gool. ACDC: The adverse conditions dataset with  
761 correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF*  
762 *International Conference on Computer Vision (ICCV)*, October 2021.
- 763 Erik Scheurer, Jenny Schmalfluss, Alexander Lis, and Andrés Bruhn. Detection defenses: An empty  
764 promise against adversarial patch attacks on optical flow. In *Proc. IEEE/CVF Winter Conference*  
765 *on Applications of Computer Vision (WACV)*, 2024.
- 766 Jenny Schmalfluss, Lukas Mehl, and Andrés Bruhn. Attacking motion estimation with adversarial  
767 snow. In *Proc. ECCV Workshop on Adversarial Robustness in the Real World (AROW)*,  
768 2022a. doi: 10.48550/arXiv.2210.11242. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2210.11242)  
769 [2210.11242](https://doi.org/10.48550/arXiv.2210.11242).
- 770 Jenny Schmalfluss, Philipp Scholze, and Andrés Bruhn. A perturbation-constrained adversarial at-  
771 tack for evaluating the robustness of optical flow. In *Proc. European Conference on Computer*  
772 *Vision (ECCV)*, LNCS 13682, pp. 183–200, 2022b.
- 773 Jenny Schmalfluss, Lukas Mehl, and Andrés Bruhn. Distracting downpour: Adversarial weather  
774 attacks for motion estimation. In *Proc. IEEE/CVF International Conference on Computer Vision*  
775 (*ICCV*), pp. 10106–10116, 2023.
- 776 Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Ad-  
777 versarially robust generalization requires more data. *NIPS*, 31, 2018.
- 778 Simon Schrodi, Tonmoy Saikia, and Thomas Brox. Towards understanding adversarial robustness  
779 of optical flow networks. In *CVPR*, pp. 8916–8924, 2022.
- 780 Hao Shi, Yifan Zhou, Kailun Yang, Xiaoting Yin, and Kaiwei Wang. Csflow: Learning optical flow  
781 via cross strip correlation for autonomous driving. In *IEEE Intelligent Vehicles Symposium (IV)*,  
782 pp. 1851–1858, 2022.
- 783 Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon  
784 See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. VideoFlow: Exploiting temporal cues for  
785 multi-frame optical flow estimation. In *Proc. IEEE/CVF International Conference on Computer*  
786 *Vision (ICCV)*, pp. 12469–12480, 2023a.
- 787 Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei  
788 Qin, Jifeng Dai, and Hongsheng Li. FlowFormer++: Masked cost volume autoencoding for  
789 pretraining optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and*  
790 *Pattern Recognition (CVPR)*, pp. 1599–1610, 2023b.
- 791 Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei  
792 Qin, Jifeng Dai, and Hongsheng Li. FlowFormer++: Masked cost volume autoencoding for  
793 pretraining optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and*  
794 *Pattern Recognition (CVPR)*, pp. 1599–1610, 2023b.
- 795 Xiuchao Sui, Shaohua Li, Xue Geng, Yan Wu, Xinxing Xu, Yong Liu, Rick Goh, and Hongyuan  
796 Zhu. CRAFT: Cross-attentional flow transformer for robust optical flow. In *Proc. IEEE/CVF*  
797 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17602–17611, 2022.
- 798 Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using  
799 pyramid, warping, and cost volume. In *Proc. IEEE/CVF Conference on Computer Vision and*  
800 *Pattern Recognition (CVPR)*, pp. 8934–8943, 2018.
- 801 Shangkun Sun, Yuanqi Chen, Yu Zhu, Guodong Guo, and Ge Li. SKFlow: Learning optical flow  
802 with super kernels. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:11313–  
803 11326, 2022.
- 804 Shiyu Tang, Ruihao Gong, Yan Wang, Aishan Liu, Jiakai Wang, Xinyun Chen, Fengwei Yu, Xiang-  
805 long Liu, Dawn Song, Alan Yuille, Philip H.S. Torr, and Dacheng Tao. Robustart: Benchmarking  
806 robustness on architecture design and training techniques. <https://arxiv.org/pdf/2109.05211.pdf>,  
807 2021.

- 810 Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *Proc.*  
811 *European Conference on Computer Vision (ECCV)*, LNCS 12347, pp. 402–419, 2020.
- 812
- 813 Dustin Tran, Jeremiah Liu, Michael W Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han,  
814 Zi Wang, Zelda Mariet, Huiyi Hu, et al. Plex: Towards reliability using pretrained large model  
815 extensions. *arXiv preprint arXiv:2207.07411*, 2022.
- 816 Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry.  
817 Robustness may be at odds with accuracy. In *ICLR*, 2019.
- 818
- 819 Bo Wang, Yifan Zhang, Jian Li, Yang Yu, Zhenping Sun, Li Liu, and Dewen Hu. Splatflow: Learning  
820 multi-frame optical flow via splatting. *International Journal of Computer Vision*, pp. 1–23, 2024.
- 821 Hengli Wang, Peide Cai, Yuxiang Sun, Lujia Wang, and Ming Liu. Learning interpretable end-to-  
822 end vision-based motion planning for autonomous driving with optical flow distillation. In *2021*  
823 *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 13731–13737. IEEE,  
824 2021.
- 825 Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li. Displacement-  
826 invariant matching cost learning for accurate optical flow estimation. *Advances in Neural Infor-*  
827 *mation Processing Systems (NeurIPS)*, 33:15220–15231, 2020.
- 828
- 829 Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training.  
830 *ArXiv*, abs/2001.03994, 2020.
- 831 J. Wulff, D. J. Butler, G. B. Stanley, and M. J. Black. Lessons and insights from creating a synthetic  
832 optical flow benchmark. In A. Fusiello et al. (Eds.) (ed.), *ECCV Workshop on Unsolved Prob-*  
833 *lems in Optical Flow and Stereo Estimation*, Part II, LNCS 7584, pp. 168–177. Springer-Verlag,  
834 October 2012.
- 835
- 836 Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Seg-  
837 former: Simple and efficient design for semantic segmentation with transformers. *Advances in*  
838 *neural information processing systems*, 34:12077–12090, 2021.
- 839 Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution optical flow  
840 from 1d attention and correlation. In *Proc. IEEE/CVF International Conference on Computer*  
841 *Vision (ICCV)*, pp. 10498–10507, 2021a.
- 842 Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, and Dacheng Tao. GMFlow: Learning  
843 optical flow via global matching. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern*  
844 *Recognition (CVPR)*, pp. 8121–8130, 2022.
- 845
- 846 Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, Fisher Yu, Dacheng Tao, and Andreas  
847 Geiger. Unifying flow, stereo and depth estimation. *IEEE Transactions on Pattern Analysis and*  
848 *Machine Intelligence (TPAMI)*, 45(11):13941–13958, 2023a.
- 849 Jiawei Xu, Zongqing Lu, and Qingmin Liao. LLA-Flow: A lightweight local aggregation on cost  
850 volume for optical flow estimation. In *IEEE International Conference on Image Processing*  
851 *(ICIP)*, pp. 3220–3224, 2023b.
- 852
- 853 Xiaogang Xu, Hengshuang Zhao, and Jiaya Jia. Dynamic divide-and-conquer adversarial training  
854 for robust semantic segmentation. In *IEEE/CVF International Conference on Computer Vision*  
855 *(ICCV)*, pp. 7466–7475, 2021b. doi: 10.1109/ICCV48922.2021.00739.
- 856 Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. *Ad-*  
857 *vances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- 858
- 859 Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for  
860 match density estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recog-*  
861 *nition (CVPR)*, pp. 6044–6053, 2019.
- 862 Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable Flow:  
863 Learning motion cost volumes for optical flow estimation. In *Proc. IEEE/CVF International*  
*Conference on Computer Vision (CVPR)*, pp. 10807–10817, 2021.

864 Zhiyong Zhang, Huaizu Jiang, and Hanumant Singh. NeufLOW: Real-time, high-accuracy optical  
865 flow estimation on robots using edge devices. *arXiv preprint arXiv:2403.10425*, 2024.  
866

867 Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. MaskFlowNet: Asymmetric  
868 feature matching with learnable occlusion mask. In *Proc. IEEE/CVF Conference on Computer  
869 Vision and Pattern Recognition (CVPR)*, pp. 6278–6287, 2020.

870 Shiyu Zhao, Long Zhao, Zhixing Zhang, Enyu Zhou, and Dimitris Metaxas. Global matching with  
871 overlapping attention for optical flow estimation. In *Proc. IEEE/CVF Conference on Computer  
872 Vision and Pattern Recognition (CVPR)*, pp. 17592–17601, 2022.  
873

874 Zihua Zheng, Ni Nie, Zhi Ling, Pengfei Xiong, Jiangyu Liu, Hao Wang, and Jiankun Li. DIP: Deep  
875 inverse patchmatch for high-resolution optical flow. In *Proc. IEEE/CVF Conference on Computer  
876 Vision and Pattern Recognition (CVPR)*, pp. 8925–8934, 2022.  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

# FlowBench: A Robustness Benchmark for Optical Flow Estimation

## Paper #1055 Supplementary Material

### TABLE OF CONTENT

The supplementary material covers the following information:

- Appendix A: [Do Synthetic Corruptions Represent The Real World?: Yes](#)
- Appendix B: [Reasons For Categorizing Methods To Their Respective Families](#)
- Appendix C: Details for the datasets used.
  - Appendix C.1: FlyingThings3D
  - Appendix C.2: KITTI2015
  - Appendix C.3: MPI Sintel
  - Appendix C.4: Spring
- Appendix D: Additional implementation details for the evaluated benchmark.
- Appendix E: In detail description of the attacks.
- Appendix F: A comprehensive look-up table for all the optical flow estimation model weight and datasets pair available in FLOWBENCH and used for evaluating the benchmark.
- Appendix G: In detail explanation of the available functionalities of the FLOWBENCH benchmarking tool and description of the arguments for each function.
- Appendix H: Here we provide additional results from the benchmark evaluated using FlowBench. For all evaluations except Adversarial Weather, the datasets used are KITTI2015, MPI Sintel (clean), and MPI Sintel (final).
  - Appendix H.1.1: Evaluations for all models against FGSM attack under  $\ell_\infty$ -norm bound and  $\ell_2$ -norm bound, as targeted (both targets  $\vec{0}$  and  $\vec{-f}$ ) and non-targeted attack.
  - Appendix H.1.2: Evaluations for all models against BIM attack under  $\ell_\infty$ -norm bound and  $\ell_2$ -norm bound, as targeted (both targets  $\vec{0}$  and  $\vec{-f}$ ) and non-targeted attack, over multiple attack iterations.
  - Appendix H.1.3: Evaluations for all models against PGD attack under  $\ell_\infty$ -norm bound and  $\ell_2$ -norm bound, as targeted (both targets  $\vec{0}$  and  $\vec{-f}$ ) and non-targeted attack, over multiple attack iterations.
  - Appendix H.1.4: Evaluations for all models against CosPGD attack under  $\ell_\infty$ -norm bound and  $\ell_2$ -norm bound, as targeted (both targets  $\vec{0}$  and  $\vec{-f}$ ) and non-targeted attack, over multiple attack iterations.
  - Appendix H.1.5: Evaluations for all models against PCFA attack under  $\ell_2$ -norm bound, as targeted (both targets  $\vec{0}$  and  $\vec{-f}$ ) attack, over multiple attack iterations.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

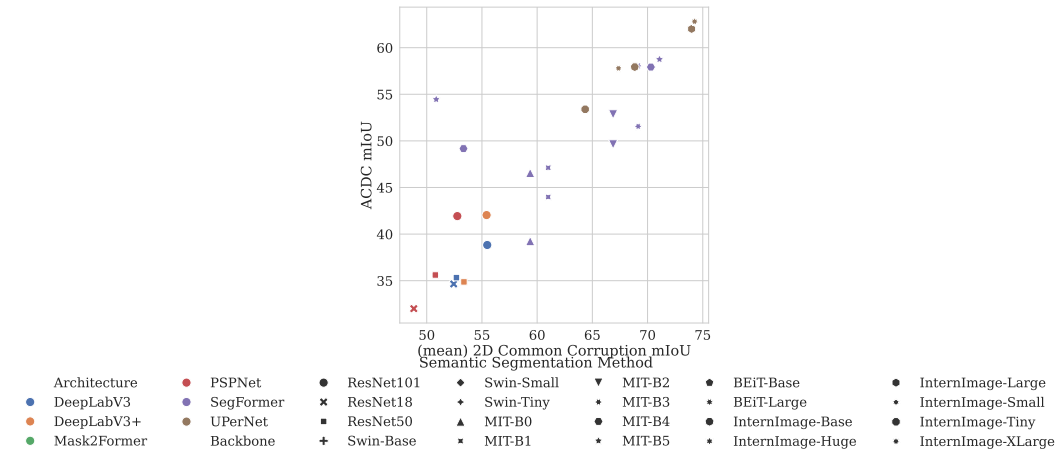


Figure 7: Results from work by Anonymous. Here they find a **very strong positive correlation between mean mIoU over the ACDC evaluation dataset (Sakaridis et al., 2021) and mean mIoU over each 2D Common Corruption (Hendrycks & Dietterich, 2019) over the Cityscapes dataset (Cordts et al., 2016)**. All models were trained using the training subset of the Cityscapes dataset. ACDC is the Adverse Conditions Dataset with Correspondences for Semantic Driving Scene Understanding captured in similar scenes are cityscapes but under four different domains: Day/Night, Rain, Snow, and Fog in the wild. ACDC is a community-used baseline for evaluating the performance of semantic segmentation methods on domain shifts observed in the wild.

- Appendix H.1.6: Evaluations for all models against Adversarial Weather attack, all four conditions: Fog, Rain, Snow, and Sparks, as targeted (both targets  $\vec{0}$  and  $-\vec{f}$ ) and non-targeted attack.
  - Appendix H.2: Evaluations for all models under 2D Common Corruptions and 3D Common Corruptions at severity 3, for KITTI2015, MPI Sintel (clean) and MPI Sintel (final) datasets.
- Appendix I: We share the initial prototype of the future website.
  - Appendix J: We discuss the limitations of FLOWBENCH.
  - Appendix K: We discuss the reproducibility of our evaluations and show that the variance in metrics is extremely low and our analysis comfortably holds under these variances.

## A DO SYNTHETIC CORRUPTIONS REPRESENT THE REAL WORLD?

NEW

In their work Anonymous, they find the correlation between mean mIoU over the ACDC evaluation dataset (Sakaridis et al., 2021) and mean mIoU over each 2D Common Corruption (Hendrycks & Dietterich, 2019) over the Cityscapes dataset (Cordts et al., 2016). We include Figure 7 from their work here for ease of understanding. All models were trained using the training subset of the Cityscapes dataset. ACDC is the Adverse Conditions Dataset with Correspondences for Semantic Driving Scene Understanding captured in similar scenes are cityscapes but under four different domains: Day/Night, Rain, Snow, and Fog in the wild. ACDC is a community-used baseline for evaluating the performance of semantic segmentation methods on domain shifts observed in the wild. They find that there exists a very strong positive correlation between the two. This shows, that **yes, synthetic corruptions can serve as a proxy for the real world**. Unfortunately, a similar “in the wild” captured dataset does not exist for optical flow estimation to evaluate the effect of domain

shifts on the performance of optical flow methods. However, given that for the task of semantic segmentation, we find a very high positive correlation between the performance on real-world corruptions and synthetic corruptions, it is a safe assumption that the same would hold true for optical flow estimation as well. Thus, in this work, we evaluate against synthetic 2D Common Corruptions (Hendrycks & Dietterich, 2019) and synthetic 3D Common Corruptions (Kar et al., 2022).

## B REASONS FOR CATEGORIZING METHODS TO THEIR RESPECTIVE FAMILIES

NEW

Over the years, various modifications have been proposed for DL-based optical flow estimation methods. These can be based on the training strategy used or new architectures. However, barring DIP Zheng et al. (2022) and StarFlow Godet et al. (2021) that appear to have significantly different architectures, all other optimal flow methods can be categorized into four major families: FlowNet (Dosovitskiy et al., 2015), PWC (Pyramid, Wrapping, and Cost Volume) (Sun et al., 2018), RAFT (Teed & Deng, 2020), and FlowFormer (Huang et al., 2022). Following we discuss the reasoning for the categorization of each method considered.

### B.1 FLOWNET FAMILY

NEW

Dosovitskiy et al. (2015) were the first to propose an end-to-end differentiable DL-based architecture for optical flow estimation, FlowNet. Many further works were inspired by FlowNet, making changes to FlowNet to propose novel optical flow estimation methods. These methods include:

- **FlowNet2.0** (Ilg et al., 2017): They improve upon FlowNet by changes to the schedule of training data usage, using a stacked architecture to include the warping of the second image with intermediate optical flow, and a sub-network to focus on small displacements.
- **LiteFlowNet** (Hui et al., 2018): Compared to FlowNet2.0 they use a more effective flow inference approach at each pyramid level through a lightweight cascaded network. They also use a flow regularization layer to ameliorate the issue of outliers and vague flow boundaries by using a feature-driven local convolution, and they use feature warping instead of image warping. They use the same training schedule as FlowNet2.0 but they train their network stage-wise.
- **LiteFlowNet2** (Hui et al., 2020): They improve the accuracy and latency from LiteFlowNet by making minor architectural changes to LiteFlowNet. They follow the training schedule of FlowNet2.0 to some extent and perform stage-wise training.
- **LiteFlowNet3** (Hui & Loy, 2020): They further improve upon the LiteFlowNet2.0 by amending each cost vector using an adaptive modulation before the flow decoding to alleviate the issue of outliers in the cost volume. Additionally, they replace each potentially inaccurately predicted optical flow with an accurate one from a near position through a warping of the flow field. They follow a special training schedule, first training the LiteFlowNet2 modules as mentioned in Hui et al. (2020), and then training the entire architecture again with the LiteFlowNet3 modifications to LiteFlowNet2 following the training protocol from FlowNet2.0.

