

# 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 point matching method 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

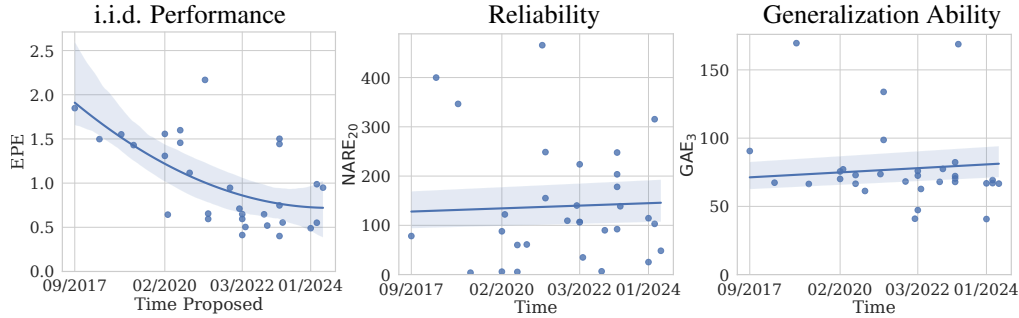


Figure 1: Optical flow estimation methods proposed over time and their reliability and generalization ability. In all three plots, the y-axis represents error, i.e., lower is better. The error of optical flow estimation methods on independent and identically distributed data samples (i.i.d.) has decreased 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

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

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

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

The main contributions of this work are as follows:

- We provide a benchmarking tool FLOWBENCH to evaluate the performance of most DL-based optical flow estimation methods over different datasets and make 91 checkpoints over different datasets publicly available for streamlined benchmarking while enabling the research community to add further checkpoints.
- We benchmark the aforementioned models against SotA and other commonly used adversarial attacks and common corruptions that can be easily queried using FLOWBENCH.
- We perform an in-depth analysis using FLOWBENCH and present interesting findings showing that methods that are SotA on i.i.d. are remarkably less reliable and generalize worse than other non-SotA methods.
- We analyze correlations between performance, reliability, and generalization abilities of optical flow estimation methods, under various lenses such as point matching methods used, and the number of learnable parameters.
- We show that the optimization of white-box adversarial attacks for optical flow estimation can be performed even without the availability of ground truth predictions, furthering the 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 B 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 D.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 D.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 D. Please refer to Appendix F 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 E 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 F.1 for details regarding the arguments):

```
from flowbench.evals import evaluate
model = 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 F.3 (2D Common Corruptions) and Appendix F.4 (3D Common Corruption) for details regarding the arguments):

```
from flowbench.evals import evaluate
model = 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 estimation 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 estimation 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 G 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 G*, 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 Schmalfuss et al. (2022b), 20 iterations are enough to optimize an attack to truly understand the performance of the attacked method.

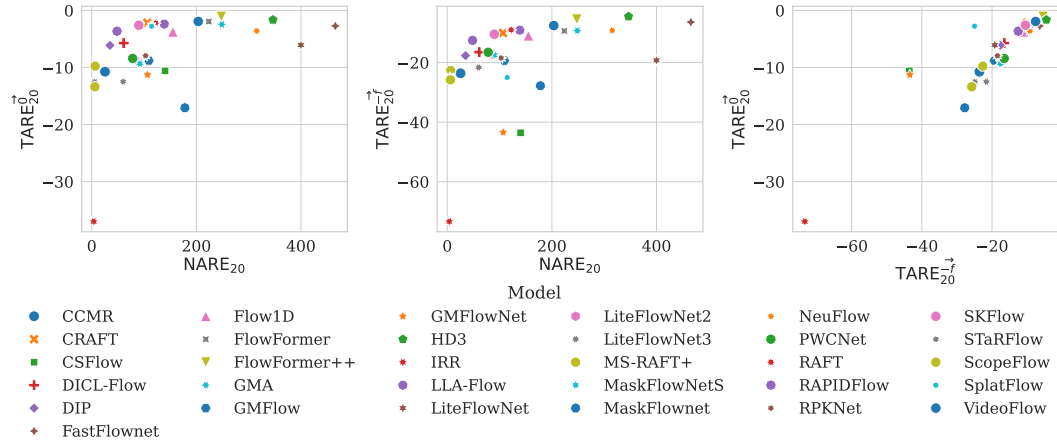


Figure 2: Analysing correlations between Targeted and Non-targeted adversarial attacks. A model is more reliable if it has a low NARM value and a high TARM value.

## 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. Following, we discuss the observations made in the comprehensive robustness benchmark created using FLOWBENCH. Please refer to Appendix B for details on the dataset, Appendix C for additional implementation details, and Appendix G 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, a higher NARE and a lower TARE value indicates low reliability and vice versa. Please refer to Appendix C for more implementation details. We observe in Fig. 2 that there is a very high correlation between the  $TARE_0^T$  and  $TARE_{20}^T$  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 stands out having poor performance and poor generalization ability. Interestingly, as shown by Fig. 3, there is a correlation between the generalization ability (GAE<sub>3</sub> values, introduced in Sec. 4.1, higher GAE

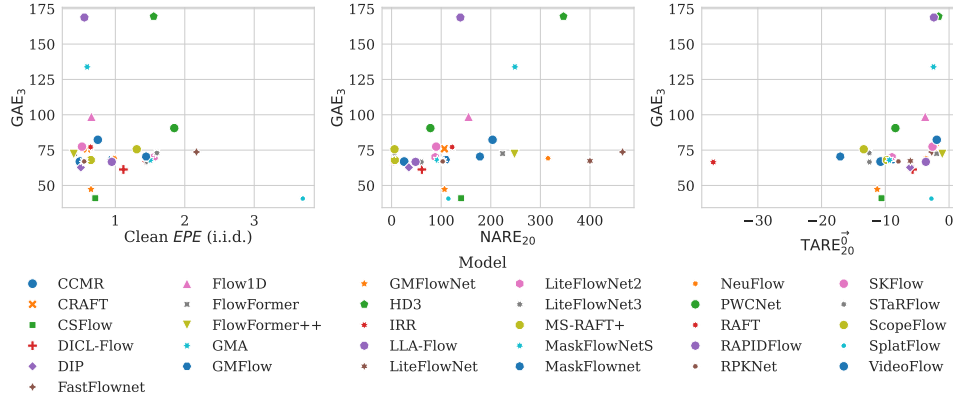


Figure 3: Analysing correlations between reliability and generalization ability of optical flow estimation methods.