### B.2 PWC FAMILY

NEW

While still using features from different at different scales, warping, and cost volume, Sun et al. (2018) proposed PWC-Net which with its architectural changes, presented a significant shift in architectures from the traditional FlowNet. Sun et al. (2018) describe, “PWC-Net uses the current optical flow estimate to warp the CNN features of the second image. It then uses the warped features and features of the first image to construct a cost volume, which is processed by a CNN to estimate the optical flow.” This was faster than FlowNet2.0, easier to train, and significantly outperformed it on established benchmarks like KITTI2015 (Menze & Geiger, 2015) and MPI Sintel (Butler et al., 2012). PWC-Net uses a similar training schedule and protocol as FlowNet2.0. Many other works followed PWC-Net either changing the training strategy or making architectural changes to PWC-Net to further improve on i.i.d. performance. These methods include:



- 1080 • **FastFlowNet** (Kong et al., 2021): They replace the dual convolution feature pyramid  
1081 in PWC-Net with the head enhanced pooling pyramid (HEPP) for enhancing the high-  
1082 resolution pyramid feature and reducing model size, then, they propose center dense di-  
1083 lated correlation layer (MFC) for constructing compact cost volume while keeping the  
1084 large search radius. followed by shuffle block decoders (SBD) at each pyramid level to  
1085 regress optical flow with significantly cheaper computation. They follow the same training  
1086 protocol mentioned by FlowNet2.0.
- 1087 • **DICL** (Wang et al., 2020): They improve upon PWC-Net by decoupling the connection  
1088 between 2D displacements and learn the matching costs at each 2D displacement hypothe-  
1089 sis independently, i.e., displacement-invariant cost learning. They apply the same 2D  
1090 convolution-based matching net independently on each 2D displacement hypothesis to  
1091 learn a 4D cost volume and avoid learning a 5D feature volume, thus saving computing  
1092 resources. They use the same training protocol as PWC-Net and FlowNet2.0, and use the  
1093 data augmentations proposed by VCN.
- 1094 • **HD3** (Yin et al., 2019): They adapt a PWC-Net-like architecture for the decomposition of  
1095 the discrete probability distribution instead of the feature representations allowing them to  
1096 learn probabilistic pixel correspondences in both optical flow and stereo matching. They  
1097 decompose the full match density into multiple scales hierarchically and estimate the local  
1098 matching distributions at each scale conditioned on the matching and warping at coarser  
1099 scales. This allows the local distributions to be composed together to form the global  
1100 match density. They essentially follow the same training protocol as FlowNet2.0 while  
1101 omitting some hard examples. Additionally, they use ImageNet1k (Russakovsky et al.,  
1102 2015)-pre-trained weights for their pyramid feature extractor.
- 1103 • **IRR** (Hur & Roth, 2019): They take inspiration from classical energy minimization ap-  
1104 proaches, as well as residual networks to propose an iterative residual refinement, they  
1105 show that their proposed IRR can be combined with both FlowNets and PWC-Net. In our  
1106 work, we consider their adaptation to PWC-Net as that has better i.i.d. performance. They  
1107 use the same training procedure as PWC-Net but additionally set out-of-bound pixels (after  
1108 applying augmentations the same as those in FlowNet2.0) as occluded.
- 1109 • **MaskFlowNet** (Zhao et al., 2020): Zhao et al. (2020) apart their proposed Occlusion-  
1110 Aware Feature Matching Module (OFMM) and Asymmetric Occlusion-Aware Feature  
1111 Matching Module (AsymOFMM) in PWC-Net and consists to two cascaded subnetworks  
1112 for obtaining dual feature pyramids. Their proposed method helps them overcome the am-  
1113 biguity caused due to occlusions in images that induce inaccuracies in the flow fields during  
1114 warping. They use the same training protocol as IRR-PWC-Net. However, first, they train  
1115 the MaskFlowNetS, then keep its weights frozen while training the entire MaskFlowNet.  
1116 They use additional data from KITTI2015 and HD1k dataset (Kondermann et al., 2014) for  
1117 fine-tuning on MPI-Sintel.
- 1118 • **MaskFlowNetS** (Zhao et al., 2020): Proposed as the first stage of MaskFlowNet, Mask-  
1119 FlowNetS inherits the network architecture from PWC-Net, but replaces the feature match-  
1120 ing modules (FMMs) by their proposed AsymOFMMs. They use the same training proce-  
1121 dure as IRR-PWC-Net.
- 1122 • **ScopeFlow** (Bar-Haim & Wolf, 2020): Bar-Haim & Wolf (2020) improve upon IRR-PWC-  
1123 net by improving the data sampling process while testing the regularization and augmenta-  
1124 tions used to mitigate the bias induced by the training protocols. They keep some aspects  
1125 of the training protocols from FlowNet2.0 intact while changing a few like cropping, and  
1126 regularization at different stages of the multi-phase training.
- 1127 • **VCN** Yang & Ramanan (2019): They improve upon the 4D cost volume used by variants  
1128 of the PWC family by proposing volumetric encoder-decoder architectures that efficiently  
1129 capture large receptive fields, multi-channel cost volumes that capture multi-dimensional  
1130 notions of pixel similarities, and separable volumetric filtering that significantly reduces  
1131 computation and parameters while preserving i.i.d. performance. They use a very similar  
1132 training procedure as FlowNet2.0 and PWC-Net, however, with fewer iterations.
- 1133

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

## B.3 RAFT FAMILY

NEW

Teed & Deng (2020) proposed Recurrent All-Pairs Field Transforms (RAFT) to extract per-pixel features to build a multi-scale 4D correlation volume for all pairs of pixels. Here a recurrent unit is used to perform lookups on these correlation volumes. They use additional data and fine-tuning compared to FlowNet2.0. RAFT was a significant architectural change from PWC-Nets and inspired many future works that made modifications to RAFT to further improve i.i.d. performance. These methods include:

- **CCMR** (Jahedi et al., 2024a): They propose adapting RAFT to use attention-based motion grouping concepts for multi-scale optical flow estimation. CCMR first computes global multi-scale context features and then uses them to guide the actual motion grouping. While iterating both steps over all coarse-to-fine scales, Jahedi et al. (2024a) adapt cross-covariance image transformers to allow for an efficient realization while maintaining scale-dependent properties. They use a training procedure similar to MS-RAFT+, after the traditionally followed training procedure of FlowNet2.0, they additionally finetune on a mixed set from KITTI and Viper dataset (Richter et al., 2017).
- **CRAFT** (Sui et al., 2022): CRAFT inherits the flow estimation pipeline of RAFT and revitalizes the correlation volume computation part with two proposed components: the Semantic Smoothing Transformer on the features from the second frame, and a Cross-Frame Attention Layer to compute the correlation volume. Sui et al. (2022) propose that these two components help suppress spurious correlations in the correlation volume. They use the same training procedure as RAFT.
- **CSFlow** (Shi et al., 2022): They propose, “Cross Strip Correlation module (CSC) and Correlation Regression Initialization module (CRI). CSC utilizes a striping operation across the target image and the attended image to encode global context into correlation volumes while maintaining high efficiency. CRI is used to maximally exploit the global context for optical flow initialization”. They take inspiration from RAFT and adapt the multi-layer GRU from the stereo estimation task to optical flow. They follow a training procedure very similar to RAFT.
- **Flow1D** (Xu et al., 2021a): They take inspiration from transformers (Bao et al., 2022) and propose a 1D attention operation that is first applied in the vertical direction of the target image, and then a simple 1D correlation in the horizontal direction of the attended image to achieve 2D correspondence modeling effect. The directions of attention and correlation can also be exchanged, resulting in two 3D cost volumes that are concatenated for optical flow regression, where they adopt RAFT’s framework to estimate the optical flow iteratively. They follow a very similar training procedure to RAFT, however for harder datasets, they use additional data for fine-tuning.
- **GMA** (Jiang et al., 2021a): They adapt an RAFT architecture to include their proposed global motion aggregation (GMA) module, a transformer-based approach to find long-range dependencies between pixels in the first image, and perform global aggregation on the corresponding motion features. This modified RAFT architecture with a GMA further inspired other architectures and works for optical flow estimation. GMA has a very similar training procedure to RAFT.
- **GMFlow** (Xu et al., 2022): They adapt RAFT to identify correspondences in image pairs by comparing their feature similarities. They use transformer-based modules to enhance extracted features, followed by self-attention modules for feature matching and flow propagation. Their feature extraction and feature upsampling modules are identical to RAFT. They follow a very similar training procedure to RAFT.
- **GMFlowNet** (Zhao et al., 2022): They adopt the iterative update operator of RAFT as the optimization step for their proposed GMFlowNet. They use their proposed patch-based overlapping attention (POLA) instead of multi-headed self-attention of transformer blocks to extract large context features to improve the matching step. They follow a very similar training procedure to RAFT.
- **LCV** (Khairi et al., 2024): They propose a lightweight module for learnable cost volume that adds onto RAFT to improve i.i.d. performance. For training, they initialize their learnable cost volume kernels to be identity and directly load the pre-trained weights from

- RAFT, and then they follow a similar training schedule as RAFT but with significantly fewer iterations.
- **LLA-Flow** (Xu et al., 2023b): They propose the local similarity aggregation for 4D cost volume and present lightweight operations to diminish the impact of outliers caused by lack of texture. They apply their module on RAFT to improve i.i.d. performance. They follow a very similar training procedure to RAFT.
  - **MS-RAFT+** (Jahedi et al., 2024b): MS-RAFT adapted RAFT for combining hierarchical concepts at multiple scales. MS-RAFT+ builds on top of MS-RAFT by exploiting an additional finer scale for estimating the flow, which is made feasible using the on-demand cost computation proposed by RAFT. They follow a very particular training schedule which is in parts similar to RAFT, however, due to an overhead of a number of learnable parameters, requires more data and training time.
  - **MatchFlow** (Dong et al., 2023): They propose a different feature matching extractor (FME) to be used with RAFT and GMA module, this proposed FME is pre-trained on a different dataset, which allows for increased i.i.d. performance due to better feature extraction and matching. After incorporating the pre-trained FME, the resultant MatchFlow is trained very similarly to RAFT.
  - **RapidFlow** (Morimitsu et al., 2024a): Inspired by RAFT, Morimitsu et al. (2024a) propose Recurrent Adaptable Pyramids with Iterative Decoding. They propose a recurrent feature encoder that uses a single shared block with efficient 1D layers (NeXt1D) to generate feature pyramids of variable levels. Their decoder is similar to RAFT, with a few changes inspired by SKFlow (Sun et al., 2022). They follow a very similar training procedure to RAFT.
  - **RPKNet** (Morimitsu et al., 2024b): They adapt RAFT to use their proposed Partial Kernel Convolution (PKConv) layers and Separable Large kernels (SLK). PKConv is used to produce variable multi-scale features with a single shared block, while SLK is used to capture large context information with low computational cost. They follow a very similar training procedure to RAFT.
  - **SCV** (Jiang et al., 2021b): They adapt RAFT to use a sparse correlation volume instead of a dense correlation volume. They follow a very similar training procedure to RAFT.
  - **SeparableFlow** (Zhang et al., 2021): They propose a separable cost volume module, a drop-in replacement to RAFT’s correlation cost volumes, that uses non-local aggregation layers to exploit global context cues and prior knowledge, to disambiguate motions in poorly constrained ambiguous regions. They follow a training procedure the same as RAFT.
  - **SKFlow** (Sun et al., 2022): They propose using Super Kernels that allow for larger receptive fields allowing it to recover occluded motions. Finally, they use the non-local GMA module from GMA for optical flow estimations. They follow a similar training procedure to RAFT.
  - **SplatFlow** (Wang et al., 2024): They essentially propose to use splatting for feature matching in architectures like RAFT and GMA. As SplatFlow is proposed to be a multi-frame optical flow estimation method, it requires three frames at a time for training as opposed to the two frames used by RAFT and GMA. We use their modified version with GMA. This requires first loading the pr-trained weights of GMA as proposed by GMA, freezing them, and training the GPU prediction and convex upsampling networks introduced by Splatflow. Then, all parameters are fine-tuned using dataset-specific finetuning procedures very similar to RAFT.
  - **VideoFlow** (Shi et al., 2023a): They propose a multi-frame optical flow estimation method and use the same iterative flow refinement module as other methods in the RAFT family, specifically they use the SKBlocks from SKFlow. For feature extractors, they take inspiration from FlowFormer (Huang et al., 2022) and use ImageNet1k pre-trained Twins-SVT (Chu et al., 2021). They use three and five-image frames during training while following training procedures slightly similar to RAFT.
  - **NeuFlow** (Zhang et al., 2024): Inspired by GMFlow, they use transformer-based blocks to implement global cross-attention, however, they use Flash Attention (Dao et al., 2022)

for slight speed improvements. They use a very similar upsampling module as GMFlow and RAFT. However, to obtain feature maps with finer details, they directly extract features from the original images using a CNN block, instead of using features for matching at the  $\frac{1}{16}$ <sup>th</sup> and  $\frac{1}{8}$ <sup>th</sup> scale like RAFT and GMFlow. They use a very similar training procedure as RAFT.

#### B.4 FLOWFORMER FAMILY

NEW

Proposed by Huang et al. (2022), FlowFormer marks a significant shift in the architecture of optical flow estimation methods compared to the RAFT family.

- **FlowFormer** (Huang et al., 2022): It is a transformer-based neural network architecture for optical flow estimation. Huang et al. (2022) describe, FlowFormer tokenizes the 4D cost volume built from an image pair, encodes the cost tokens into a cost memory with alternate group transformer (AGT) layers in a latent space, and decodes the cost memory via a recurrent transformer decoder with dynamic positional cost queries. The two-stage Twins-SVT (Chu et al., 2021) feature extractor is pre-trained on the ImageNet1k dataset. After that, the training procedure of the entire FlowFormer is similar to RAFT’s training procedure.
- **FlowFormer++** (Shi et al., 2023b): This is built upon FlowFormer to include Masked Cost Volume Autoencoding (MCVA) to improve the i.i.d. performance of FlowFormer by pre-training the cost-volume encoder with a mask encoding strategy proposed by them. FlowFormer++ requires significantly different pre-training, while the training and fine-tuning procedures are similar to RAFT.

### C DATASET DETAILS

FLOWBENCH supports a total of four distinct optical flow datasets. Following, we describe these datasets in detail.

#### C.1 FLYINGTHINGS3D

This is a synthetic dataset proposed by Mayer et al. (2016) largely used for training and evaluation of optical flow estimation methods. This dataset consists of 25000 stereo frames, of everyday objects such as chairs, tables, cars, etc. flying around in 3D trajectories. The idea behind this dataset is to have a large volume of trajectories and random movements rather than focus on a real-world application. In their work, Dosovitskiy et al. (2015) showed models trained on FlyingThings3D can generalize to a certain extent to other datasets.

#### C.2 KITTI2015

Proposed by Menze & Geiger (2015), this dataset is focused on the real-world driving scenario. It contains a total of 400 pairs of image frames, split equally for training and testing. The image frames were captured in the wild while driving around on the streets of various cities. The ground-truth labels were obtained by an automated process.

#### C.3 MPI SINTEL

Proposed by Butler et al. (2012) and Wulff et al. (2012), this dataset is derived from an open-source animated short film and consists of a total of 1064 synthetic frames for training and 564 synthetic frames for testing, both at a resolution of  $1024 \times 436$ . The intention of this dataset is to enforce realism while having a dataset at scale. This dataset is provided as two datasets, which are passes with more transformations and effects on the frames that originally have constant albedo over time, these passes are,

- **MPI Sintel (clean)**: This is the clean pass that adds some realism to the images by adding some spectral effects, like illumination, shadows, and smooth shading.

- MPI Sintel (final): This is the final pass that adds more realism by adding effects such as blur due to depth and camera focus, blur due to motion and atmospheric effects such as snow during snow storms, etc.

#### C.4 SPRING

Similar to MPI Sintel, Mehl et al. (2023) proposed a new dataset and benchmark for optical flow estimation which is much larger than any other dataset before. It consists of frames from the open-source Blender movie “Spring” and consists of 6000 stereo image pairs from 47 sequences with SotA visual effects at full HD resolution ( $1920 \times 1080$  pixels).

## D IMPLEMENTATION DETAILS OF THE BENCHMARK

Following we provide details regarding the experiments done for creating the benchmark used in the analysis.

**Compute Resources.** Most experiments were done on a single 40 GB NVIDIA Tesla V100 GPU each, however, MS-RAFT+, FlowFormer, and FlowFormer++ are more compute-intensive, and thus 80GB NVIDIA A100 GPUs or NVIDIA H100 were used for these models, a single GPU for each experiment.

**Datasets Used.** Performing adversarial attacks and OOD robustness evaluations are very expensive and compute-intensive. Thus, performing evaluation using all model-dataset pairs is not possible given the limited computing resources at our disposal. Thus, for the benchmark, we only use KITTI2015, MPI Sintel (clean), and MPI Sintel (final) as these are the most commonly used datasets for evaluation (Ilg et al., 2017; Huang et al., 2022; Schmalfluss et al., 2022b; Schrodi et al., 2022; Agnihotri et al., 2024).

**Metrics Calculation.** In Sec. 4 we introduce three new metrics for better understanding our analysis, given the large scale of the benchmark created. For calculating TARE and NARE values we used BIM, PGD, and CosPGD attack with step size  $\alpha=0.01$ , perturbation budget  $\epsilon = \frac{8}{255}$  under the  $\ell_\infty$ -norm bound, as targeted and non-targeted attacks respectively. We use  $\ell_\infty$ -norm bound as we observe in Appendix H that there is a high correlation between the performance of optical flow estimation methods when attacked using  $\ell_\infty$ -norm bounded attacks and  $\ell_2$ -norm bounded attacks. We use 20 attack iterations for calculating TARE and NARE as we observe in *Appendix H*, that at a lower number of iterations, the gap in performance of different optical flow estimation methods is small, thus an in-depth analysis would be difficult, and we do not go beyond 20 attack iterations as computing each attack step for an adversarial attack is very expensive, and as shown by Agnihotri et al. (2024) and Schmalfluss et al. (2022b), 20 iterations are enough to optimize an attack to truly understand the performance of the attacked method. For calculating GAE, we use all 15 2D Common Corruptions: ‘Gaussian Noise’, ‘Shot Noise’, ‘Impulse Noise’, ‘Defocus Blur’, ‘Frosted Glass Blur’, ‘Motion Blur’, ‘Zoom Blur’, ‘Snow’, ‘Frost’, ‘Fog’, ‘Brightness’, ‘Contrast’, ‘Elastic Transform’, ‘Pixelate’, ‘JPEG Compression’, and eight 3D Common Corruptions: ‘Color Quantization’, ‘Far Focus’, ‘Fog 3D’, ‘ISO Noise’, ‘Low Light’, ‘Near Focus’, ‘XY Motion Blur’, and ‘Z Motion Blur’. All the common corruptions are at severity 3. Kar et al. (2022) offers more 3D Common Corruptions, however computing them is resource intensive. Thus, given our limited resources and an overlap in the corruptions between 2D Common Corruptions and 3D Common Corruptions, we focus on generating 3D Common Corruptions that might be unique from their 2D counterpart, require fewer sources to generate, and are interesting from an optical flow estimation perspective.

**Calculating the EPE.** *EPE* is the Euclidean distance between the two vectors, where one vector is the predicted flow by the optical flow estimation method and the other vector is the ground truth in case of i.i.d. performance evaluations, non-targeted attacks evaluations, and OOD robustness evaluations, while it is the target flow vector, in case of targeted attacks. For each dataset, the *EPE* value is calculated over all the samples of the evaluation set of the respective dataset and then the mean *EPE* value is used as the mean-*EPE* of the respective method over the respective dataset.

NEW



**Other Metrics.** Apart from EPE, FLOWBENCH also enables calculating a lot of other interesting metrics, such as  $\ell_0$ ,  $\ell_2$ ,  $\ell_\infty$ , distance between the perturbations of each image before and after a threat. Apart from these, in all scenarios, we also capture the outlier error, 1-px error, 3-px error, 5-px error and cosine distance between two vectors. These vectors are the same as that in the case of EPE calculations. We limited the analysis in this work to use EPE, since it is the most commonly used metric for evaluation, moreover, most works on optical flow estimation (Agnihotri et al., 2024; Schmalfluss et al., 2022b; Schrodi et al., 2022; Teed & Deng, 2020; Jahedi et al., 2024b) show a very high correlation between performance evaluations using different metrics.

**Models Used.** All available checkpoints, as shown in Tab. 1 for MPI Sintel and KITTI2015 dataset were used for creating the benchmark, except the following four models: Separableflow, SCV, VCN, Unimatch as due to special operations used in these models, they required specific libraries which were creating conflicts with all the others models, and as most of these models are very old and do not have performance close to SotA performance, we did not include them.

**Adversarial Weather** For generating adversarial weather attacks, we followed the implementation proposed by Schmalfluss et al. (2023). However, generating this attack is highly compute-intensive, and thus doing so for all models was not possible. Thus, based on the performance and reliability of all the models, we identified a few (eight) interesting models and only attacked them using the four different attacks curtailed within adversarial weather. This was done to demonstrate the capability of FLOWBENCH to perform this attack. The following are the specifications for the weather attacks:

- Adversarial Weather: **Snow** (random snowflakes)
  - Number of Particles: 3000
  - Number of optimization steps: 750
- Adversarial Weather: **Rain** (rain streaks of length 0.15 with motion blur )
  - Number of Particles: 20
  - Number of optimization steps: 750
- Adversarial Weather: **Fog** (random large less opacity particles)
  - Number of Particles: 60
  - Number of optimization steps: 750
- Adversarial Weather: **Sparks** (random red sparks)
  - Number of Particles: 3000
  - Number of optimization steps: 750

Please note, that these specifications are identical to the optimal ones proposed by Schmalfluss et al. (2023).

## E DESCRIPTION OF FLOWBENCH

Following, we describe the benchmarking tool, FLOWBENCH. It is built using pltflow (Morimitsu, 2021), and supports 36 unique architectures and 4 distinct datasets, namely FlyingThings3D (Mayer et al., 2016), KITTI2015 (Menze & Geiger, 2015), MPI Sintel (Butler et al., 2012) (clean and final) and Spring (Mehl et al., 2023) datasets (please refer Appendix C for additional details on the datasets). It enables training and evaluations on all aforementioned datasets including evaluations using SotA adversarial attacks such as CosPGD (Agnihotri et al., 2024) and PCFA (Schmalfluss et al., 2022b), Adversarial weather (Schmalfluss et al., 2023), and other commonly used adversarial attacks like BIM (Kurakin et al., 2018), PGD (Kurakin et al., 2017), FGSM (Goodfellow et al., 2015), under various lipshitz ( $l_p$ ) norm bounds.

Additionally, it enables evaluations for Out-of-Distribution (OOD) robustness by corrupting the inference samples using 2D Common Corruptions (Hendrycks & Dietterich, 2019) and 3D Common Corruptions (Kar et al., 2022).



We follow the nomenclature set by RobustBench (Croce et al., 2021) and use “threat\_model” to define the kind of evaluation to be performed. When “threat\_model” is defined to be “None”, the evaluation is performed on unperturbed and unaltered images, if the “threat\_model” is defined to be an adversarial attack, for example “PGD”, “CosPGD” or “PCFA”, then FLOWBENCH performs an adversarial attack using the user-defined parameters. We elaborate on this in Appendix E.1. Whereas, if “threat\_model” is defined to be “2DCommonCorruptions” or “3DCommonCorruptions”, the FLOWBENCH performs evaluations after perturbing the images with 2D Common Corruptions and 3D Common Corruptions respectively. We elaborate on this in Appendix E.2.

If the queried evaluation already exists in the benchmark provided by this work, then FLOWBENCH simply retrieves the evaluations, thus saving computation.

## E.1 ADVERSARIAL ATTACKS

FLOWBENCH enables the use of all the attacks mentioned in Sec. 2.3 to help users better study the reliability of their optical flow methods. We choose to specifically include these white-box adversarial attacks as they either serve as the common benchmark for adversarial attacks in classification literature (FGSM, BIM, PGD, APGD) or they are unique attacks proposed specifically for pixel-wise prediction tasks (CosPGD) and optical flow estimation (PCFA and Adversarial Weather). These attacks can either be *Non-targeted* which are designed to simply fool the model into making incorrect predictions, irrespective of what the model eventually predicts, or can be *Targeted*, where the model is fooled to make a certain prediction. Most attacks can be, designed to be either Targeted or Non-targeted, these include, FGSM, BIM, PGD, APGD, CosPGD and Adversarial Weather. However, by design, some attacks are limited to being only one of the two, for example, PCFA which is a targeted attack. Following, we discuss these attacks in detail and highlight their key differences.

**FGSM.** Assuming a non-targeted attack, given a model  $f_\theta$  and an unperturbed input sample  $\mathbf{X}^{\text{clean}}$  and ground truth label  $\mathbf{Y}$ , FGSM attack adds noise  $\delta$  to  $\mathbf{X}^{\text{clean}}$  as follows,

$$\mathbf{X}^{\text{adv}} = \mathbf{X}^{\text{clean}} + \alpha \cdot \text{sign} \nabla_{\mathbf{X}^{\text{clean}}} L(f_\theta(\mathbf{X}^{\text{clean}}), \mathbf{Y}), \quad (1)$$

$$\delta = \phi^\epsilon(\mathbf{X}^{\text{adv}} - \mathbf{X}^{\text{clean}}), \quad (2)$$

$$\mathbf{X}^{\text{adv}} = \phi^r(\mathbf{X}^{\text{clean}} + \delta). \quad (3)$$

Here,  $L(\cdot)$  is the loss function (differentiable at least once) which calculates the loss between the model prediction and ground truth,  $\mathbf{Y}$ .  $\alpha$  is a small value of  $\epsilon$  that decides the size of the step to be taken in the direction of the gradient of the loss w.r.t. the input image, which leads to the input sample being perturbed such that the loss increases.  $\mathbf{X}^{\text{adv}}$  is the adversarial sample obtained after perturbing  $\mathbf{X}^{\text{clean}}$ . To make sure that the perturbed sample is semantically indistinguishable from the unperturbed clean sample to the human eye, steps from Eq. (2) and Eq. (3) are performed. Here, function  $\phi^\epsilon$  is clipping the  $\delta$  in  $\epsilon$ -ball for  $\ell_\infty$ -norm bounded attacks or the  $\epsilon$ -projection in other  $l_p$ -norm bounded attacks, complying with the  $\ell_\infty$ -norm or other  $l_p$ -norm constraints, respectively. While function  $\phi^r$  clips the perturbed sample ensuring that it is still within the valid input space. FGSM, as proposed, is a single step attack. For targeted attacks,  $\mathbf{Y}$  is the target and  $\alpha$  is multiplied by -1 so that a step is taken to minimize the loss between the model’s prediction and the target prediction.

**BIM.** This is the direct extension of FGSM into an iterative attack method. In FGSM,  $\mathbf{X}^{\text{clean}}$  was perturbed just once. While in BIM,  $\mathbf{X}^{\text{clean}}$  is perturbed iteratively for time steps  $t \in [0, T]$ , such that  $t \in \mathbb{Z}^+$ , where  $T$  are the total number of permissible attack iterations. This changes the steps of the attack from FGSM to the following,

$$\mathbf{X}^{\text{adv}_{t+1}} = \mathbf{X}^{\text{adv}_t} + \alpha \cdot \text{sign} \nabla_{\mathbf{X}^{\text{adv}_t}} L(f_\theta(\mathbf{X}^{\text{adv}_t}), \mathbf{Y}), \quad (4)$$

$$\delta = \phi^\epsilon(\mathbf{X}^{\text{adv}_{t+1}} - \mathbf{X}^{\text{clean}}), \quad (5)$$

$$\mathbf{X}^{\text{adv}_{t+1}} = \phi^r(\mathbf{X}^{\text{clean}} + \delta). \quad (6)$$

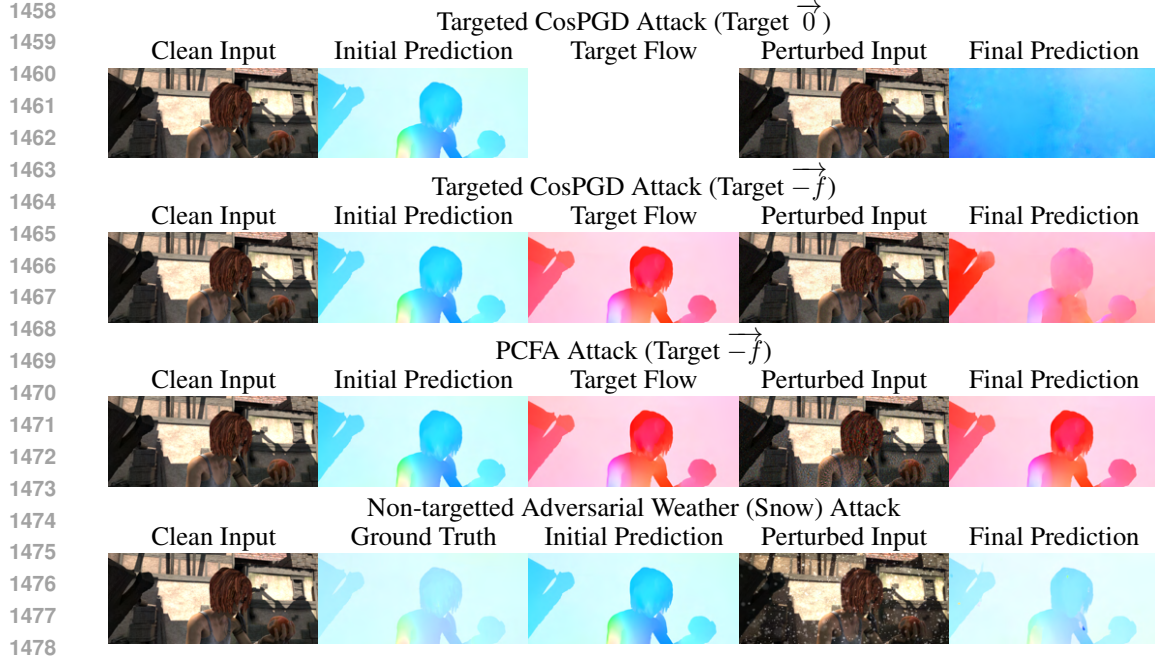


Figure 8: Examples of MPI Sintel images perturbed by the mentioned adversarial attacks and the optical flow predictions using FlowFormer++. These examples are intended to show the versatility of FLOWBENCH.

Here, at  $t=0$ ,  $\mathbf{X}^{\text{adv}_t} = \mathbf{X}^{\text{clean}}$ .

**PGD.** Since in BIM, the initial prediction always started from  $\mathbf{X}^{\text{clean}}$ , the attack required a significant amount of steps to optimize the adversarial noise and yet it was not guaranteed that in the permissible  $\epsilon$ -bound,  $\mathbf{X}^{\text{adv}_{t+1}}$  was far from  $\mathbf{X}^{\text{clean}}$ . Thus, PGD proposed introducing stochasticity to ensure random starting points for attack optimization. They achieved this by perturbing  $\mathbf{X}^{\text{clean}}$  with  $\mathcal{U}(-\epsilon, \epsilon)$ , a uniform distribution in  $[-\epsilon, \epsilon]$ , before making the first prediction, such that, at  $t=0$

$$\mathbf{X}^{\text{adv}_t} = \phi^r(\mathbf{X}^{\text{clean}} + \mathcal{U}(-\epsilon, \epsilon)). \quad (7)$$

**APGD.** Auto-PGD is an effective extension to the PGD attack that effectively scales the step size  $\alpha$  over attack iterations considering the compute budget and the success rate of the attack.

**CosPGD.** All previously discussed attacks were proposed for the image classification task. Here, the input sample is a 2D image of resolution  $H \times W$ , where  $H$  and  $W$  are the height and width of the spatial resolution of the sample, respectively. Pixel-wise information is inconsequential for image classification. This led to the pixel-wise loss  $\mathcal{L}(\cdot)$  being aggregated to  $L(\cdot)$ , as follows,

$$L(f_\theta(\mathbf{X}^{\text{adv}_t}), \mathbf{Y}) = \frac{1}{H \times W} \sum_{i \in H \times W} \mathcal{L}(f_\theta(\mathbf{X}^{\text{adv}_t})_i, \mathbf{Y}_i). \quad (8)$$

This aggregation of  $\mathcal{L}(\cdot)$  fails to account for pixel-wise information available in tasks other than image classification, such as pixel-wise prediction tasks like Optical Flow estimation. Thus, in their work Agnihotri et al. (2024) propose an effective extension of the PGD attack that takes pixel-wise information into account by scaling  $\mathcal{L}(\cdot)$  by the alignment between the distribution of the predictions and the distributions of  $\mathbf{Y}$  before aggregating leading to a better-optimized attack, modifying Eq. (4) as follows,

$$\mathbf{X}^{\text{adv}_{t+1}} = \mathbf{X}^{\text{adv}_t} + \alpha \cdot \text{sign} \nabla_{\mathbf{X}^{\text{adv}_t}} \sum_{i \in H \times W} \cos(\psi(f_\theta(\mathbf{X}^{\text{adv}_t})_i), \Psi(\mathbf{Y}_i)) \cdot \mathcal{L}(f_\theta(\mathbf{X}^{\text{adv}_t})_i, \mathbf{Y}_i). \quad (9)$$

Where, functions  $\psi$  and  $\Psi$  are used to obtain the distribution over the predictions and  $\mathbf{Y}_i$ , respectively, and the function  $\text{cos}$  calculates the cosine similarity between the two distributions. CosPGD is the unified SotA adversarial attack for pixel-wise prediction tasks.

**PCFA.** Recently proposed by Schmalfluss et al. (2022b), is the SotA targeted adversarial attack specifically designed for optical flow estimation. It optimizes the input perturbation  $\delta = \mathbf{X}^{\text{adv}_t} - \mathbf{X}^{\text{clean}}$  within a given  $l_2$  bound to obtain a given target flow  $\mathbf{Y}^{\text{targ}}$ . Mathematically, PCFA transforms the constrained optimization problem to find the most destructive perturbation under an  $l_2$  constraint  $\varepsilon_2$  into an unconstrained optimization problem by adding a term that penalizes deviations from the  $l_2$  constraint:

$$\mathbf{X}^{\text{adv}_{t+1}} = \mathbf{X}^{\text{adv}_t} + \underset{\hat{\delta}}{\text{argmin}}(\mathcal{L}(f_\theta(\mathbf{X}^{\text{adv}_t}), \mathbf{Y}^{\text{targ}}) + \mu \cdot \text{ReLU}(\|\hat{\delta}\|_2^2 - (\varepsilon_2 \sqrt{2 \times H \times W})^2)) \quad (10)$$

Here,  $\mathcal{L}(\cdot)$  is a generic loss function, like EPE or cosine distance. The penalty scaling parameter  $\mu$  influences how severely deviations from the per-pixel  $l_2$  bound  $\varepsilon_2$  are penalized. The optimization problem  $\text{argmin}(\cdot)$  is solved with an L-BFGS optimizer.

**Adversarial Weather.** Unlike the previous attacks which introduced per-pixel modifications, adversarial weather Schmalfluss et al. (2023; 2022a) attacks optical flow methods through optimizing the motion trajectories of rendered weather particles  $\mathcal{P}$  like snow flakes, rain drops or fog clouds. The particle trajectories are modelled as positions  $\mathbf{P} = \{\mathbf{P}_1, \mathbf{P}_2\}$  in the two frames  $I_1, I_2$ . Consequently,  $\mathbf{X}^{\text{adv}}(\mathbf{P})$  is generated by differentially rendering the particles with their respective 3D positions to the 2D images. The update step optimizes the particle positions to achieve a certain target flow  $\mathbf{Y}^{\text{targ}}$  while simultaneously limiting the position offset size  $\delta_{\mathbf{P}^t} = \mathbf{P}^{\text{init}} - \mathbf{P}^t$ :

$$\mathbf{X}^{\text{adv}}(\mathbf{P}^{t+1}) = \mathbf{X}^{\text{adv}}\left(\mathbf{P}^t + \alpha \cdot \nabla_{\mathbf{P}^t}\left(\text{EPE}(f_\theta(\mathbf{X}^{\text{adv}}(\mathbf{P}^t)), \mathbf{Y}^{\text{targ}}) + \sum_{I \in \{1,2\}} \frac{\beta_I}{|\mathcal{P}|} \sum_{j \in \mathcal{P}} \frac{\|\delta_{\mathbf{P}_I^t}^j\|_2^2}{d_I^j}\right)\right). \quad (11)$$

Here,  $\beta$  balances the two optimization goals of reaching the target flow and limiting trajectory offsets. The allowed trajectory offsets are further scaled with the particle depth  $d$  in the scene, to generate visually pleasing results.

Fig. 8, shows adversarial examples created using the SotA attacks and how they affect the model predictions.

## E.2 OUT-OF-DISTRIBUTION ROBUSTNESS

While adversarial attacks help explore vulnerabilities of inefficient feature representations learned by a model, another important aspect of reliability is generalization ability. Especially, generalization to previously unseen samples or samples from significantly shifted distributions compared to the distribution of the samples seen while learning model parameters. As one cannot cover all possible scenarios during model training, a certain degree of generalization ability is expected from models. However, multiple works (Hendrycks & Dietterich, 2019; Kar et al., 2022; Hoffmann et al., 2021) showed that models are surprisingly less robust to distribution shifts, even those that can be caused by commonly occurring phenomena such as weather changes, lighting changes, etc. This makes the study of Out-of-Distribution (OOD) robustness an interesting avenue for research. Thus, to facilitate the study of robustness to such commonly occurring corruptions, FLOWBENCH enables evaluating against prominent image corruption methods. Following, we describe these methods in detail.

**2D Common Corruptions.** Hendrycks & Dietterich (2019) propose introducing distribution shift in the input samples by perturbing images with a total of 15 synthetic corruptions that could occur in the real world. These corruptions include weather phenomena such as fog, and frost, digital corruptions such as jpeg compression, pixelation, and different kinds of blurs like motion, and zoom blur, and noise corruptions such as Gaussian and shot noise amongst others corruption types. Each of these corruptions can perturb the image at 5 different severity levels between 1 and 5. The final performance of the model is the mean of the model’s performance on all the corruptions, such that

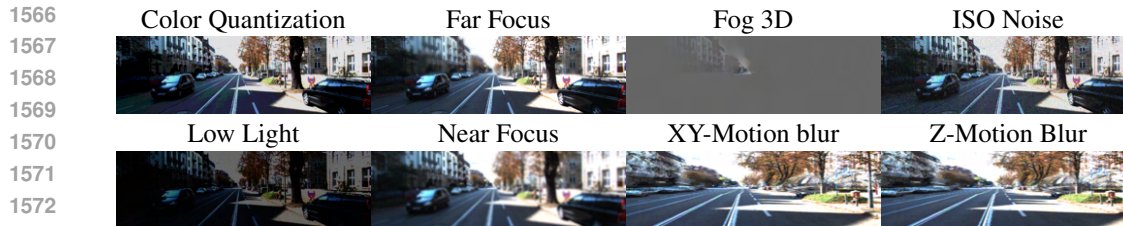


Figure 9: Examples of images from KITTI2015 corrupted using 3D Common Corruptions for evaluation of OOD robustness.

every corruption is used to perturb each image in the evaluation dataset. Since these corruptions are applied to a 2D image, they are collectively termed 2D Common Corruptions.

**3D Common Corruptions.** Since the real world is 3D, Kar et al. (2022) extend 2D Common Corruptions to formulate more realistic-looking corruptions by leveraging depth information (synthetic depth information when real depth is not readily available) and luminescence angles. They name these image corruptions as 3D Common Corruptions. Fig. 9, shows examples of KITTI2015 images corrupted using 3D Common Corruptions.

## F MODEL ZOO

The trained checkpoints for all models available in FLOWBENCH can be obtained using the following lines of code:

```
from flowbench.evals import load_model
model = load_model(model_name='RAFT', dataset='KITTI2015')
```

Each model checkpoint can be retrieved with the pair of ‘model\_name’, the name of the model, and ‘dataset’, the dataset for which the checkpoint was last fine-tuned. In Table 1, we provide a comprehensive look-up table for all ‘model\_name’ and ‘dataset’ pairs for which trained checkpoints are available in FlowBench.

NEW

## G FLOWBENCH USAGE DETAILS

Following we provide a detailed description of the evaluation functions and their arguments provided in FlowBench.

### G.1 ADVERSARIAL ATTACKS

To evaluate a model for a given dataset, on an attack, the following lines of code are required.

```
from flowbench.evals import evaluate
model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
                          threat_model='CosPGD', iterations=20, alpha=0.01,
                          epsilon=8/255, lp_norm='Linf', targeted=True,
                          optim_wrt='ground_truth', retrieve_existing=True)
```

The argument description is as follows:

- ‘model\_name’ is the name of the optical flow estimation method to be used, given as a string.
- ‘dataset’ is the name of the dataset to be used also given as a string.
- ‘threat\_model’ is the name of the adversarial attack to be used, given as a string.
- ‘iterations’ are the number of attack iterations, given as an integer.

Table 1: Overview of all available model checkpoints (model X, trained for dataset Y) in FLOWBENCH.