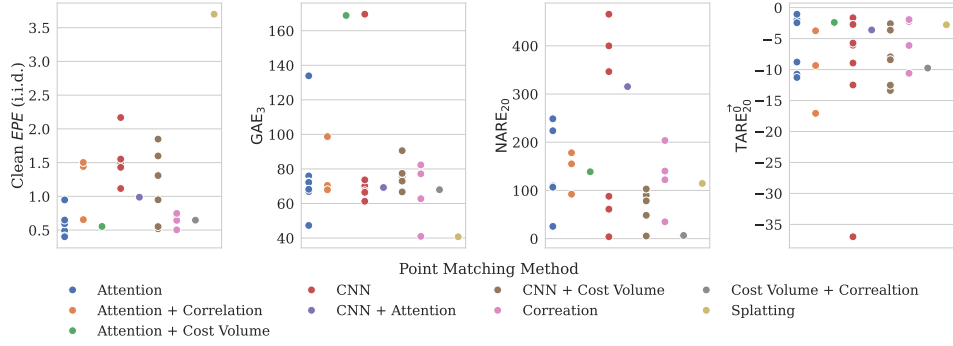


Figure 4: Analysing correlations between the point matching method used by an optical flow estimation method and its corresponding performance, reliability and generalization ability.

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.

### 5.3 ANALYSING POINT MATCHING METHODS

Optical flow estimation methods proposed over the years use different methods for matching points from the first frame to the next. For point matching, all works use either an Attention-based method (Huang et al., 2022; Shi et al., 2023b; Jahedi et al., 2024a), or a Correlation-based method (Shi et al., 2022; Jiang et al., 2021b), or a CNN based method (Dosovitskiy et al., 2015; Ilg et al., 2017; Hui et al., 2018), or a Cost Volume based method (Khairi et al., 2024), or a combination of the two, such as Attention and Cost Volume (Xu et al., 2023b) or Attention and Correlation (Zhao et al., 2020) and others (please refer Tab. 1 for detailed categorization of each method). Thus, based on the observations made in Sec. 5.1 and Sec. 5.2, we determined it would be interesting to observe the relation between the point matching method used by an optical flow estimation method and its performance, reliability, and generalization ability. In Fig. 4 we observe that Attention-based methods have relatively better performance and generalization ability but are also less reliable. Whereas, some CNN and Cost Volume-based methods might not have the best performance but they are reliable and have relatively better generalization abilities. However, some CNN-based methods are highly unreliable as well.

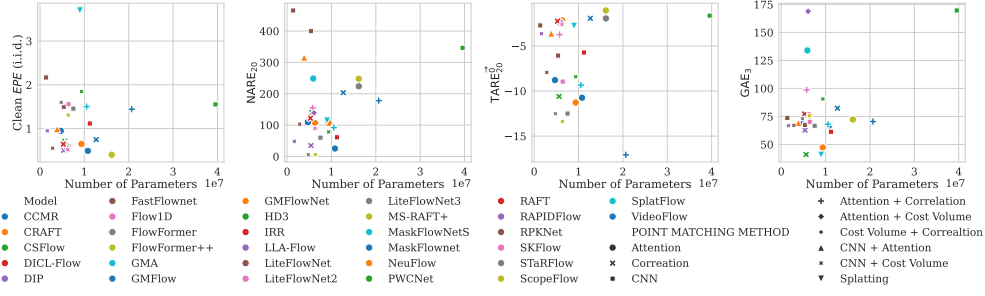


Figure 5: Analysing correlation between the number of learnable parameters in a DL-based optical flow estimation method and its performance, reliability, and generalization ability. Colors show the different optical flow methods while marker styles show the point-matching method used by them.

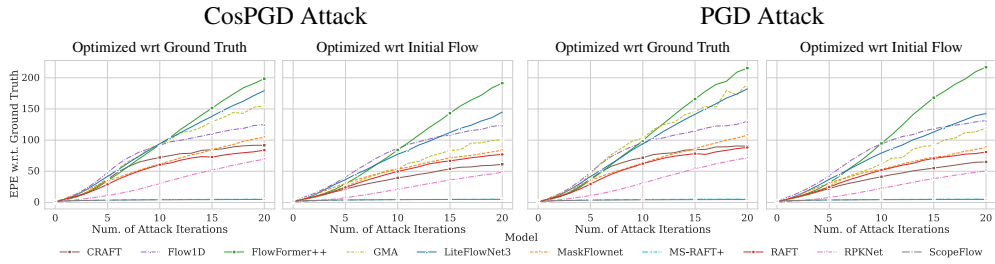


Figure 6: Performance of interesting optical flow estimation methods under different non-targeted adversarial attacks optimized using initial flow predictions on the KITTI2015 dataset.

#### 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++, and VideoFlow have relatively more parameters than other methods however they are less reliable and have a poor generalization ability. On the other hand, methods like STaRFlow, and LiteFlowNet3 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 Schmalfuss 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 predictions (predictions when the method is not attacked), independent of the ground truth. Thus, we attempt to optimize the adversarial attack w.r.t. to the initial flow prediction on the unperturbed input sample before any attacks, as access to this is almost guaranteed. This helps us ascertain if initial flow predictions can be used as a proxy to ground truth while optimizing attacks. Thus, in Eq. (4), Eq. (8), Eq. (9) and there places where applicable  $Y = X^{\text{clean}}$  (please refer Appendix D.1). However, this optimization is only possible for attacks that introduce certain randomness in the initial input sample, as shown by Eq. (7). This allows for there to exist a non-zero loss between the predictions



on the clean input samples and the perturbed input samples allowing for optimization. We report the evaluations for CosPGD and PGD attack using the KITTI2015 dataset for 10 interesting methods in Fig. 6. We choose the optical flow estimation methods on the basis of their performance in Sec. 5 and their performance on i.i.d. samples. For additional evaluation using more models please refer to Appendix G. We observe in Fig. 6 that there appears a high correlation in the performance of all considered methods under attack when optimized using the ground truth flow and the initial flow prediction. Thus, initial flow predictions from methods do serve as a strong proxy to the ground truth for optimizing attacks. This new finding over a big sample, helps advance study in the reliability of optical flow methods, even when ground truth predictions are not available.

## 6 CONCLUSION

NEW

FLOWBENCH is the first robustness benchmarking tool and a novel benchmark for optical flow estimation methods. It currently supports 91 model checkpoints, over distinct datasets, and all relevant robustness evaluation methods including SotA adversarial attacks and image corruptions. We discuss the unique features of FLOWBENCH in detail and demonstrate that the library is user-friendly. Adding new evaluation methods or optical flow estimation methods to FLOWBENCH is easy and intuitive. In Sec. 5.1, we find that there is a high correlation in the performance of optical flow estimation methods against targeted attacks using different targets, thus saving compute for future works as they need to evaluate only against one target. In Sec. 5.2, we observe the methods known to be SotA on i.i.d. samples are not reliable, and do not generalize well to image corruptions, demonstrating the gap in current research when considering real-world applications. Additionally, we observe here that there is no apparent correlation between generalization abilities and the reliability of optical flow estimation methods. In Sec. 5.3, we show that methods using attention-based pointing matching are marginally more reliable than methods using other matching techniques, while methods using CNN and Cost Volume-based matching have marginally better generalization abilities. This in conjecture with the previous observation helps us conclude that based on current works, different approaches might be required to attain reliability under attacks and generalization ability to image corruptions. In Sec. 5.4, we show that, unlike image classification, increasing the number of learnable parameters does not help increase the robustness of optical flow estimation methods. Lastly, we show that white-box adversarial attacks on optical flow estimation methods can be independent of the availability of ground truth information, and can harness the information in the initial flow predictions to optimize attacks, thus overcoming a huge limitation in the field. Such an in-depth understanding of reliability and generalization abilities to optical flow estimation methods can only be obtained using our proposed FLOWBENCH. We are certain that FLOWBENCH will be immensely helpful to gather more such interesting findings and its comprehensive and consolidated nature would make things easier for the research community.

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

**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.



## 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.

## 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.
- 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.
- MMFlow Contributors. MMFlow: Openmmlab optical flow toolbox and benchmark. <https://github.com/open-mmlab/mmdetection>, 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.
- Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11963–11975, 2022.
- Qiaole Dong, Chenjie Cao, and Yanwei Fu. Rethinking optical flow from geometric matching consistent perspective. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1337–1347, 2023.
- Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *Proc. of the IEEE international conference on computer vision*, pp. 2758–2766, 2015.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2018.

- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, nov 2020.
- Pierre Godet, Alexandre Boulch, Aurélien Plyer, and Guy Le Besnerais. STaRFlow: A spatiotemporal recurrent cell for lightweight multi-frame optical flow estimation. In *Proc. IEEE International Conference on Pattern Recognition (ICPR)*, pp. 2462–2469, 2021.
- Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proc. International Conference on Learning Representations (ICLR)*, 2015.
- Julia Grabinski, Paul Gavrikov, Janis Keuper, and Margret Keuper. Robust models are less overconfident. *NeurIPS*, 2022.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *Proc. International Conference on Learning Representations (ICLR)*, 2019.
- Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. AugMix: A simple data processing method to improve robustness and uncertainty. *Proc. of the International Conference on Learning Representations (ICLR)*, 2020.
- J Hoffmann, S Agnihotri, Tonmoy Saikia, and Thomas Brox. Towards improving robustness of compressed cnns. In *ICML Workshop on Uncertainty and Robustness in Deep Learning (UDL)*, 2021.
- Berthold KP Horn and Brian G Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3): 185–203, 1981.
- 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.
- Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. LiteFlowNet: A lightweight convolutional neural network for optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8981–8989, 2018.
- Tak-Wai Hui, Xiaoou Tang, and Chen Change Loy. A lightweight optical flow CNN—revisiting data fidelity and regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 43(8):2555–2569, 2020.
- Junhwa Hur and Stefan Roth. Iterative residual refinement for joint optical flow and occlusion estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5754–5763, 2019.
- Eddy Ilg, Nikolaus Mayer, Tonmoy Saikia, Margret Keuper, Alexey Dosovitskiy, and Thomas Brox. FlowNet 2.0: Evolution of optical flow estimation with deep networks. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2462–2470, 2017.
- Azin Jahedi, Maximilian Luz, Marc Rivinius, and Andrés Bruhn. CCMR: High resolution optical flow estimation via coarse-to-fine context-guided motion reasoning. In *Proc. IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 6899–6908, 2024a.
- Azin Jahedi, Maximilian Luz, Marc Rivinius, Lukas Mehl, and Andrés Bruhn. MS-RAFT+: high resolution multi-scale raft. *International Journal of Computer Vision (IJCV)*, 132(5):1835–1856, 2024b.
- Shihao Jiang, Dylan Campbell, Yao Lu, Hongdong Li, and Richard Hartley. Learning to estimate hidden motions with global motion aggregation. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9772–9781, 2021a.

- Shihao Jiang, Yao Lu, Hongdong Li, and Richard Hartley. Learning optical flow from a few matches. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16592–16600, 2021b.
- Steffen Jung, Jovita Lukasik, and Margret Keuper. Neural architecture design and robustness: A dataset. In *ICLR. OpenReview. net*, 2023.
- Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18963–18974, 2022.
- Sarra Khairi, Etienne Meunier, Renaud Fraisse, and Patrick Bouthemy. Efficient local correlation volume for unsupervised optical flow estimation on small moving objects in large satellite images. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 440–448, 2024.
- Hoki Kim. Torchattacks: A pytorch repository for adversarial attacks. *arXiv preprint arXiv:2010.01950*, 2020.
- Lingtong Kong, Chunhua Shen, and Jie Yang. Fastflownet: A lightweight network for fast optical flow estimation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 10310–10316, 2021.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 and cifar-100 datasets, learning multiple layers of features from tiny images. URL: <https://www.cs.toronto.edu/kriz/cifar.html>, 6(1): 1, 2009.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *Proc. International Conference on Learning Representations (ICLR)*, 2017.
- Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *Artificial Intelligence Safety and Security*, pp. 99–112. Chapman and Hall/CRC, 2018.
- Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proc. of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11976–11986, 2022.
- Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *IJCAI’81: 7th international joint conference on Artificial intelligence*, volume 2, pp. 674–679, 1981.
- Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4040–4048, 2016.
- Lukas Mehl, Jenny Schmalfuss, Azin Jahedi, Yaroslava Nalivayko, and Andrés Bruhn. Spring: A high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4981–4991, 2023.
- Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3061–3070, 2015.
- Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. In *Proc. of the IEEE conference on computer vision and pattern recognition*, pp. 2574–2582, 2016.
- Henrique Morimitsu. Pytorch lightning optical flow. <https://github.com/hmorimitsu/ptlflow>, 2021.
- Henrique Morimitsu, Xiaobin Zhu, Roberto M Cesar, Xiangyang Ji, and Xu-Cheng Yin. RAPID-Flow: Recurrent adaptable pyramids with iterative decoding for efficient optical flow estimation. In *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2946–2952, 2024a.