Model	Dataset			Method Family	Time
	FlyingThings3D (Mayer et al., 2016)	KITTI2015 (Menze & Geiger, 2015)	MPI Sintel (Butler et al., 2012)		
CCMR (Jahedi et al., 2024a)	X	✓	✓	RAFT	January 2024
CRAFT (Sui et al., 2022)	✓	✓	✓	RAFT	March 2022
CSFlow (Shi et al., 2022)	✓	✓	X	RAFT	February 2022
DICL (Wang et al., 2020)	✓	✓	✓	PWC	October 2020
DIP (Zheng et al., 2022)	✓	✓	✓	Deep Inverse Patchmatch	April 2022
FastFlowNet (Kong et al., 2021)	✓	✓	✓	PWC	March 2021
Flow1D (Xu et al., 2021a)	✓	✓	✓	RAFT	April 2021
FlowFormer (Huang et al., 2022)	✓	✓	✓	FlowFormer	March 2022
FlowFormer++ (Shi et al., 2023b)	✓	✓	✓	FlowFormer	March 2023
FlowNet2.0 (Ilg et al., 2017)	✓	X	X	FlowNet	December 2016
GMA (Jiang et al., 2021a)	✓	✓	✓	RAFT	April 2021
GMFlow (Xu et al., 2022)	✓	✓	✓	RAFT	November 2021
GMFlowNet (Zhao et al., 2022)	✓	✓	✓	RAFT	March 2022
HD3 (Yin et al., 2019)	✓	✓	✓	PWC	December 2018
IRR (Hur & Roth, 2019)	✓	✓	✓	PWC	April 2019
LCV (Khairi et al., 2024)	✓	X	X	RAFT	July 2020
LiteFlowNet (Hui et al., 2018)	✓	✓	✓	FlowNet	May 2018
LiteFlowNet2 (Hui et al., 2020)	X	✓	✓	FlowNet	February 2020
LiteFlowNet3 (Hui & Loy, 2020)	X	✓	✓	FlowNet	July 2020
LLA-Flow (Xu et al., 2023b)	✓	✓	✓	RAFT	April 2023
MaskFlowNetS (Zhao et al., 2020)	✓	X	✓	PWC	March 2023
MaskFlowNet (Zhao et al., 2020)	X	✓	✓	PWC	March 2023
MS-RAFT+ (Jahedi et al., 2024b)	✓	✓	✓	RAFT	October 2022
MatchFlow (Dong et al., 2023)	✓	✓	✓	RAFT	March 2023
NeuFlow (Zhang et al., 2024)	X	X	✓	FlowNet	March 2024
PWC-Net (Sun et al., 2018)	✓	X	✓	PWC	September 2017
RapidFlow (Morimitsu et al., 2024a)	✓	✓	✓	RAFT	May 2024
RAFT (Teed & Deng, 2020)	✓	✓	✓	RAFT	March 2020
RPKNet (Morimitsu et al., 2024b)	✓	✓	✓	RAFT	March 2024
ScopeFlow (Bar-Haim & Wolf, 2020)	✓	✓	✓	PWC	February 2020
SCV (Jiang et al., 2021b)	✓	✓	✓	RAFT	April 2021
SeparableFlow (Zhang et al., 2021)	✓	✓	✓	RAFT	October 2021
SKFlow (Sun et al., 2022)	✓	✓	✓	RAFT	November 2022
SplatFlow (Wang et al., 2024)	X	✓	X	RAFT	January, 2024
StarFlow (Godet et al., 2021)	✓	✓	✓	SSRFlow	July 2020
Unimatch (Xu et al., 2023a)	✓	X	X	RAFT	November 2022
VCN (Yang & Ramanan, 2019)	✓	X	X	PWC	December 2019
VideoFlow (Shi et al., 2023a)	✓	✓	✓	RAFT	March 2023

- ‘epsilon’ is the permissible perturbation budget  $\epsilon$  given a floating point (float).
- ‘alpha’ is the step size of the attack,  $\alpha$ , given as a floating point (float).
- ‘lp\_norm’ is the Lipschitz continuity norm ( $l_p$ -norm) to be used for bounding the perturbation, possible options are ‘Linf’ and ‘L2’ given as a string.
- ‘targeted’ is a boolean flag that decides if the attack must be targeted or not. If targeted=‘True’, then by default the target is  $\vec{0}$ , passed as target=‘zero’, this can be changed to negative of the initial flow by passing target=‘negative’.
- ‘optim\_wrt’ decides wrt what attack should be optimized, available choices are ‘ground\_truth’ and ‘initial\_flow’ as string. Please note, this only works well with attacks that utilize Eq. (7).
- ‘retrieve\_existing’ is a boolean flag, which when set to ‘True’ will retrieve the evaluation from the benchmark if the queried evaluation exists in the benchmark provided by this work, else FLOWBENCH will perform the evaluation. If the ‘retrieve\_existing’ boolean flag is set to ‘False’ then FLOWBENCH will perform the evaluation even if the queried evaluation exists in the provided benchmark.

## G.2 ADVERSARIAL WEATHER

As an attack, adversarial weather works slightly different compared to other adversarial attacks, thus we additionally mention the commands for using adversarial weather.

```

from flowbench.evals import evaluate
model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
                          threat_model='Adversarial_Weather', weather='snow',
                          num_particles=10000, targeted=True,
                          retrieve_existing=True)

```

The argument description is as follows:

- 1674 • ‘model\_name’ is the name of the optical flow estimation method to be used, given as a
- 1675 string.
- 1676 • ‘dataset’ is the name of the dataset to be used also given as a string.
- 1677 • ‘threat\_model’ is the name of the adversarial attack to be used, given as a string.
- 1678 • ‘weather’ is the name of the weather condition in adversarial weather attack to be used,
- 1679 given as a string, options include ‘snow’, ‘fog’, ‘rain’ and ‘sparks’.
- 1680 • ‘num\_particles’ is the number of particles per frame to be used, given as an integer.
- 1681 • ‘targeted’ is a boolean flag that decides if the attack must be targeted or not. If tar-
- 1682 geted=‘True’, then by default the target is  $\vec{0}$ , passed as target=‘zero’, this can be changed
- 1683 to negative of the initial flow by passing target=‘negative’.
- 1684 • ‘optim\_wrt’ decides wrt what attack should be optimized, available choices are
- 1685 ‘ground\_truth’ and ‘initial\_flow’ as string. Please note, this only works well with attacks
- 1686 that utilize Eq. (7).
- 1687 • ‘retrieve\_existing’ is a boolean flag, which when set to ‘True’ will retrieve the evaluation
- 1688 from the benchmark if the queried evaluation exists in the benchmark provided by this
- 1689 work, else FLOWBENCH will perform the evaluation. If the ‘retrieve\_existing’ boolean
- 1690 flag is set to ‘False’ then FLOWBENCH will perform the evaluation even if the queried
- 1691 evaluation exists in the provided benchmark.
- 1692
- 1693

### 1694 G.3 2D COMMON CORRUPTIONS

1695 To evaluate a model for a given dataset, with 2D Common Corruptions, the following lines of code

1696 are required.

```
1698 from flowbench.evals import evaluate
1699 model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
1700                         threat_model='2DCommonCorruption',
1701                         severity=3, retrieve_existing=True)
1702
```

1703 The argument description is as follows:

- 1704 • ‘model\_name’ is the name of the optical flow estimation method to be used, given as a
- 1705 string.
- 1706 • ‘dataset’ is the name of the dataset to be used also given as a string.
- 1707 • ‘threat\_model’ is the name of the common corruption to be used, given as a string.
- 1708 • ‘severity’ is the severity of the corruption, given as an integer between 1 and 5 (both inclu-
- 1709 sive).
- 1710 • ‘retrieve\_existing’ is a boolean flag, which when set to ‘True’ will retrieve the evaluation
- 1711 from the benchmark if the queried evaluation exists in the benchmark provided by this
- 1712 work, else FLOWBENCH will perform the evaluation. If the ‘retrieve\_existing’ boolean
- 1713 flag is set to ‘False’ then FLOWBENCH will perform the evaluation even if the queried
- 1714 evaluation exists in the provided benchmark.
- 1715
- 1716

1717 FLOWBENCH supports the following 2D Common Corruption: ‘gaussian\_noise’, ‘shot\_noise’, ‘im-

1718 pulse\_noise’, ‘defocus\_blur’, ‘frosted\_glass\_blur’, ‘motion\_blur’, ‘zoom\_blur’, ‘snow’, ‘frost’, ‘fog’,

1719 ‘brightness’, ‘contrast’, ‘elastic’, ‘pixelate’, ‘jpeg’. For the evaluation, FLOWBENCH will evaluate

1720 the model on the validation images from the respective dataset corrupted using each of the afore-

1721 mentioned corruptions for the given severity, and then report the mean performance over all of them.

### 1722 G.4 3D COMMON CORRUPTIONS

1723 To evaluate a model for a given dataset, with 3D Common Corruptions, the following lines of code

1724 are required.

```
1725 from flowbench.evals import evaluate
1726 model, results = evaluate(model_name='RAFT', dataset='KITTI2015',
1727
```





1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

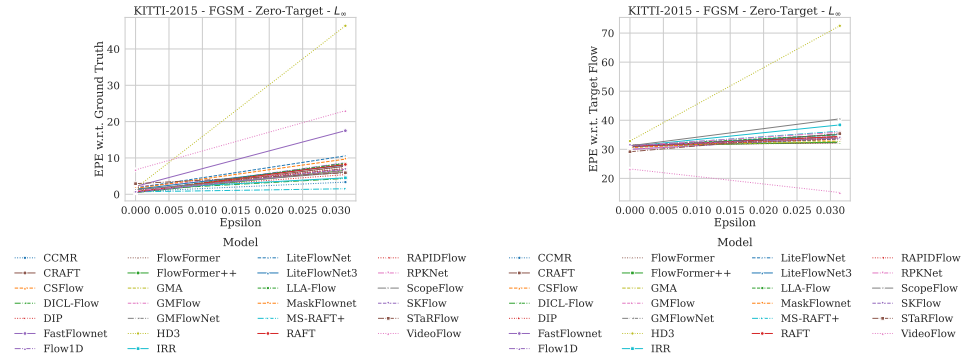


Figure 11: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

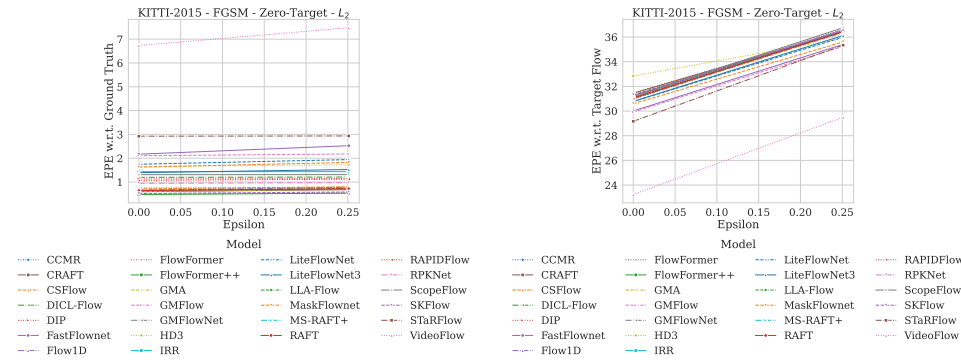


Figure 12: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

### H.1.2 BIM ATTACK

Here we report the evaluations using BIM attack, both as targeted (both targets:  $\vec{0}$  and  $-\vec{f}$ ) and non-targeted attacks optimized under the  $\ell_\infty$ -norm bound and the  $\ell_2$ -norm bound over multiple attack iterations. For  $\ell_\infty$ -norm bound, perturbation budget  $\epsilon = \frac{8}{255}$ , and step size  $\alpha=0.01$ , while for  $\ell_2$ -norm bound, perturbation budget  $\epsilon = \frac{64}{255}$  and step size  $\alpha=0.1$ . Attack evaluations include Fig. 25, Fig. 26, Fig. 27, Fig. 28, Fig. 29, Fig. 30, Fig. 31, Fig. 32, Fig. 33, Fig. 34, Fig. 35, Fig. 36, Fig. 37, Fig. 38, and Fig. 39.

### H.1.3 PGD ATTACK

Here we report the evaluations using PGD attack, both as targeted (both targets:  $\vec{0}$  and  $-\vec{f}$ ) and non-targeted attacks optimized under the  $\ell_\infty$ -norm bound and the  $\ell_2$ -norm bound over multiple attack iterations. For  $\ell_\infty$ -norm bound, perturbation budget  $\epsilon = \frac{8}{255}$ , and step size  $\alpha=0.01$ , while for  $\ell_2$ -norm bound, perturbation budget  $\epsilon = \frac{64}{255}$  and step size  $\alpha=0.1$ . Attack evaluations include Fig. 40, Fig. 41, Fig. 42, Fig. 43, Fig. 44, Fig. 45, Fig. 46, Fig. 47, Fig. 48, Fig. 49, Fig. 50, Fig. 51, Fig. 52, Fig. 53, and Fig. 54.

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

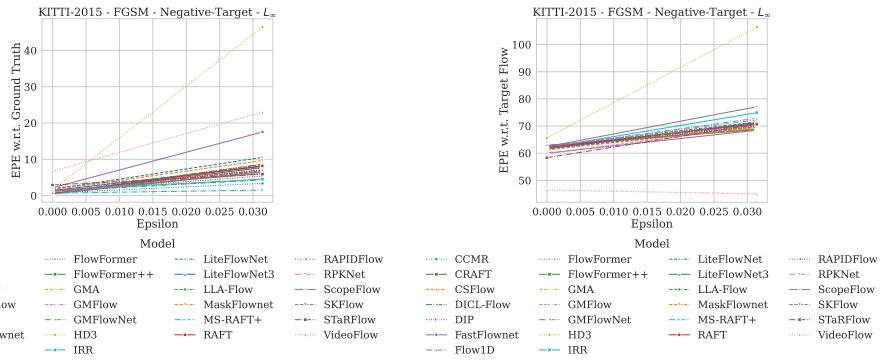


Figure 13: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

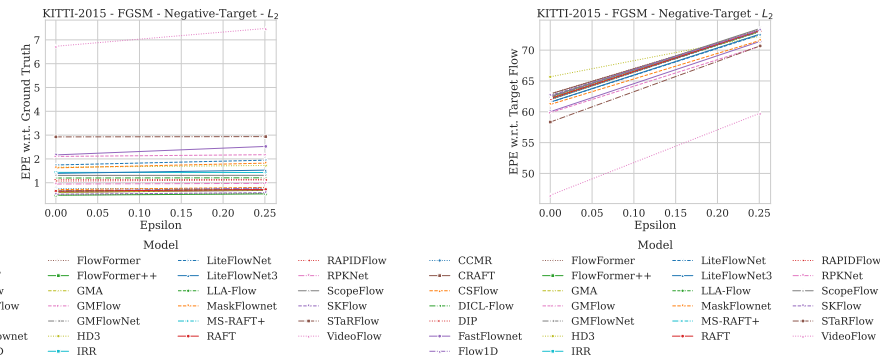


Figure 14: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

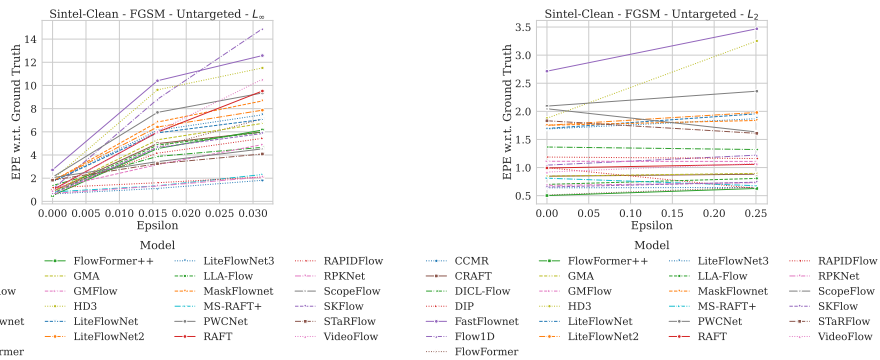


Figure 15: Evaluations for non-targeted FGSM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

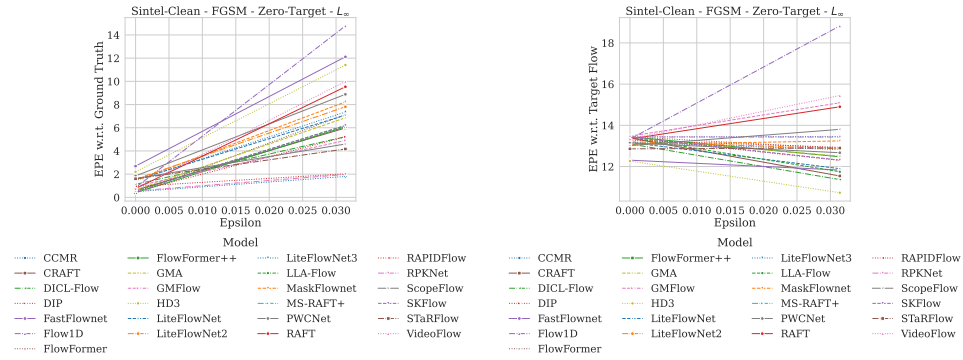


Figure 16: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

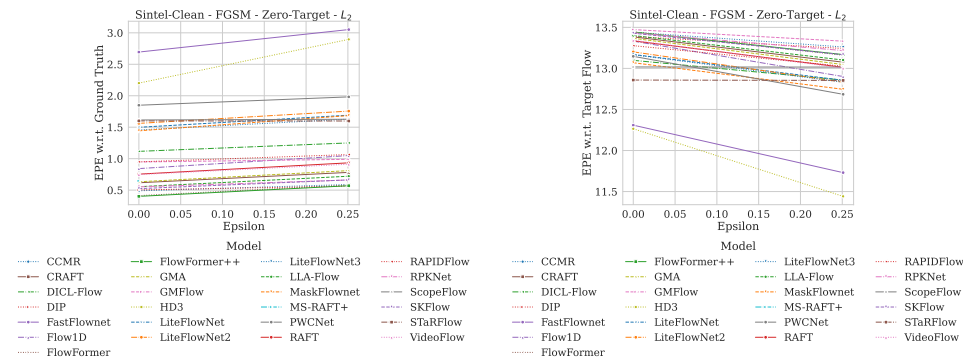


Figure 17: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

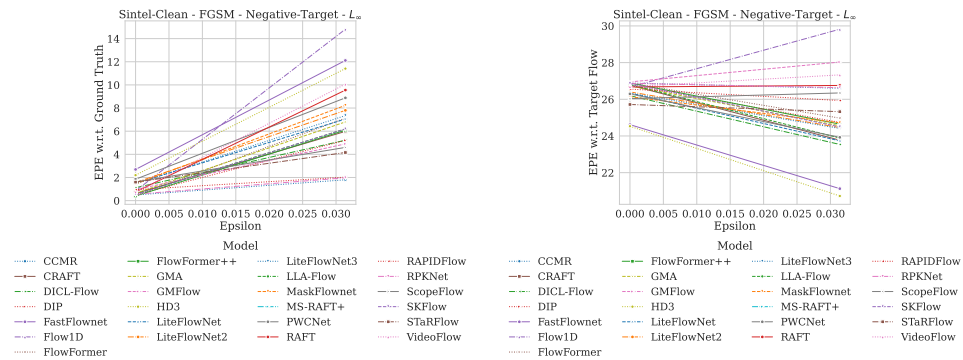


Figure 18: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

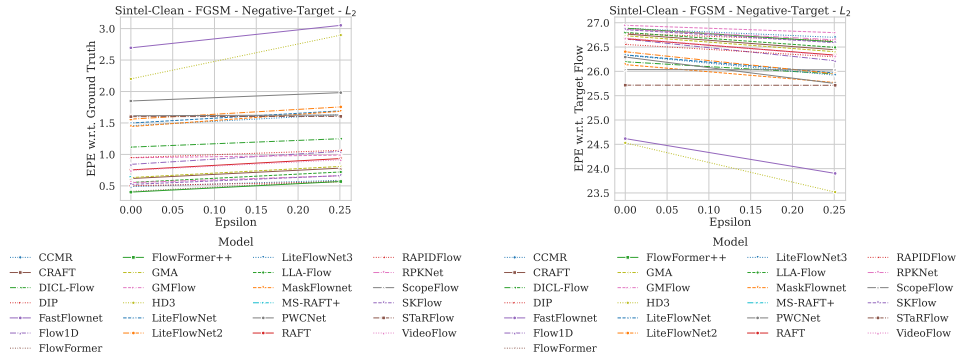


Figure 19: Evaluations for targeted FGSM attack with target  $-\hat{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

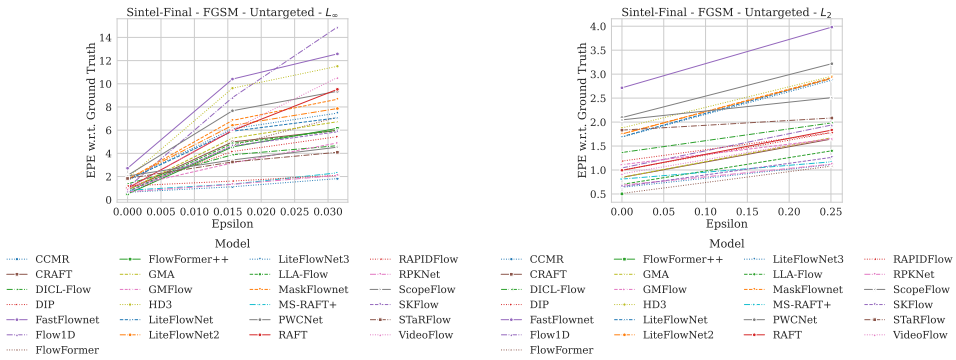


Figure 20: Evaluations for non-targeted FGSM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

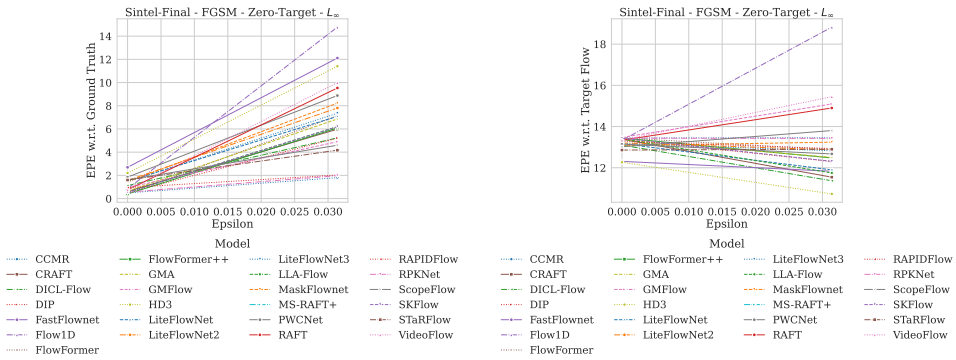


Figure 21: Evaluations for targeted FGSM attack with target  $\hat{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