- Henrique Morimitsu, Xiaobin Zhu, Xiangyang Ji, and Xu-Cheng Yin. Recurrent partial kernel network for efficient optical flow estimation. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2024b.
- Adarsh Prasad. *Towards Robust and Resilient Machine Learning*. PhD thesis, Carnegie Mellon University, 2022.
- Anurag Ranjan, Joel Janai, Andreas Geiger, and Michael J Black. Attacking optical flow. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2404–2413, 2019.
- Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel. Foolbox native: Fast adversarial attacks to benchmark the robustness of machine learning models in pytorch, tensorflow, and jax. *Journal of Open Source Software*, 5(53):2607, 2020. doi: 10.21105/joss.02607. URL <https://doi.org/10.21105/joss.02607>.
- Benoît Rosa, Valentin Bordoux, and Florent Nageotte. Combining differential kinematics and optical flow for automatic labeling of continuum robots in minimally invasive surgery. *Frontiers in Robotics and AI*, 6:86, 2019.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Christos Sakaridis, Dengxin Dai, and Luc Van Gool. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.
- Erik Scheurer, Jenny Schmalfuss, Alexander Lis, and Andrés Bruhn. Detection defenses: An empty promise against adversarial patch attacks on optical flow. In *Proc. IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2024.
- Jenny Schmalfuss, Lukas Mehl, and Andrés Bruhn. Attacking motion estimation with adversarial snow. In *Proc. ECCV Workshop on Adversarial Robustness in the Real World (AROW)*, 2022a. doi: 10.48550/arXiv.2210.11242. URL <https://doi.org/10.48550/arXiv.2210.11242>.
- Jenny Schmalfuss, Philipp Scholze, and Andrés Bruhn. A perturbation-constrained adversarial attack for evaluating the robustness of optical flow. In *Proc. European Conference on Computer Vision (ECCV)*, LNCS 13682, pp. 183–200, 2022b.
- Jenny Schmalfuss, Lukas Mehl, and Andrés Bruhn. Distracting downpour: Adversarial weather attacks for motion estimation. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10106–10116, 2023.
- Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. *NIPS*, 31, 2018.
- Simon Schrodi, Tonmoy Saikia, and Thomas Brox. Towards understanding adversarial robustness of optical flow networks. In *CVPR*, pp. 8916–8924, 2022.
- Hao Shi, Yifan Zhou, Kailun Yang, Xiaoting Yin, and Kaiwei Wang. Csflo: Learning optical flow via cross strip correlation for autonomous driving. In *IEEE Intelligent Vehicles Symposium (IV)*, pp. 1851–1858, 2022.
- Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. VideoFlow: Exploiting temporal cues for multi-frame optical flow estimation. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 12469–12480, 2023a.
- Xiaoyu Shi, Zhaoyang Huang, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. FlowFormer++: Masked cost volume autoencoding for pretraining optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1599–1610, 2023b.

- Xiuchao Sui, Shaohua Li, Xue Geng, Yan Wu, Xinxing Xu, Yong Liu, Rick Goh, and Hongyuan Zhu. CRAFT: Cross-attentional flow transformer for robust optical flow. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17602–17611, 2022.
- Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8934–8943, 2018.
- Shangkun Sun, Yuanqi Chen, Yu Zhu, Guodong Guo, and Ge Li. SKFlow: Learning optical flow with super kernels. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:11313–11326, 2022.
- Shiyu Tang, Ruihao Gong, Yan Wang, Aishan Liu, Jiakai Wang, Xinyun Chen, Fengwei Yu, Xianglong Liu, Dawn Song, Alan Yuille, Philip H.S. Torr, and Dacheng Tao. Robustart: Benchmarking robustness on architecture design and training techniques. <https://arxiv.org/pdf/2109.05211.pdf>, 2021.
- Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *Proc. European Conference on Computer Vision (ECCV)*, LNCS 12347, pp. 402–419, 2020.
- Dustin Tran, Jeremiah Liu, Michael W Dusenberry, Du Phan, Mark Collier, Jie Ren, Kehang Han, Zi Wang, Zelda Mariet, Huiyi Hu, et al. Plex: Towards reliability using pretrained large model extensions. *arXiv preprint arXiv:2207.07411*, 2022.
- Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *ICLR*, 2019.
- Bo Wang, Yifan Zhang, Jian Li, Yang Yu, Zhenping Sun, Li Liu, and Dewen Hu. Splatflow: Learning multi-frame optical flow via splatting. *International Journal of Computer Vision*, pp. 1–23, 2024.
- Hengli Wang, Peide Cai, Yuxiang Sun, Lujia Wang, and Ming Liu. Learning interpretable end-to-end vision-based motion planning for autonomous driving with optical flow distillation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 13731–13737. IEEE, 2021.
- Jianyuan Wang, Yiran Zhong, Yuchao Dai, Kaihao Zhang, Pan Ji, and Hongdong Li. Displacement-invariant matching cost learning for accurate optical flow estimation. *Advances in Neural Information Processing Systems (NeurIPS)*, 33:15220–15231, 2020.
- Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. *ArXiv*, abs/2001.03994, 2020.
- J. Wulff, D. J. Butler, G. B. Stanley, and M. J. Black. Lessons and insights from creating a synthetic optical flow benchmark. In A. Fusiello et al. (Eds.) (ed.), *ECCV Workshop on Unsolved Problems in Optical Flow and Stereo Estimation*, Part II, LNCS 7584, pp. 168–177. Springer-Verlag, October 2012.
- Haofei Xu, Jiaolong Yang, Jianfei Cai, Juyong Zhang, and Xin Tong. High-resolution optical flow from 1d attention and correlation. In *Proc. IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 10498–10507, 2021a.
- Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, and Dacheng Tao. GMFlow: Learning optical flow via global matching. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8121–8130, 2022.
- Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezaatofghi, Fisher Yu, Dacheng Tao, and Andreas Geiger. Unifying flow, stereo and depth estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 45(11):13941–13958, 2023a.
- Jiawei Xu, Zongqing Lu, and Qingmin Liao. LLA-Flow: A lightweight local aggregation on cost volume for optical flow estimation. In *IEEE International Conference on Image Processing (ICIP)*, pp. 3220–3224, 2023b.