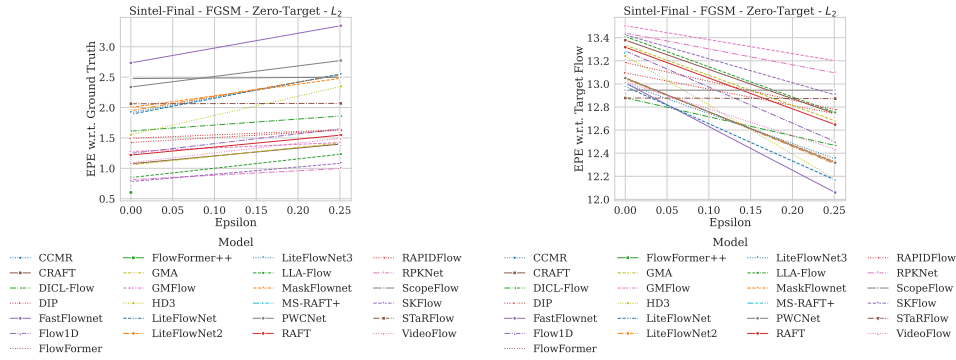


Figure 22: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

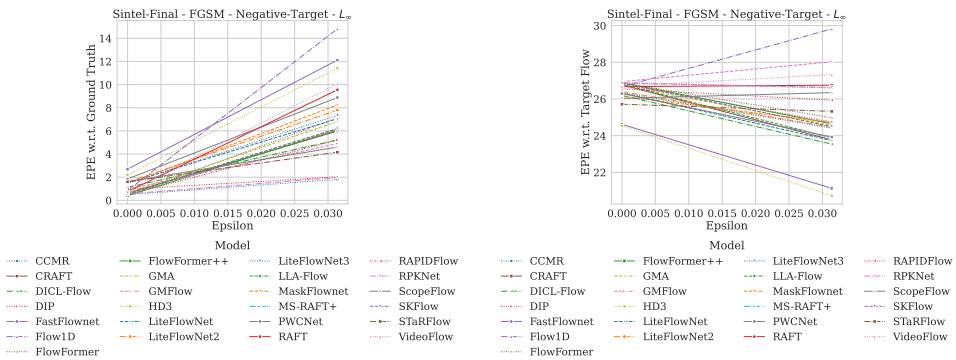


Figure 23: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

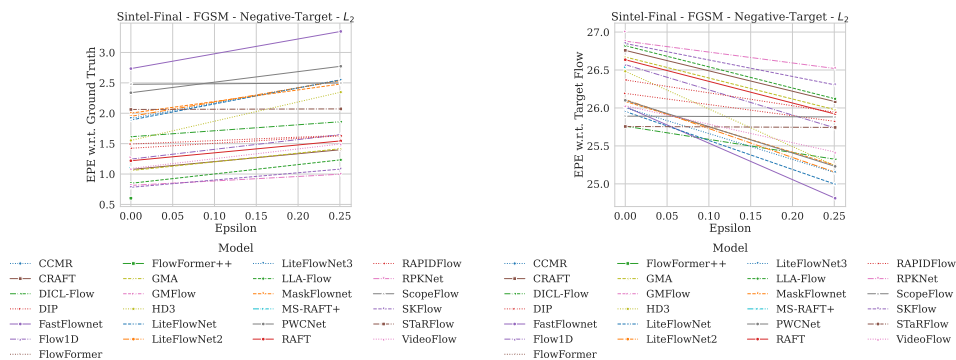


Figure 24: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

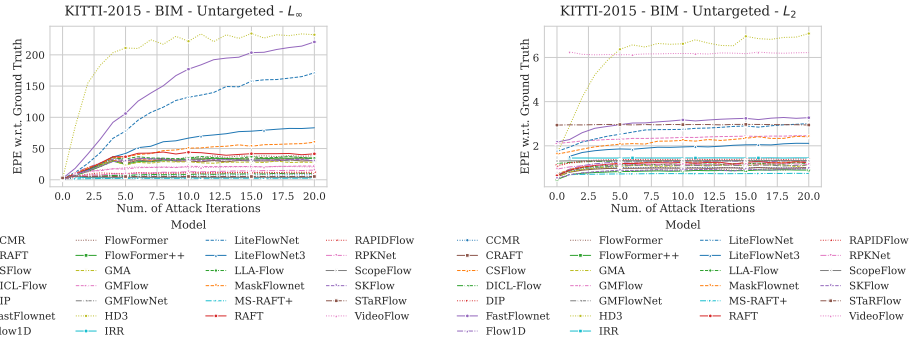


Figure 25: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

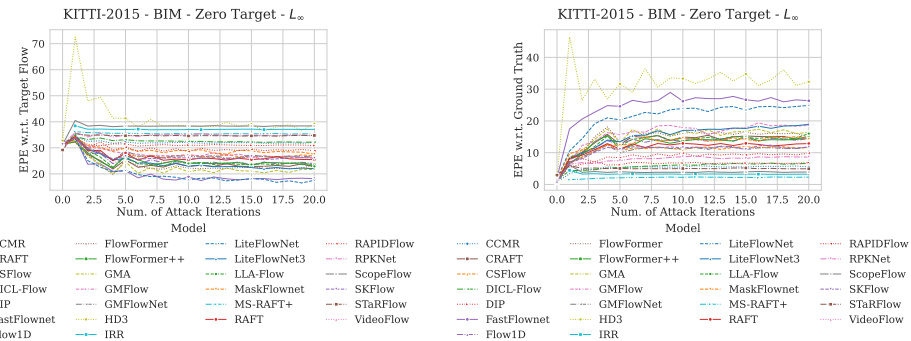


Figure 26: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

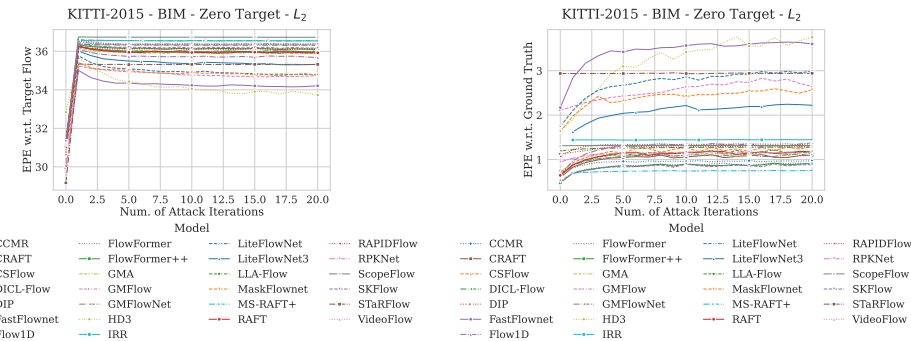


Figure 27: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.



2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

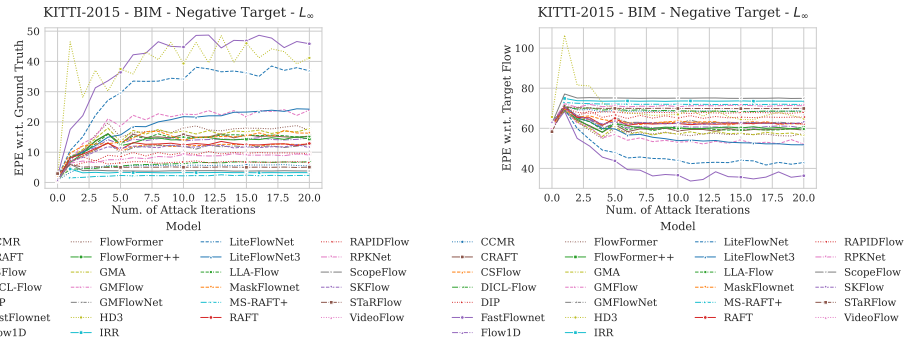


Figure 28: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

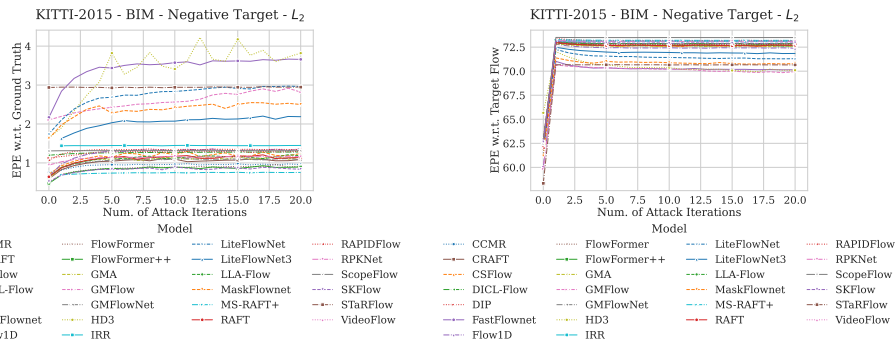


Figure 29: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

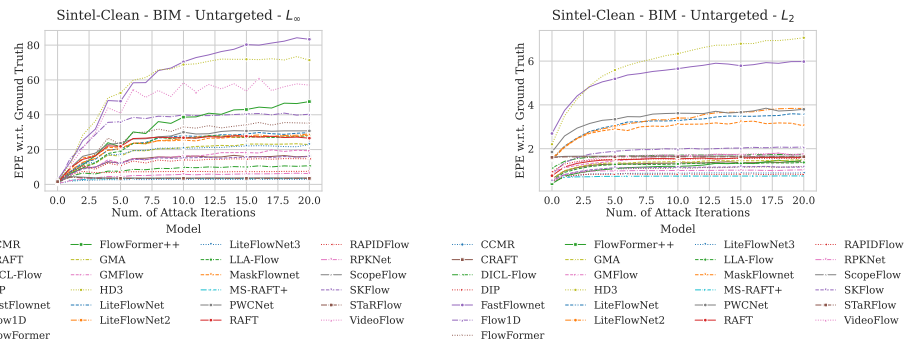


Figure 30: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

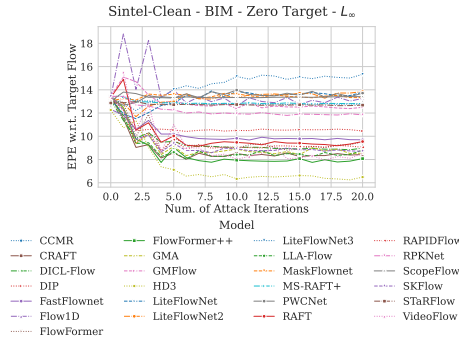


Figure 31: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

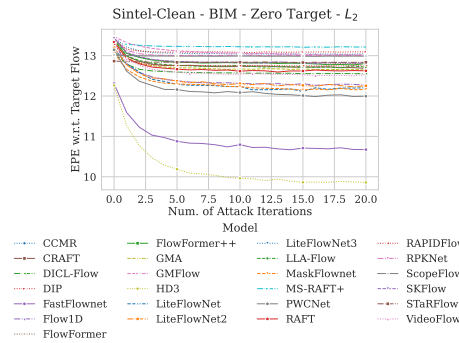


Figure 32: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

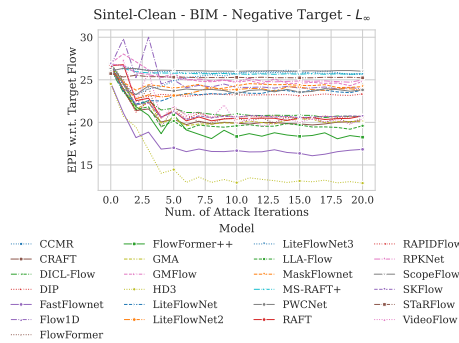


Figure 33: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

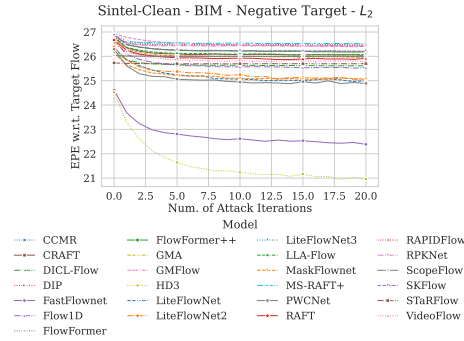


Figure 34: Evaluations for targeted BIM attack with target  $\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

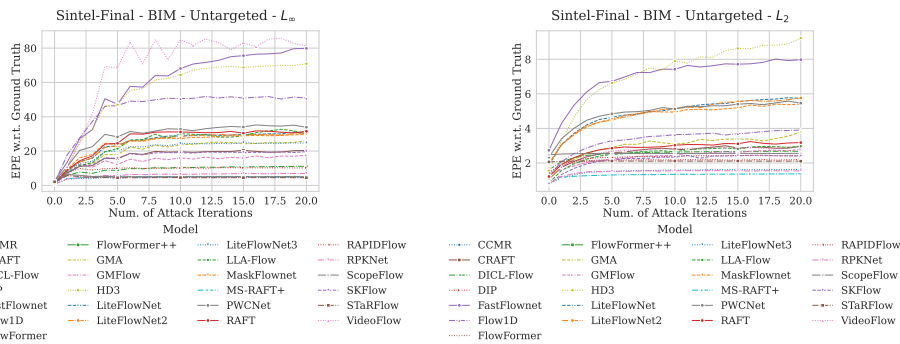


Figure 35: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

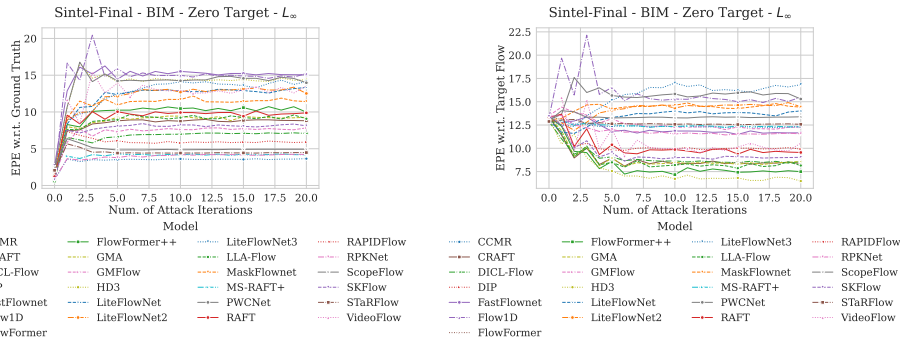


Figure 36: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

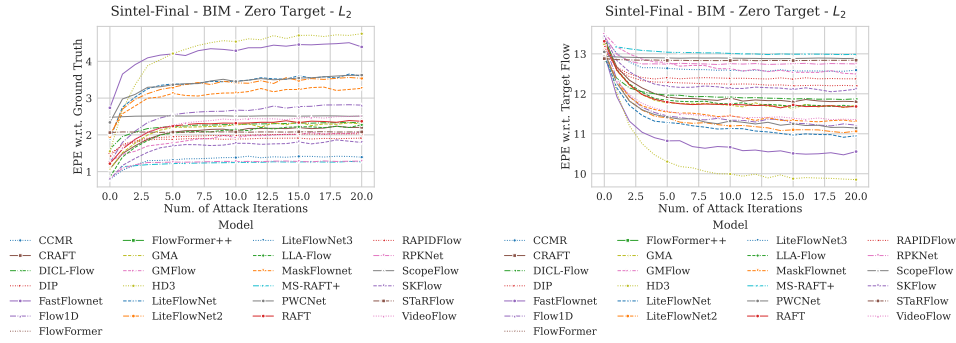


Figure 37: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

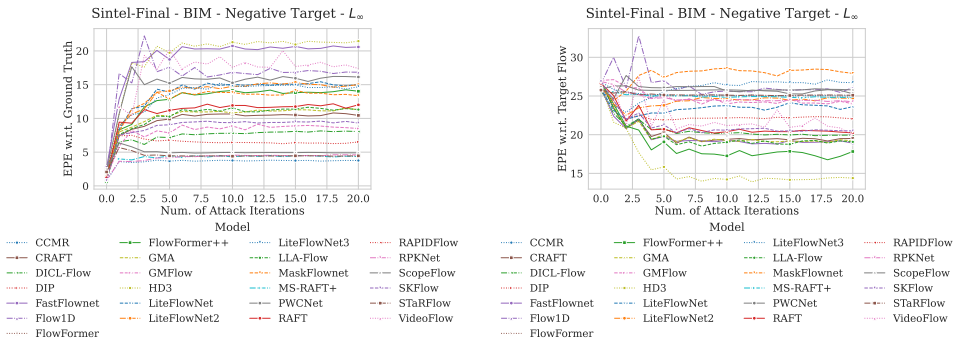


Figure 38: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

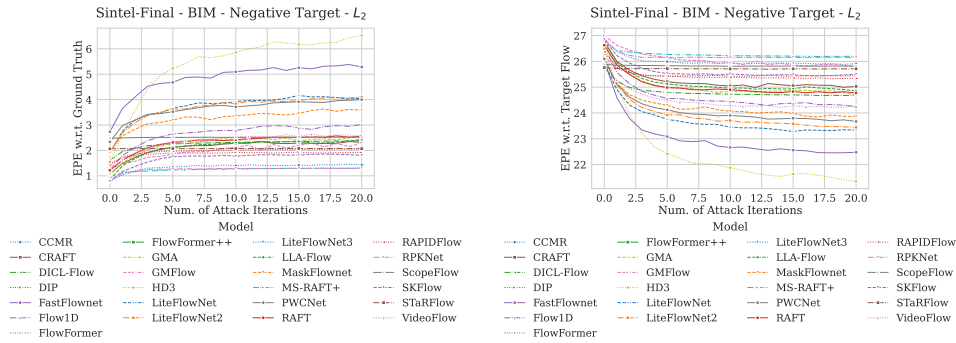


Figure 39: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.



2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389

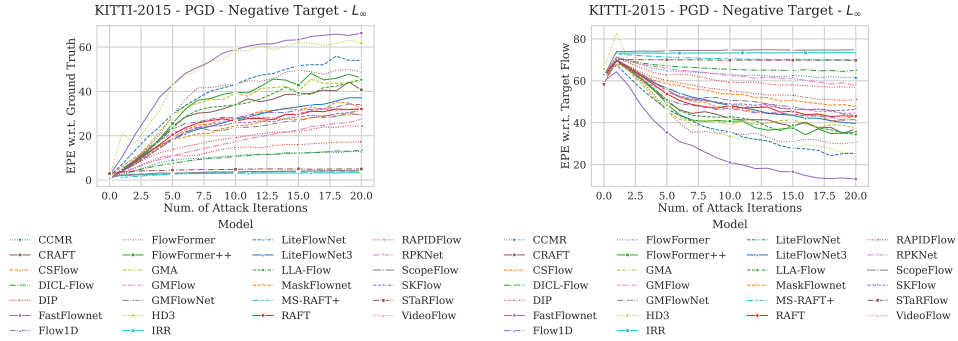


Figure 43: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

2392  
2393  
2394  
2395  
2396

2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407

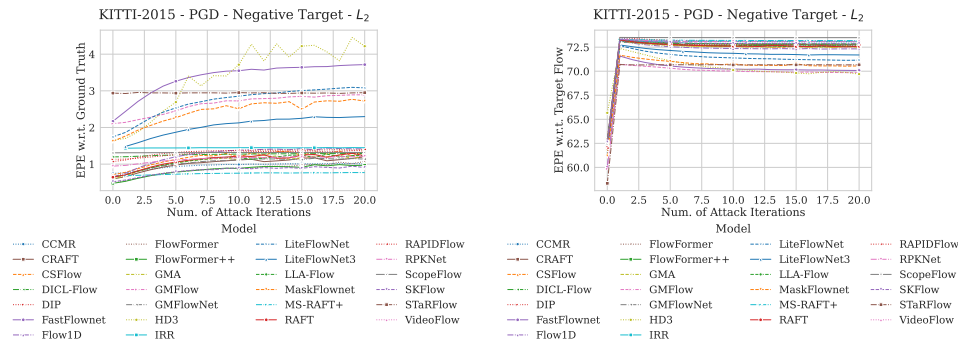


Figure 44: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

2411  
2412  
2413  
2414

2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425

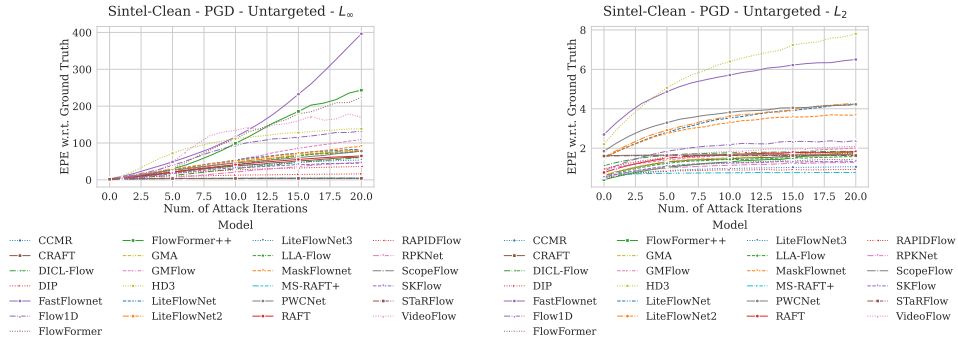


Figure 45: Evaluations for non-targeted PGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

2426  
2427  
2428  
2429

2430  
 2431  
 2432  
 2433  
 2434  
 2435  
 2436  
 2437  
 2438  
 2439  
 2440  
 2441  
 2442  
 2443  
 2444  
 2445  
 2446  
 2447  
 2448  
 2449  
 2450  
 2451  
 2452  
 2453  
 2454  
 2455  
 2456  
 2457  
 2458  
 2459  
 2460  
 2461  
 2462  
 2463  
 2464  
 2465  
 2466  
 2467  
 2468  
 2469  
 2470  
 2471  
 2472  
 2473  
 2474  
 2475  
 2476  
 2477  
 2478  
 2479  
 2480  
 2481  
 2482  
 2483

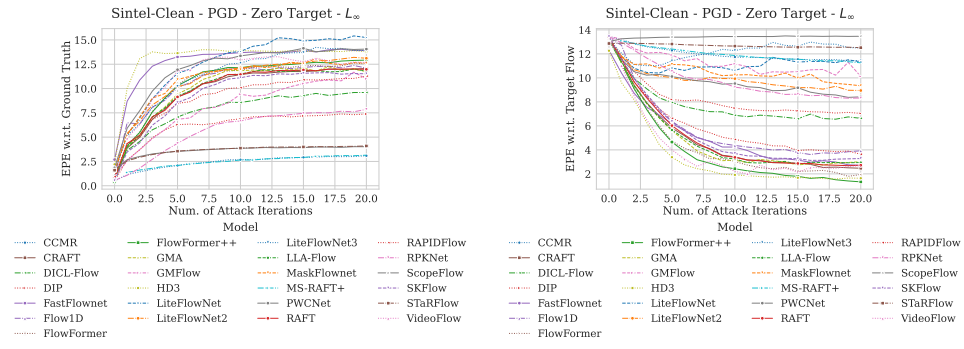


Figure 46: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

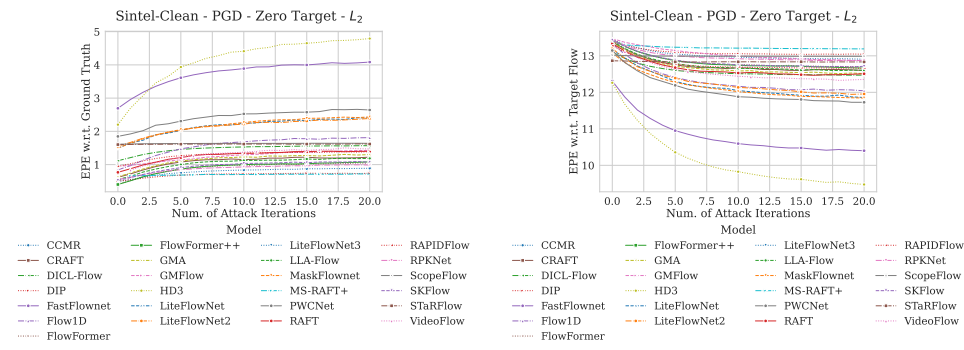


Figure 47: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

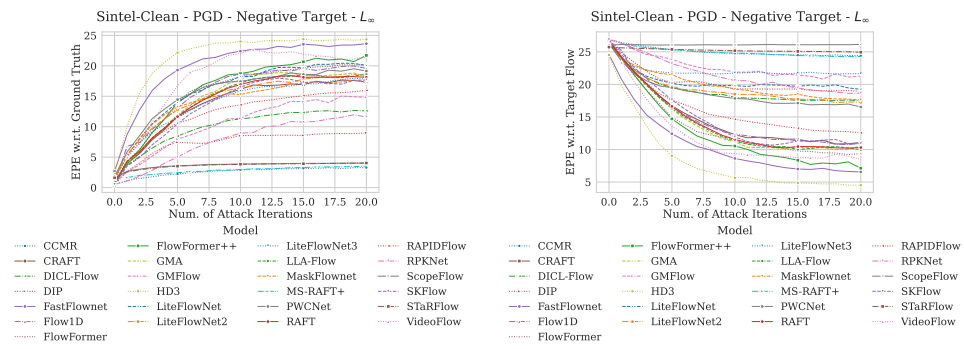


Figure 48: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.



2484  
2485  
2486  
2487  
2488  
2489  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537

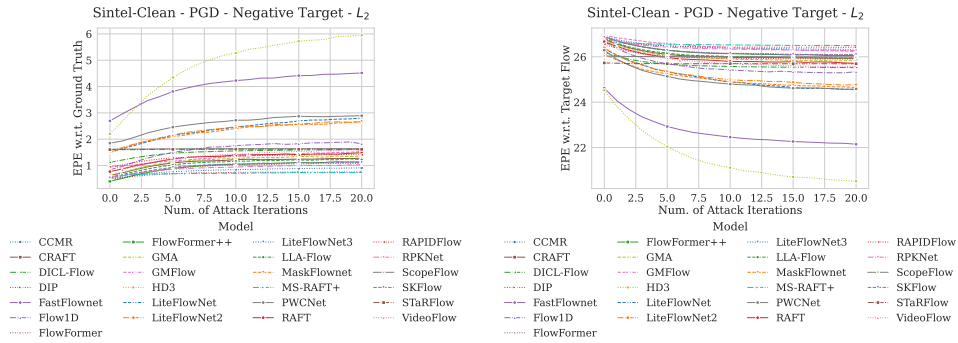


Figure 49: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

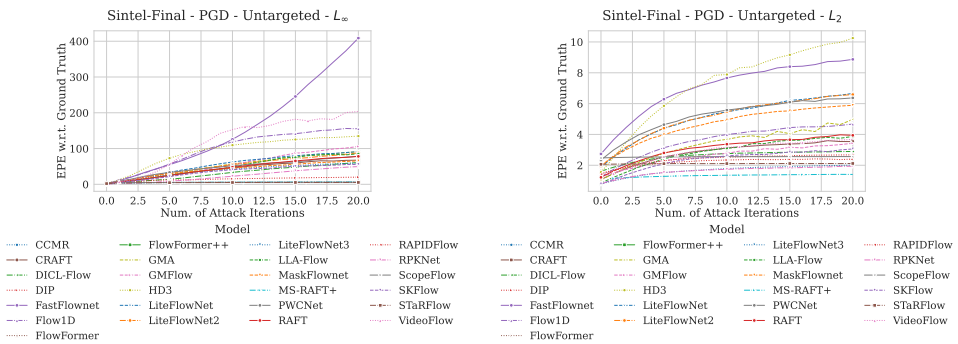


Figure 50: Evaluations for non-targeted PGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

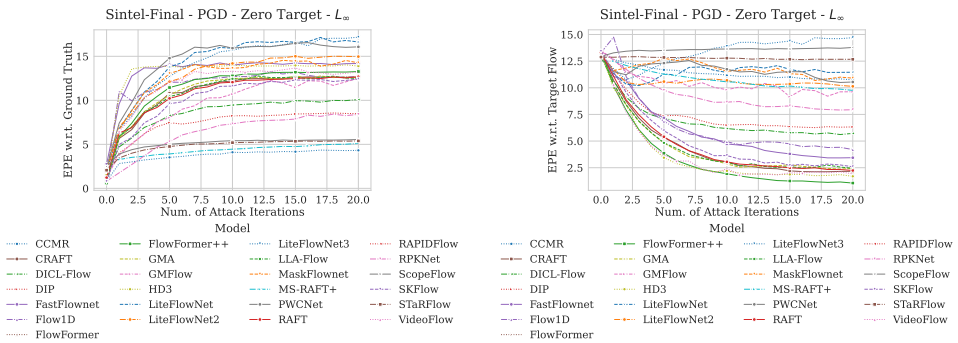


Figure 51: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

2538  
 2539  
 2540  
 2541  
 2542  
 2543  
 2544  
 2545  
 2546  
 2547  
 2548  
 2549  
 2550  
 2551  
 2552  
 2553  
 2554  
 2555  
 2556  
 2557  
 2558  
 2559  
 2560  
 2561  
 2562  
 2563  
 2564  
 2565  
 2566  
 2567  
 2568  
 2569  
 2570  
 2571  
 2572  
 2573  
 2574  
 2575  
 2576  
 2577  
 2578  
 2579  
 2580  
 2581  
 2582  
 2583  
 2584  
 2585  
 2586  
 2587  
 2588  
 2589  
 2590  
 2591

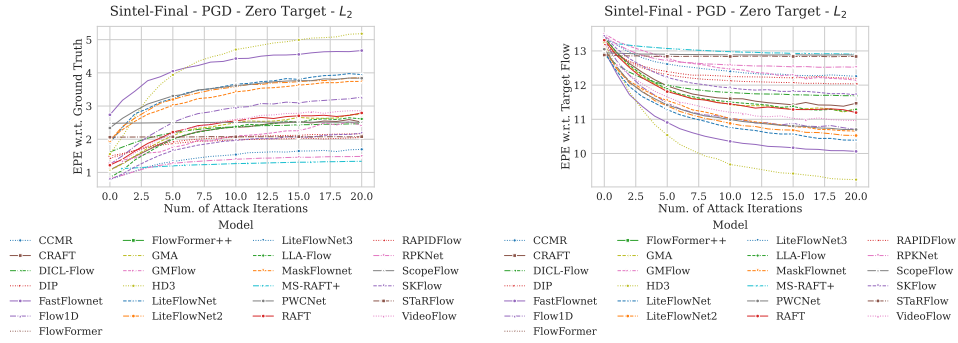


Figure 52: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

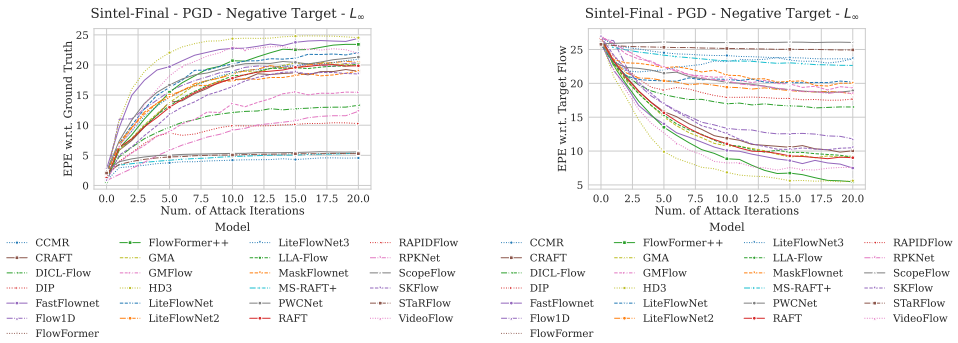


Figure 53: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

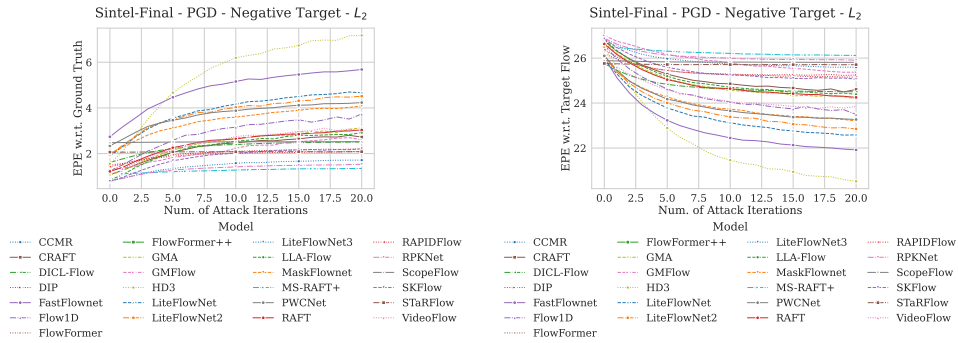


Figure 54: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.



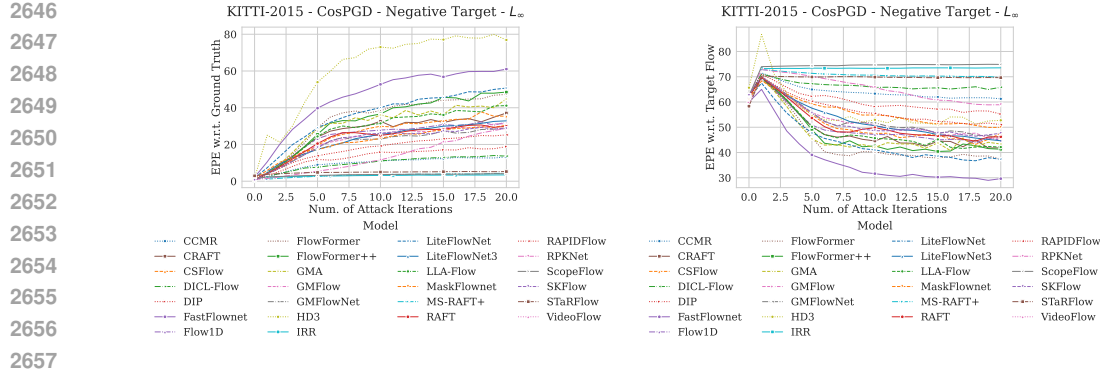


Figure 58: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

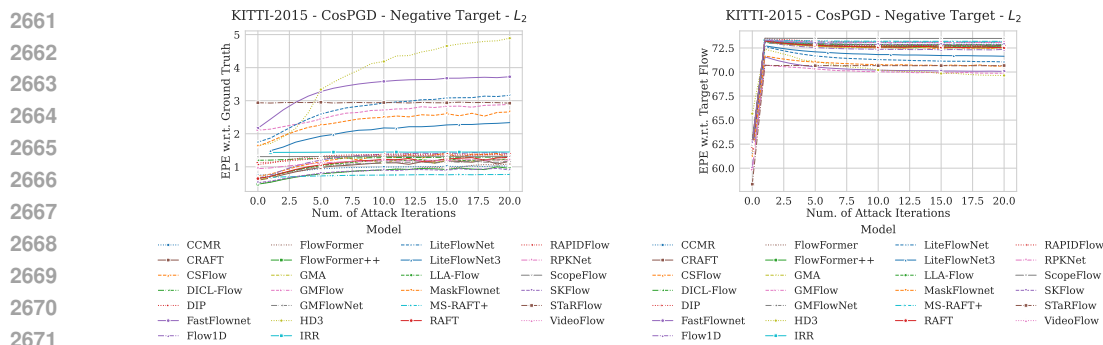


Figure 59: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

### H.1.5 PCFA ATTACK

Here we report the evaluations using PCFA attack, as targeted (both targets:  $\vec{0}$  and  $-\vec{f}$ ) optimized under the  $\ell_2$ -norm bound over multiple attack iterations. Here the perturbation budget  $\epsilon = 0.05$  and step size  $\alpha = 1e - 7$ . Attack evaluations include Fig. 70 and Fig. 71.

### H.1.6 ADVERSARIAL WEATHER ATTACK

Here we report the evaluations using different Adversarial Weather, both as targeted (both targets:  $\vec{0}$  and  $-\vec{f}$ ) and non-targeted attacks. Attack evaluations include Fig. 72, Fig. 73, Fig. 74 and Fig. 75.

## H.2 COMMON CORRUPTIONS OVERVIEW

Following we provide an overview of the performance over all corruptions. This is reported in Fig. 76.

### H.3 2D COMMON CORRUPTIONS

Here we report evaluations using different 2D common corruptions over all considered datasets. OOD Robustness evaluations with 2D Common Corruptions include Fig. 77, Fig. 78 and Fig. 79.

### H.4 3D COMMON CORRUPTIONS

Here we report evaluations using different considered 3D common corruptions over all considered datasets. OOD Robustness evaluations with 3D Common Corruptions include Fig. 80, Fig. 81 and Fig. 82.

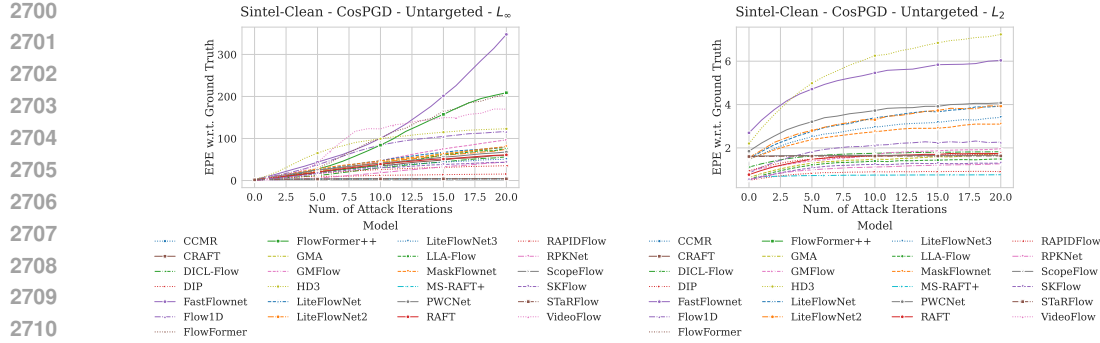


Figure 60: Evaluations for non-targeted CosPGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

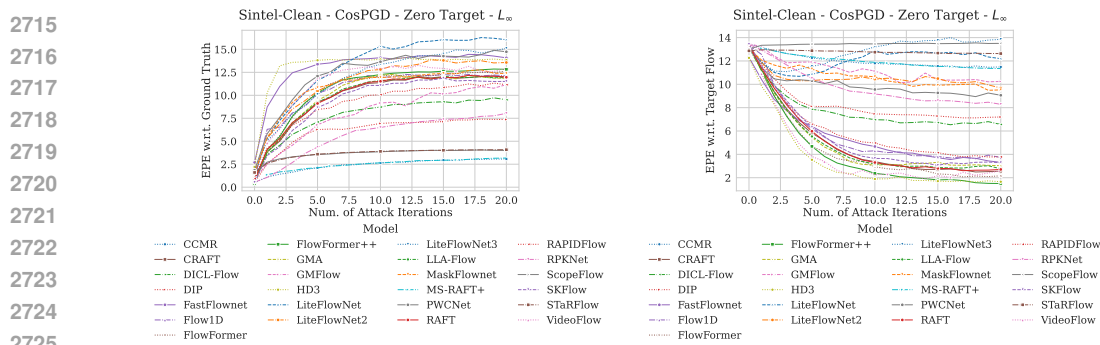


Figure 61: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

## I INITIAL PROTOTYPE OF THE FUTURE WEBSITE

NEW

In Figure 83 we share a screenshot from our prototype website currently under work, that would help better understand the metrics. In this screenshot, the methods are ranked based on their EPE w.r.t. the ground truth flow under non-targeted CosPGD attack at 20 attack iterations under the  $\ell_\infty$ -norm bound (lower means the method is more robust) evaluated using the KITTI2015 dataset. We are currently designing it to make the numbers and column headings better visible to the users, and the users can dynamically rank these based on any of the columns. We will host the website after acceptance.

## J LIMITATIONS

Benchmarking optical flow estimation methods is a compute and labor-intensive endeavor. Thus, best utilizing available resources we use FLOWBENCH to benchmark a limited number of settings. The benchmarking tool itself offers significantly more combinations that can be benchmarked. Nonetheless, the benchmark provided is comprehensive and instills interest to further utilize FLOWBENCH.