- Xiaogang Xu, Hengshuang Zhao, and Jiaya Jia. Dynamic divide-and-conquer adversarial training for robust semantic segmentation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 7466–7475, 2021b. doi: 10.1109/ICCV48922.2021.00739.
- Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019.
- Zhichao Yin, Trevor Darrell, and Fisher Yu. Hierarchical discrete distribution decomposition for match density estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6044–6053, 2019.
- Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable Flow: Learning motion cost volumes for optical flow estimation. In *Proc. IEEE/CVF International Conference on Computer Vision (CVPR)*, pp. 10807–10817, 2021.
- Zhiyong Zhang, Huaizu Jiang, and Hanumant Singh. NeufLOW: Real-time, high-accuracy optical flow estimation on robots using edge devices. *arXiv preprint arXiv:2403.10425*, 2024.
- Shengyu Zhao, Yilun Sheng, Yue Dong, Eric I Chang, Yan Xu, et al. MaskFlowNet: Asymmetric feature matching with learnable occlusion mask. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6278–6287, 2020.
- Shiyu Zhao, Long Zhao, Zhixing Zhang, Enyu Zhou, and Dimitris Metaxas. Global matching with overlapping attention for optical flow estimation. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 17592–17601, 2022.
- Zihua Zheng, Ni Nie, Zhi Ling, Pengfei Xiong, Jiangyu Liu, Hao Wang, and Jiankun Li. DIP: Deep inverse patchmatch for high-resolution optical flow. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8925–8934, 2022.



# 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: Details for the datasets used.
  - Appendix B.1: FlyingThings3D
  - Appendix B.2: KITTI2015
  - Appendix B.3: MPI Sintel
  - Appendix B.4: Spring
- Appendix C: Additional implementation details for the evaluated benchmark.
- Appendix D: In detail description of the attacks.
- Appendix E: 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 F: In detail explanation of the available functionalities of the FLOWBENCH benchmarking tool and description of the arguments for each function.
- Appendix G: 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 G.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 G.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 G.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 G.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 G.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.
  - Appendix G.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.

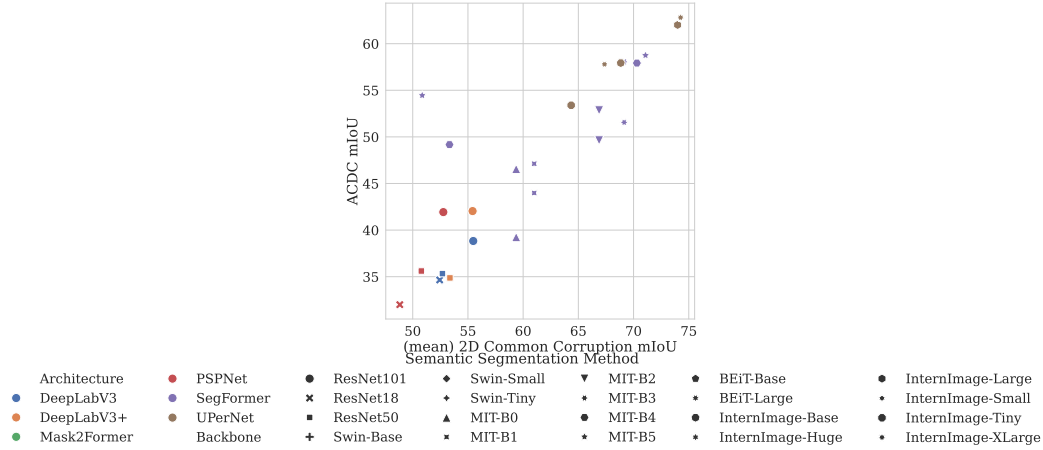


Figure 7: Results from work by Anonymous. Here they find a **very strong positive correlation between mean mIoU over the ACDC evaluation dataset (Sakaridis et al., 2021) and mean mIoU over each 2D Common Corruption (Hendrycks & Dietterich, 2019) over the Cityscapes dataset (Cordts et al., 2016)**. All 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 G.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.

## 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 for evaluating 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 DATASET DETAILS

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

## B.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.

## B.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.

## B.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.

## B.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).

# C 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 G 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 G*, 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; Schrodri 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 Schmalfuss et al. (2023).

## D 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 B 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 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 D.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 D.2.

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

### D.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,

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

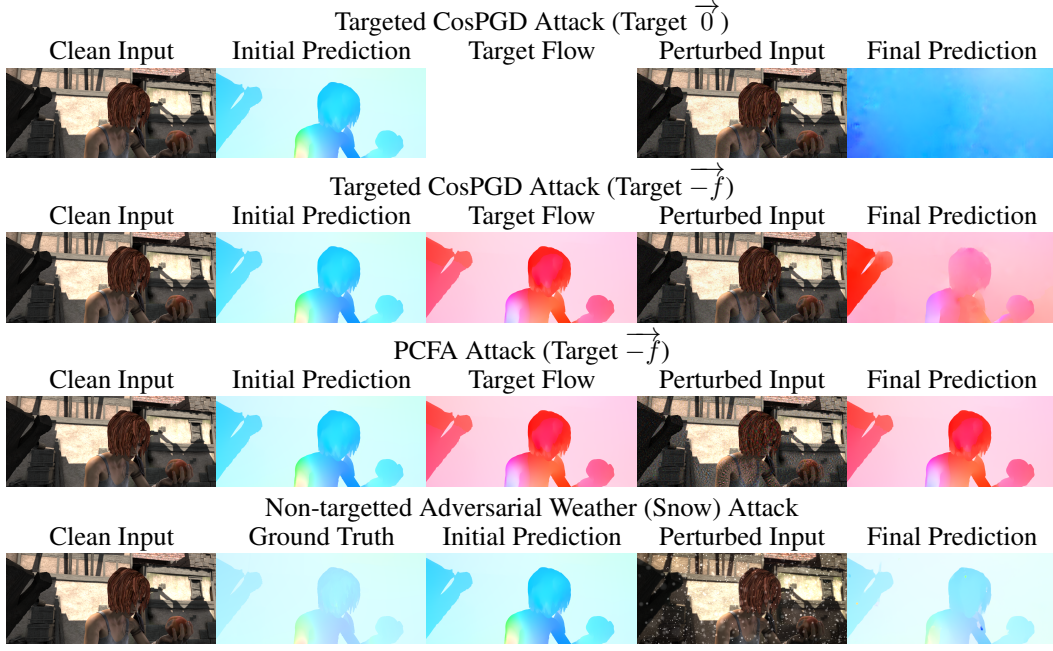


Figure 8: Examples of MPI Sintel images perturbed by the mentioned adversarial attacks and the optical flow predictions using FlowFormer++. These examples are intended to show the versatility of FLOWBENCH.

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

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

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  $\ell_p$ -norm bounded attacks, complying with the  $\ell_\infty$ -norm or other  $\ell_p$ -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,

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

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

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

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  $\cos$  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  $\underset{\hat{\delta}}{\text{argmin}}(\cdot)$  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  $\delta_{\mathbf{P}^t} = \mathbf{P}^{\text{init}} - \mathbf{P}^t$ :

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

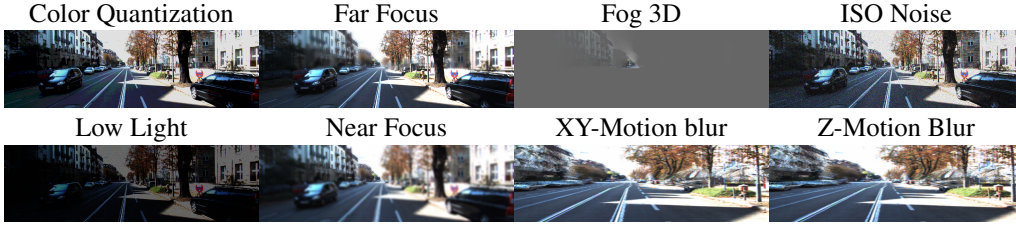


Figure 9: Examples of images from KITTI2015 corrupted using 3D Common Corruptions for evaluation of OOD robustness.

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.

## D.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 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.

## E 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

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

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

## F FLOWBENCH USAGE DETAILS

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

### F.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 = 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.
- ‘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.

## F.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 = 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:

- ‘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.
- ‘weather’ is the name of the weather condition in adversarial weather attack to be used, given as a string, options include ‘snow’, ‘fog’, ‘rain’ and ‘sparks’.
- ‘num\_particles’ is the number of particles per frame to be used, given as an integer.
- ‘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.

## F.3 2D COMMON CORRUPTIONS

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

```
from flowbench.evals import evaluate
model = evaluate(model_name='RAFT', dataset='KITTI2015',
                 threat_model='2DCommonCorruption',
                 severity=3, 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 common corruption to be used, given as a string.
- ‘severity’ is the severity of the corruption, given as an integer between 1 and 5 (both inclusive).
- ‘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.

FLOWBENCH supports the following 2D Common Corruption: ‘gaussian\_noise’, ‘shot\_noise’, ‘impulse\_noise’, ‘defocus\_blur’, ‘frosted\_glass\_blur’, ‘motion\_blur’, ‘zoom\_blur’, ‘snow’, ‘frost’, ‘fog’, ‘brightness’, ‘contrast’, ‘elastic’, ‘pixelate’, ‘jpeg’. For the evaluation, FLOWBENCH will evaluate the model on the validation images from the respective dataset corrupted using each of the aforementioned corruptions for the given severity, and then report the mean performance over all of them.

#### F.4 3D COMMON CORRUPTIONS

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

```
from flowbench.evals import evaluate
model = evaluate(model_name='RAFT', dataset='KITTI2015',
                  threat_model='3DCommonCorruption',
                  severity=3, 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 common corruption to be used, given as a string.
- ‘severity’ is the severity of the corruption, given as an integer between 1 and 5 (both inclusive).
- ‘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.

FLOWBENCH supports the following 3D Common Corruption: ‘color\_quant’, ‘far\_focus’, ‘fog\_3d’, ‘iso\_noise’, ‘low\_light’, ‘near\_focus’, ‘xy\_motion\_blur’, and ‘z\_motion\_blur’. For the evaluation, FLOWBENCH will evaluate the model on the validation images from the respective dataset corrupted using each of the aforementioned corruptions for the given severity, and then report the mean performance over all of them.

## G ADDITIONAL RESULTS

Following we include additional results from the benchmark made using FLOWBENCH.

### G.1 ADVERSARIAL ATTACKS

Here we report additional results for all adversarial attacks.

### G.1.1 FGSM ATTACK

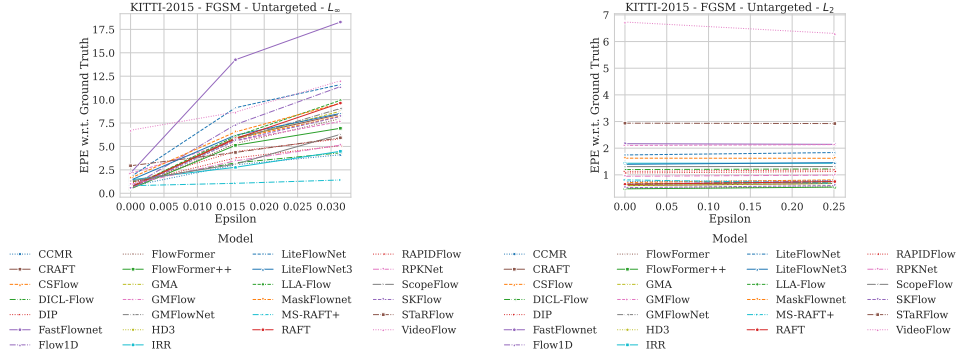


Figure 10: Evaluations for non-targeted FGSM attack under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