## K REPRODUCIBILITY OF EVALUATIONS

NEW

There always exists stochasticity when evaluating adversarial attacks, due to the randomness these attacks exploit, and also common corruptions due to variations in seeds and calculation approximations made by various python libraries. Therefore, for transparency and reproducibility, we evaluate different runs on the same seed and different runs of different seeds. We report these evaluations

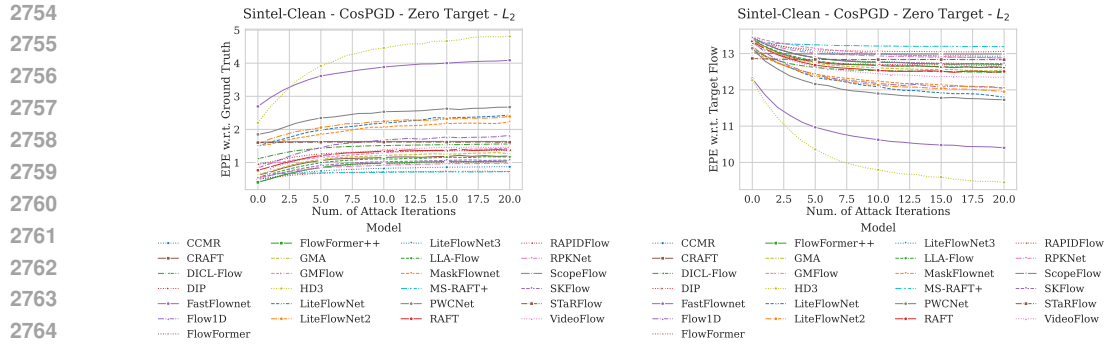


Figure 62: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

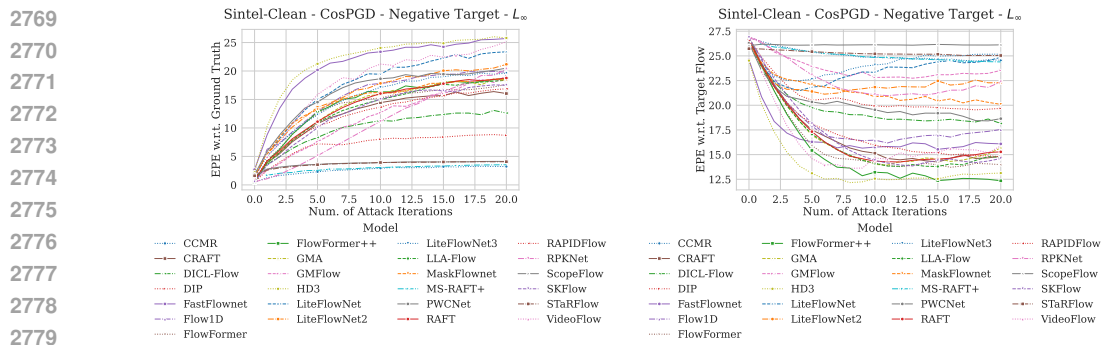


Figure 63: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

here, using Tab. 2 for adversarial attacks and Tab. 3 for common corruptions, and observe that the variance is extremely low and the analysis performed in this work still stands.

To ensure the reproducibility of our adversarial attack evaluations we repeat experiments in two ways: first, three different runs with the same seed, and second, one run each for three different seeds. We observe very minute variations in results in both cases which can be attributed to calculation approximations made by different libraries such as pytorch (Paszke et al., 2019). Due to the compute-hungry nature of these evaluations, we limit them to using one method: RAFT on the KITTI2015 dataset, and the attack used is CosPGD. We evaluate multiple settings: different  $\ell_p$ -norm bounds, different attack optimization methods (optimizing w.r.t. ground truth flow and optimizing w.r.t. initial flow prediction.), and for targeted attacks, two different targets. The attack settings are consistent with the paper.



2808  
2809  
2810  
2811  
2812  
2813  
2814  
2815  
2816  
2817  
2818  
2819  
2820  
2821  
2822  
2823  
2824  
2825  
2826  
2827  
2828  
2829  
2830  
2831  
2832  
2833  
2834  
2835  
2836  
2837  
2838  
2839  
2840  
2841  
2842  
2843  
2844  
2845  
2846  
2847  
2848  
2849  
2850  
2851  
2852  
2853  
2854  
2855  
2856  
2857  
2858  
2859  
2860  
2861

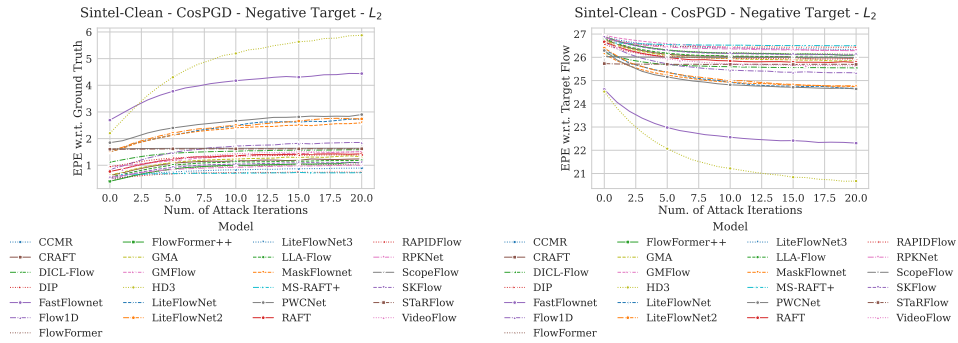


Figure 64: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

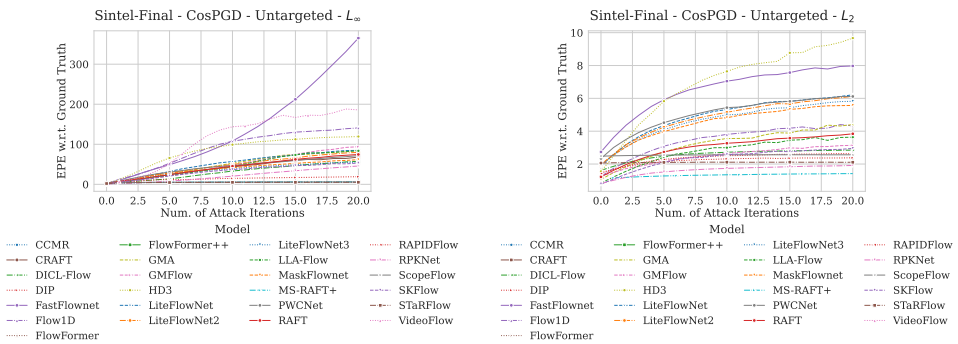


Figure 65: Evaluations for non-targeted CosPGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

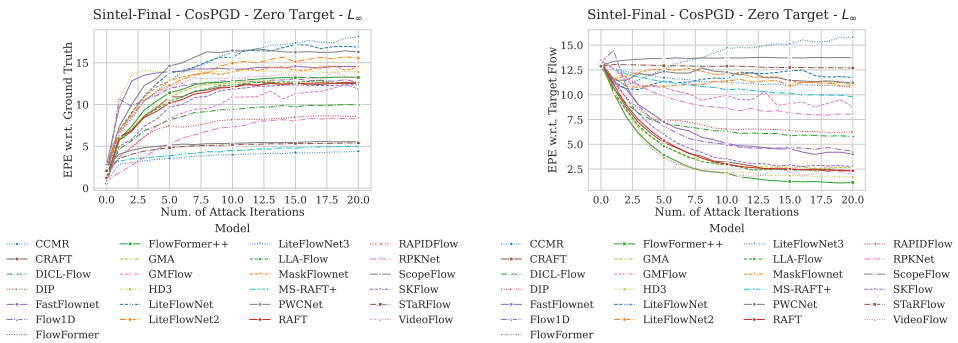
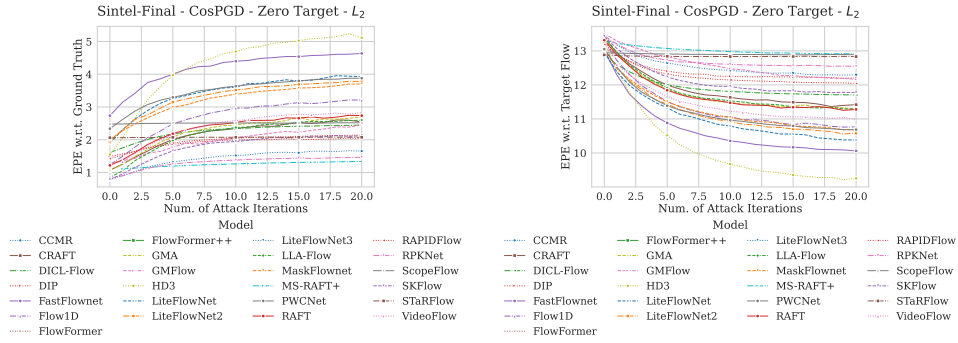


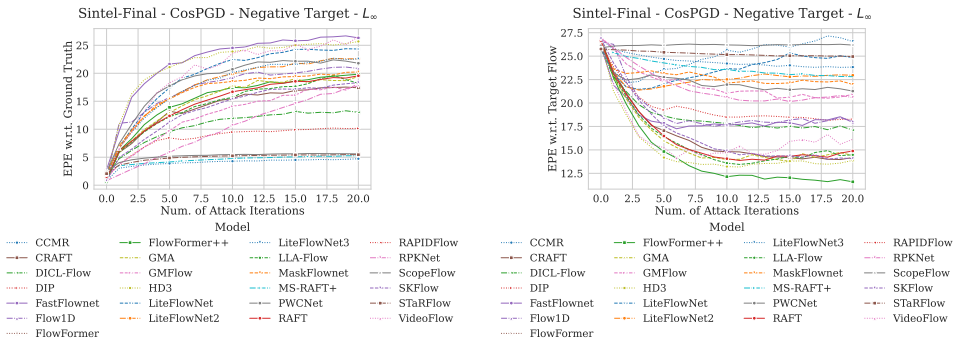
Figure 66: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875



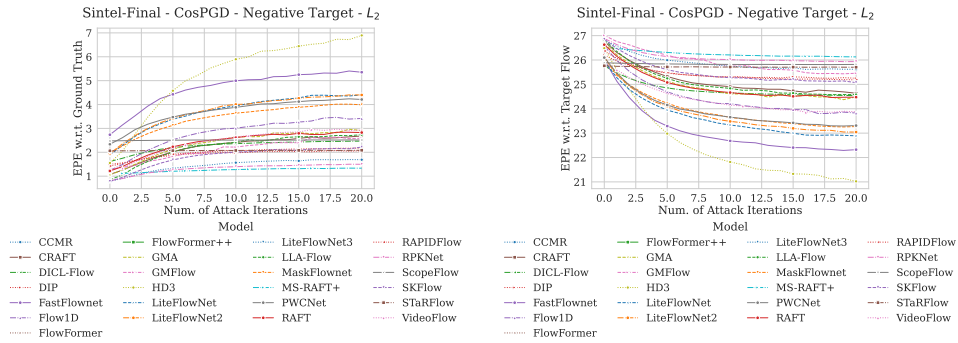
2876 Figure 67: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the  
2877 MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.  
2878  
2879  
2880

2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893



2894 Figure 68: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the  
2895 MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.  
2896  
2897  
2898  
2899

2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911



2912 Figure 69: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the  
2913 MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.  
2914  
2915

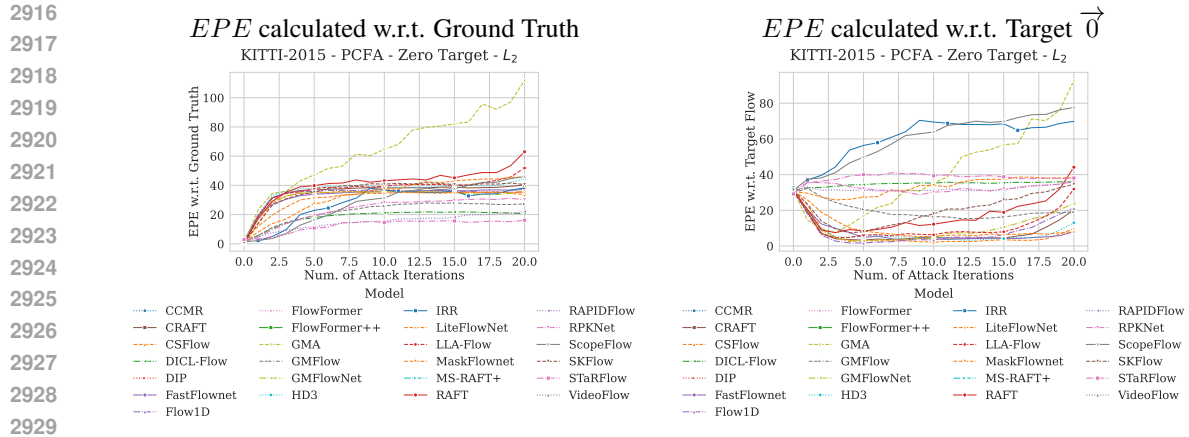


Figure 70: Evaluating all optical flow estimation methods against PCFA attack with target  $\vec{0}$  over multiple attack iterations.

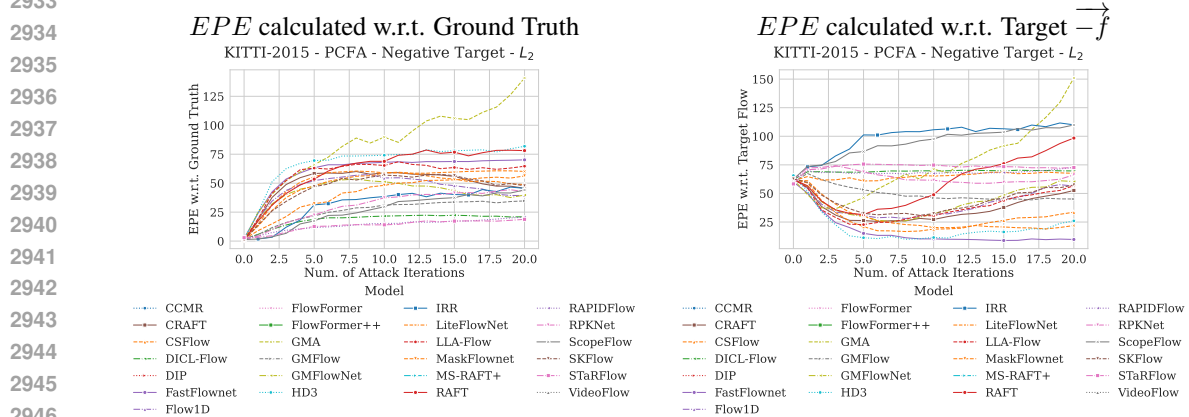


Figure 71: Evaluating all optical flow estimation methods against PCFA attack with target  $-\vec{f}$  over multiple attack iterations.

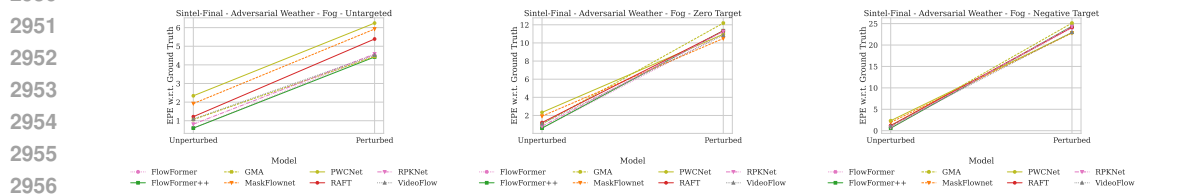


Figure 72: Evaluations for Adversarial Weather attack with Fog optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

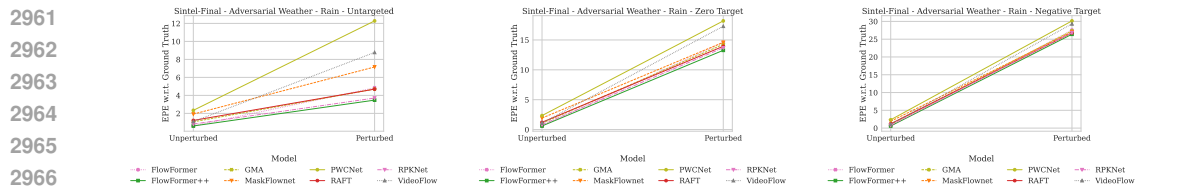


Figure 73: Evaluations for Adversarial Weather attack with Rain optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

2970  
2971  
2972  
2973  
2974  
2975  
2976  
2977  
2978  
2979  
2980  
2981  
2982  
2983  
2984  
2985  
2986  
2987  
2988  
2989  
2990  
2991  
2992  
2993  
2994  
2995  
2996  
2997  
2998  
2999  
3000  
3001  
3002  
3003  
3004  
3005  
3006  
3007  
3008  
3009  
3010  
3011  
3012  
3013  
3014  
3015  
3016  
3017  
3018  
3019  
3020  
3021  
3022  
3023

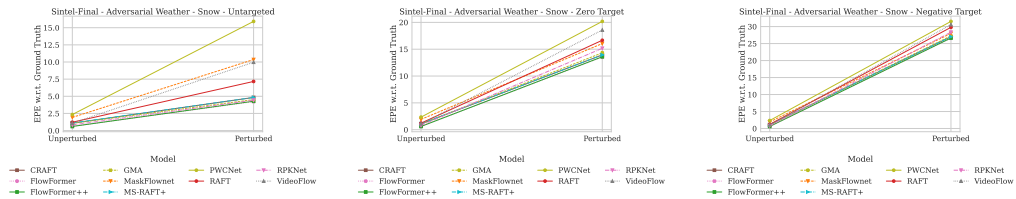


Figure 74: Evaluations for Adversarial Weather attack with Snow optimized as an non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

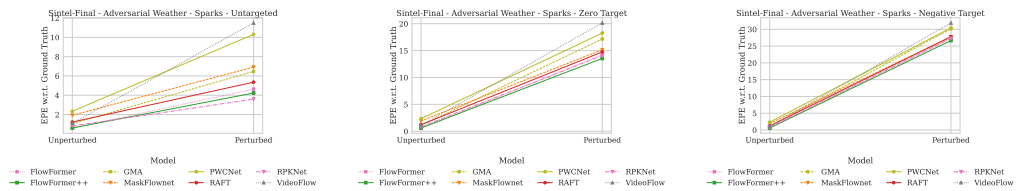


Figure 75: Evaluations for Adversarial Weather attack with Sparks optimized as an non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

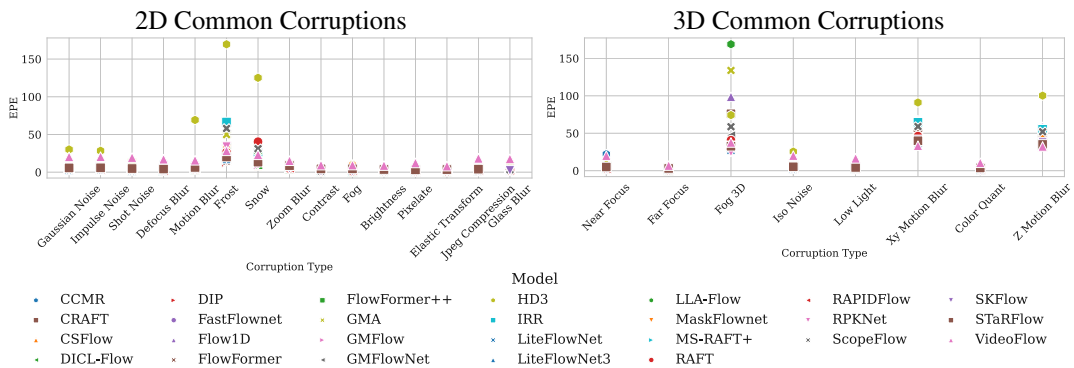


Figure 76: Performance of various optical flow estimation methods after corruptions on the KITTI2015 dataset.