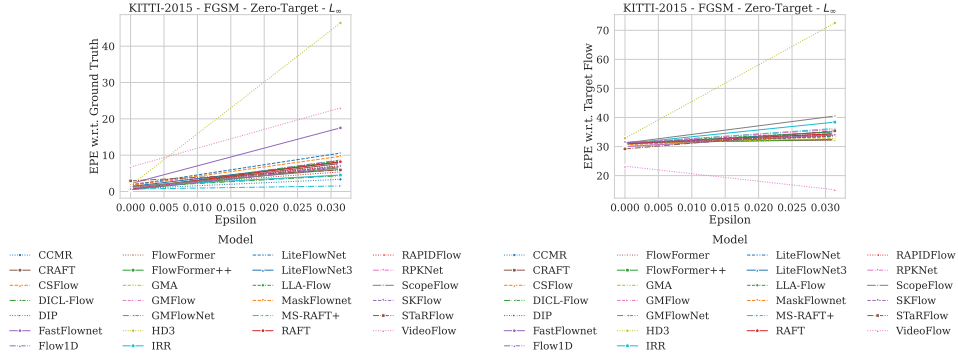


Figure 11: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

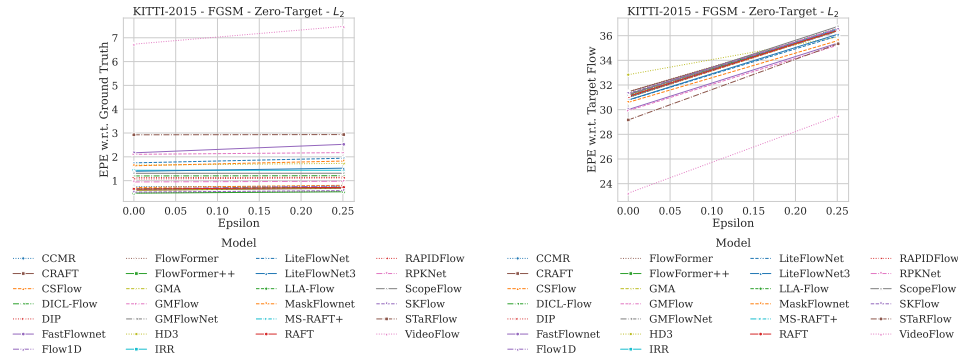


Figure 12: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

Here we report the evaluations using FGSM 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. For  $\ell_\infty$ -norm bound, perturbation budget  $\epsilon = \frac{8}{255}$ , while for  $\ell_2$ -norm bound, perturbation budget  $\epsilon = \frac{64}{255}$ .

Attack evaluations include Fig. 10, Fig. 11, Fig. 12, Fig. 13, Fig. 14, Fig. 15, Fig. 16, Fig. 17, Fig. 18, Fig. 19, Fig. 20, Fig. 21, Fig. 22, Fig. 23, and Fig. 24.



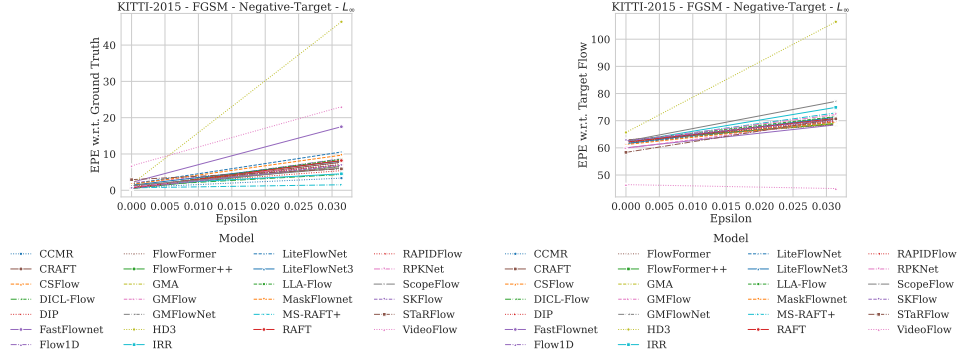


Figure 13: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

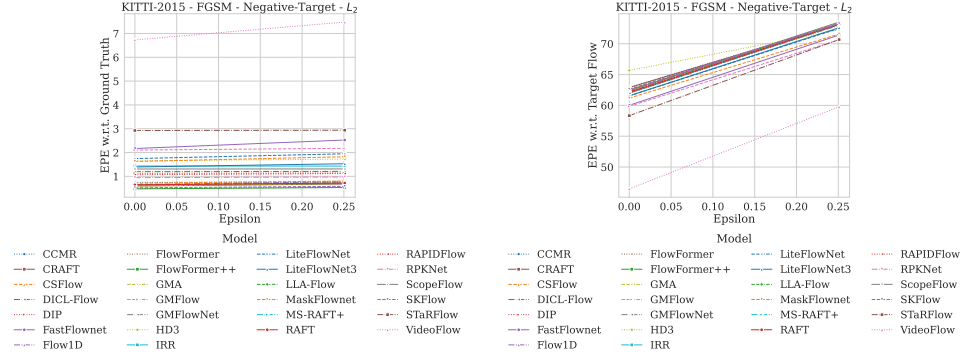


Figure 14: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

### G.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.

### G.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.

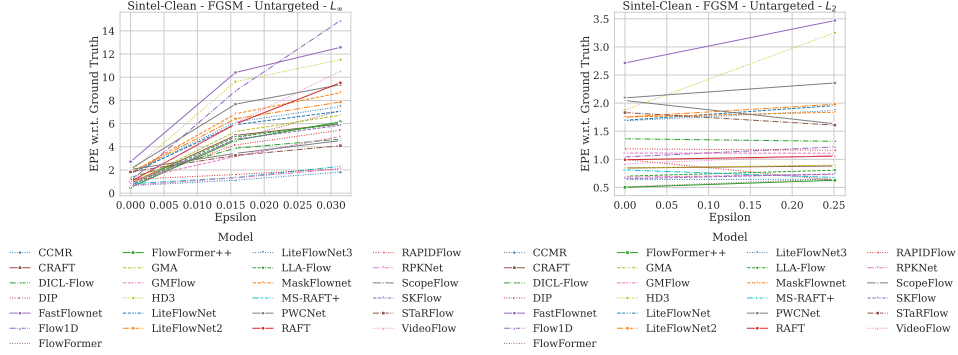


Figure 15: Evaluations for non-targeted FGSM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