3078  
3079  
3080  
3081  
3082  
3083  
3084  
3085  
3086  
3087  
3088  
3089  
3090  
3091  
3092  
3093  
3094  
3095  
3096  
3097  
3098  
3099  
3100  
3101  
3102  
3103  
3104  
3105  
3106  
3107  
3108  
3109  
3110  
3111  
3112  
3113  
3114  
3115  
3116  
3117  
3118  
3119  
3120  
3121  
3122  
3123  
3124  
3125  
3126  
3127  
3128  
3129  
3130  
3131

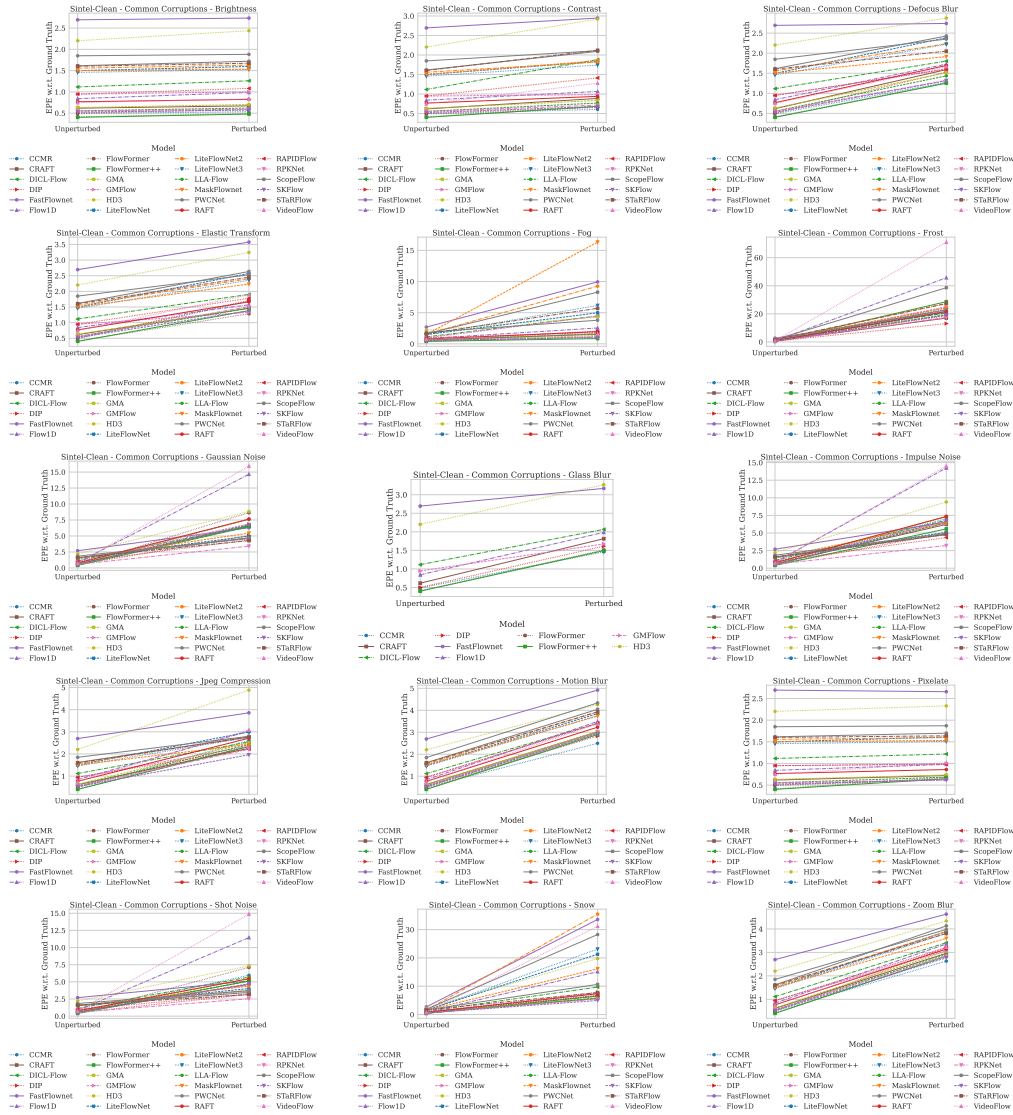


Figure 78: Evaluating optical flow estimation methods against all 2D Common Corruptions on the MPI Sintel (clean) dataset.



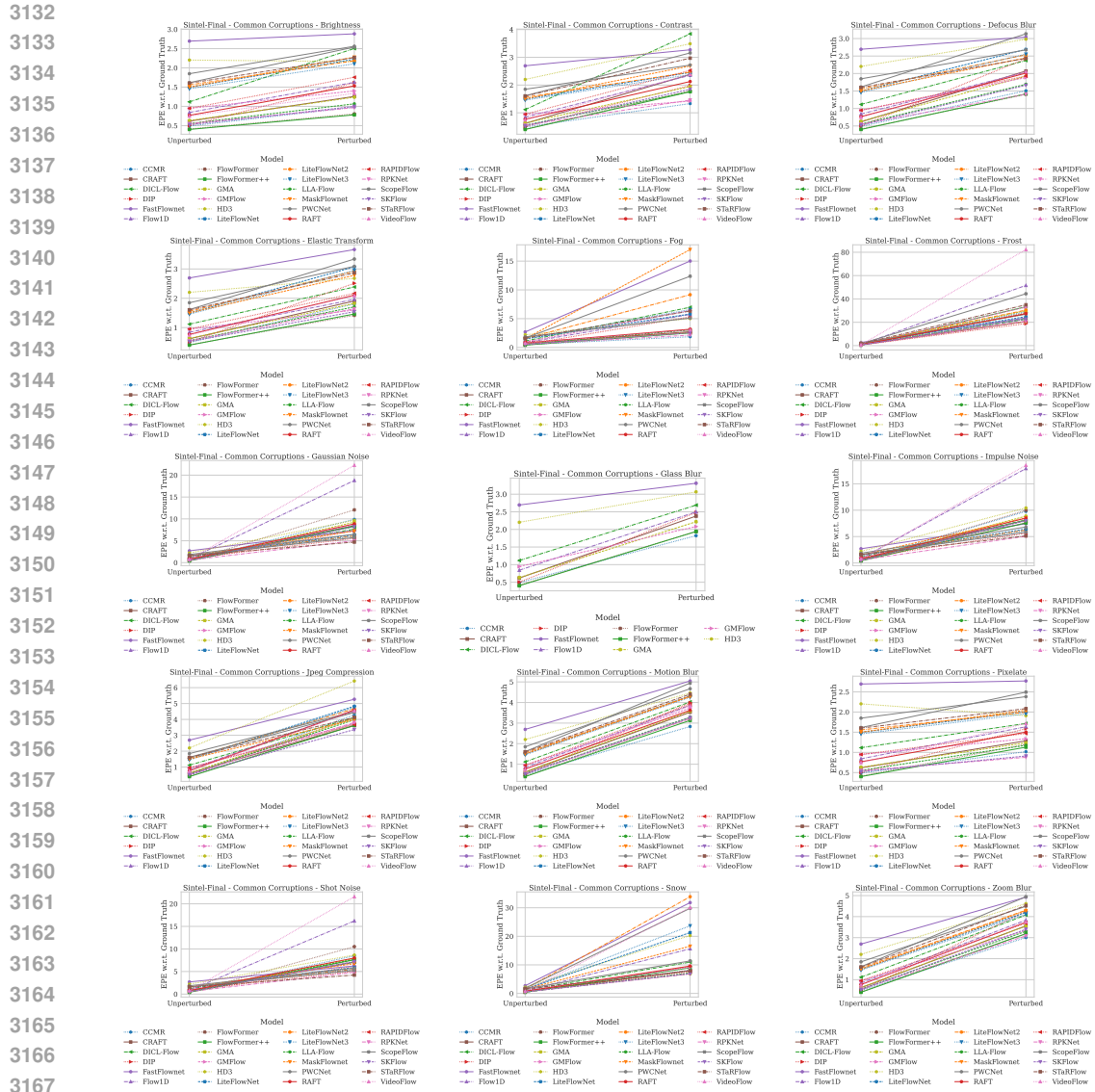


Figure 79: Evaluating optical flow estimation methods against all 2D Common Corruptions on the MPI Sintel (final) dataset.

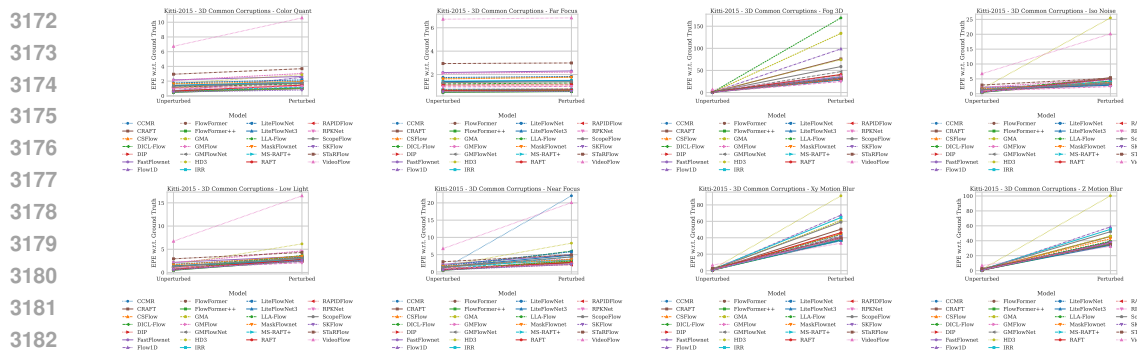


Figure 80: Evaluating optical flow estimation methods against the considered 3D Common Corruptions on the KITTI2015 dataset.

3186  
3187  
3188  
3189  
3190  
3191  
3192  
3193  
3194  
3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239

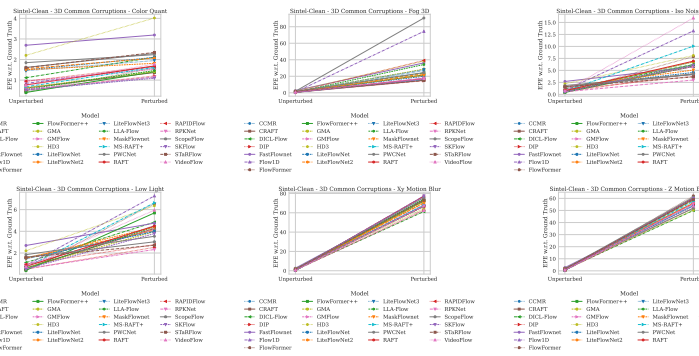


Figure 81: Evaluating optical flow estimation methods against the considered 3D Common Corruptions on the MPI Sintel (clean) dataset.

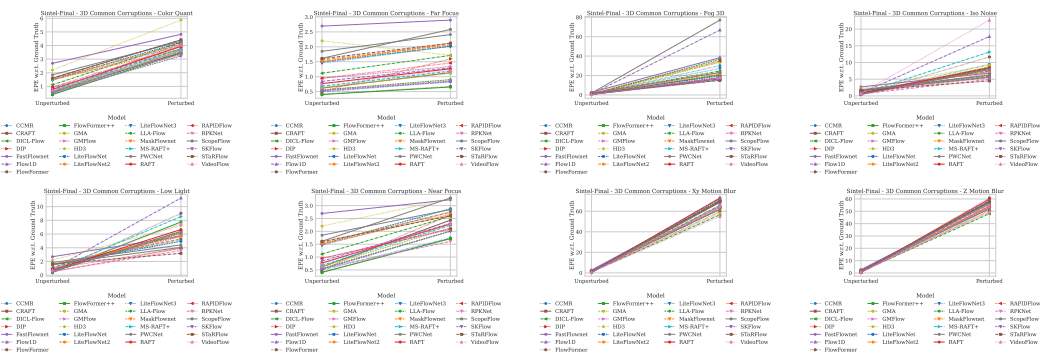


Figure 82: Evaluating optical flow estimation methods against the considered 3D Common Corruptions on the MPI Sintel (final) dataset.

### FlowBench

Optical flow estimation is a crucial computer vision task often applied to safety critical real-world scenarios like autonomous driving and medical imaging. While optical flow estimation accuracy has greatly benefited from the emergence of Deep learning, learning-based methods are also known for their lack of generalization and reliability. However, reliability is paramount when optical flow methods are employed in the real world, where safety is essential. Furthermore, a deeper understanding of the robustness and reliability of learning-based optical flow estimation methods is still lacking, hindering the research community from building methods safe for real-world deployment. Thus we propose FLOWBENCH, a robustness benchmark and evaluation tool for learning-based optical flow methods. FLOWBENCH facilitates streamlined research into the reliability of optical flow methods by benchmarking their robustness to adversarial attacks and out-of-distribution samples. With FLOWBENCH, we benchmark 91 methods across 3 different datasets under 7 diverse adversarial attacks and 23 established common corruptions, making it the most comprehensive robustness analysis of optical flow methods to date. Across this wide range of methods, we consistently find that methods with state-of-the-art performance on established standard benchmarks lack reliability and generalization ability. Moreover, we find interesting correlations between performance, reliability, and generalization ability of optical flow estimation methods, under various lenses such as design choices used, number of parameters, etc. After acceptance, FLOWBENCH will be open-source and publicly available, including the weights of all tested models.

Show rows:   
Dataset:

#### Leaderboard: Optical Flow Estimation

Rank	Architecture	CosPGD-EPE No...	PCFA-EPE target...	PGD-EPE Non-ta...	Checkpoint	Time_Proposed
1	IRR_PWC	2.4330115861	69.8680467474	3.8397940102	KITTI	Mar 2020
2	ScopeFlow	2.6299911237	77.6024067444	4.5666427672	KITTI	Nov 2023
3	MS-RAFT+	2.695769617	40.9803659793	5.0118829945	KITTI	Jul 2020
4	StarFlow	4.1756131017	38.06560606	5.4629693043	KITTI	Mar 2024
5	DICL	10.5310789597	35.934047974	21.0924557686	KITTI	Jul 2021

Figure 83: A share a screenshot from our prototype website currently under works, that would help better understand the metrics. In this screenshot, the methods are ranked based on their EPE w.r.t. the ground truth flow under non-targeted CosPGD attack at 20 attack iterations under the  $\ell_\infty$ -norm bound (lower means the method is more robust) evaluated using the KITTI2015 dataset. We are currently designing it to make the numbers and column headings better visible to the users, and the users can dynamically rank these based on any of the columns.

3240  
3241  
3242  
3243  
3244  
3245  
3246  
3247  
3248  
3249  
3250  
3251  
3252  
3253  
3254  
3255  
3256  
3257  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283  
3284  
3285  
3286  
3287  
3288  
3289  
3290  
3291  
3292  
3293

Table 2: To ensure reproducibility of our adversarial attack evaluations we repeat experiments in two ways: first, three different runs with the same seed, and second, one run each for three different seeds. We observe very minute variations in results in both cases which can be attributed to calculation approximations made by different libraries such as pytorch (Paszke et al., 2019). Due to the compute-hungry nature of these evaluations, we limit them to using one method: RAFT on the KITTI2015 dataset, and the attack used is CosPGD. We evaluate multiple settings: different  $\ell_p$ -norm bounds, different attack optimization methods (optimizing w.r.t. ground truth flow and optimizing w.r.t. initial flow prediction.), and for targeted attacks, two different targets. The attack settings are consistent with the paper. Target ‘None’ means the attack was Non-targeted.

$\ell_p$ -norm bound	Target	Attack Optimized w.r.t.	EPE mean $\pm$ std	px3 error mean $\pm$ std
Three different runs on the same seed				
$\ell_\infty$ -norm	None	Ground Truth Flow	119.504 $\pm$ 2.95E+0	0.078 $\pm$ 6.76E-3
$\ell_\infty$ -norm	$\vec{-f}$	Ground Truth Flow	45.357 $\pm$ 4.26E-1	0.200 $\pm$ 7.84E-4
$\ell_\infty$ -norm	$\vec{0}$	Ground Truth Flow	10.674 $\pm$ 2.74E-1	0.647 $\pm$ 1.06E-2
$\ell_2$ -norm	None	Ground Truth Flow	0.644 $\pm$ 2.72E-6	0.968 $\pm$ 3.38E-6
$\ell_2$ -norm	$\vec{-f}$	Ground Truth Flow	73.454 $\pm$ 3.12E-5	0.129 $\pm$ 2.86E-7
$\ell_2$ -norm	$\vec{0}$	Ground Truth Flow	36.724 $\pm$ 2.11E-5	0.170 $\pm$ 4.41E-7
$\ell_2$ -norm	None	Initial Flow Pred	0.643 $\pm$ 7.74E-6	0.968 $\pm$ 1.49E-6
One run each using three different seeds				
$\ell_\infty$ -norm	None	Ground Truth Flow	119.692 $\pm$ 1.75E+0	0.077 $\pm$ 4.27E-3
$\ell_\infty$ -norm	$\vec{-f}$	Ground Truth Flow	45.149 $\pm$ 9.16E-1	0.202 $\pm$ 1.65E-3
$\ell_\infty$ -norm	$\vec{0}$	Ground Truth Flow	11.016 $\pm$ 4.91E-1	0.625 $\pm$ 9.90E-3
$\ell_2$ -norm	None	Ground Truth Flow	0.644 $\pm$ 7.46E-6	0.968 $\pm$ 6.05E-6
$\ell_2$ -norm	$\vec{-f}$	Ground Truth Flow	73.454 $\pm$ 1.15E-4	0.129 $\pm$ 1.85E-7
$\ell_2$ -norm	$\vec{0}$	Ground Truth Flow	36.724 $\pm$ 1.10E-4	0.170 $\pm$ 7.92E-7
$\ell_2$ -norm	None	Initial Flow Pred	0.643 $\pm$ 1.44E-4	0.968 $\pm$ 1.55E-5

3294  
 3295  
 3296  
 3297  
 3298  
 3299  
 3300  
 3301  
 3302  
 3303  
 3304  
 3305  
 3306  
 3307  
 3308  
 3309  
 3310  
 3311  
 3312  
 3313  
 3314  
 3315  
 3316  
 3317  
 3318  
 3319  
 3320  
 3321  
 3322  
 3323  
 3324  
 3325  
 3326  
 3327  
 3328  
 3329  
 3330  
 3331  
 3332  
 3333  
 3334  
 3335  
 3336  
 3337  
 3338  
 3339  
 3340  
 3341  
 3342  
 3343  
 3344  
 3345  
 3346  
 3347

Table 3: To ensure reproducibility of our Common Corruptions evaluations we repeat experiments in two ways: first, three different runs with the same seed, and second, one run each for three different seeds. We observe extremely minute variations in results which can be attributed to differences in seeds and calculation approximations made by the Python libraries. Due to the compute-hungry nature of these evaluations, we limit them to using one method: RAFT on the KITTI2015 dataset, and all the fifteen 2D Common Corruptions.

2D Common Corruption Name	EPE mean $\pm$ std	px3 error mean $\pm$ std
Three different runs on the same seed		
brightness	1.235 $\pm$ 0.000	0.935 $\pm$ 0.000
contrast	1.187 $\pm$ 0.000	0.938 $\pm$ 0.000
defocus blur	2.026 $\pm$ 0.000	0.899 $\pm$ 0.000
elastic transform	1.043 $\pm$ 0.000	0.954 $\pm$ 0.000
fog	1.221 $\pm$ 0.000	0.936 $\pm$ 0.000
frost	29.640 $\pm$ 0.000	0.383 $\pm$ 0.000
gaussian noise	5.931 $\pm$ 0.000	0.732 $\pm$ 0.000
glass blur	2.409 $\pm$ 0.000	0.861 $\pm$ 0.000
impulse noise	6.098 $\pm$ 0.220	0.736 $\pm$ 0.002
jpeg compression	1.942 $\pm$ 0.000	0.892 $\pm$ 0.000
motion blur	5.515 $\pm$ 0.000	0.549 $\pm$ 0.000
pixelate	0.785 $\pm$ 0.000	0.960 $\pm$ 0.000
shot noise	4.435 $\pm$ 0.000	0.780 $\pm$ 0.000
snow	41.974 $\pm$ 0.000	0.354 $\pm$ 0.000
zoom blur	4.808 $\pm$ 0.000	0.746 $\pm$ 0.000
One run each using three different seeds		
brightness	1.235 $\pm$ 0.000	0.935 $\pm$ 0.000
contrast	1.187 $\pm$ 0.000	0.938 $\pm$ 0.000
defocus blur	2.026 $\pm$ 0.000	0.899 $\pm$ 0.000
elastic transform	1.041 $\pm$ 0.009	0.953 $\pm$ 0.001
fog	1.282 $\pm$ 0.076	0.937 $\pm$ 0.001
frost	28.783 $\pm$ 0.827	0.391 $\pm$ 0.019
gaussian noise	5.877 $\pm$ 0.026	0.735 $\pm$ 0.001
glass blur	2.532 $\pm$ 0.105	0.859 $\pm$ 0.002
impulse noise	5.856 $\pm$ 0.067	0.737 $\pm$ 0.002
jpeg compression	1.942 $\pm$ 0.000	0.892 $\pm$ 0.000
motion blur	4.135 $\pm$ 0.141	0.565 $\pm$ 0.010
pixelate	0.785 $\pm$ 0.000	0.960 $\pm$ 0.000
shot noise	4.336 $\pm$ 0.143	0.781 $\pm$ 0.004
snow	42.984 $\pm$ 4.786	0.362 $\pm$ 0.007
zoom blur	4.808 $\pm$ 0.000	0.746 $\pm$ 0.000