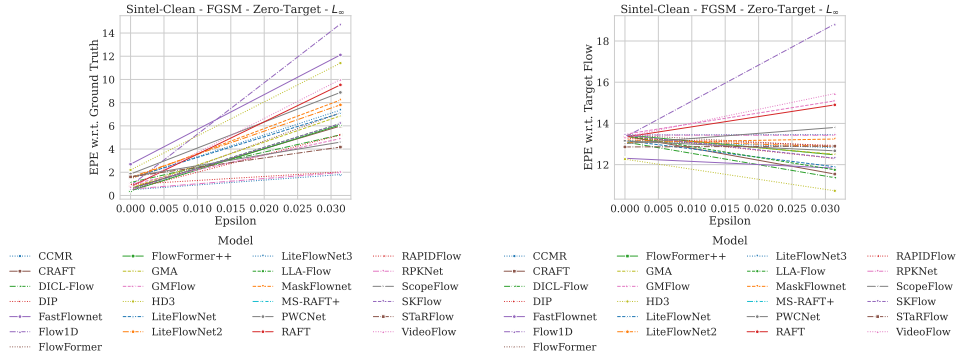


Figure 16: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

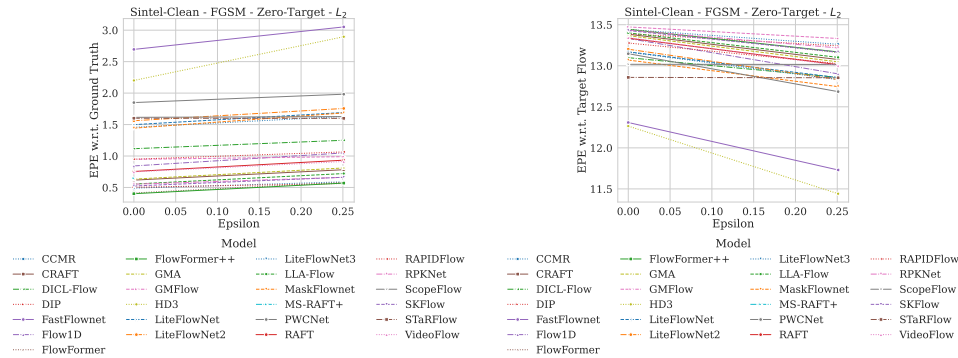


Figure 17: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

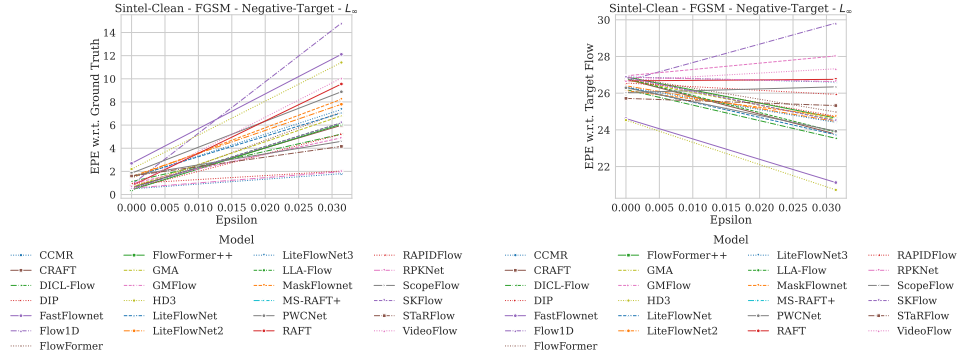


Figure 18: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

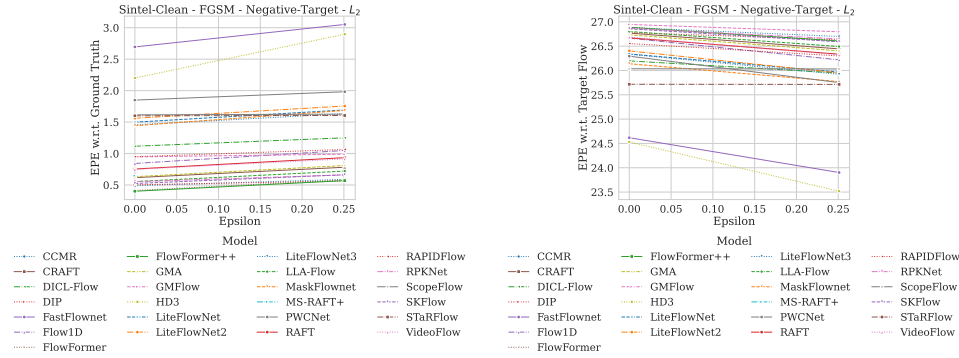


Figure 19: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

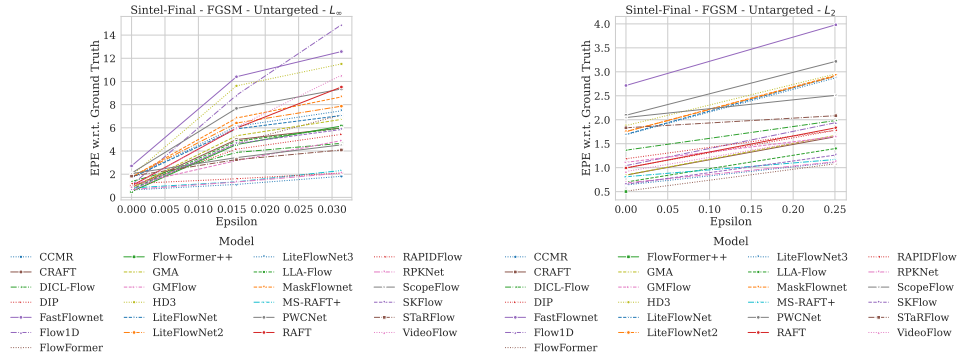


Figure 20: Evaluations for non-targeted FGSM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

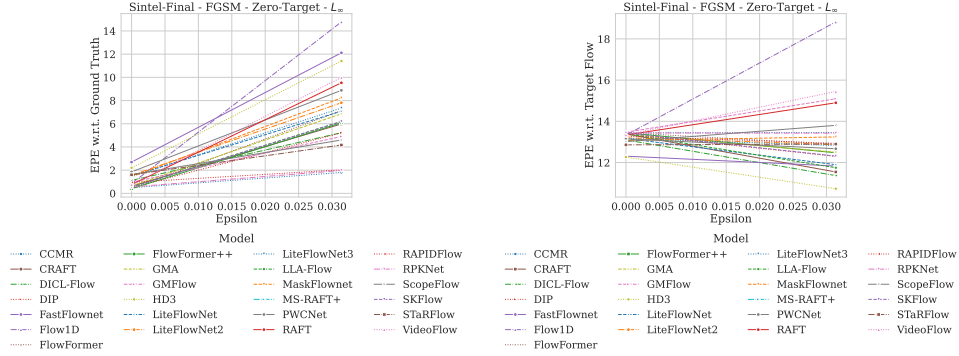


Figure 21: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

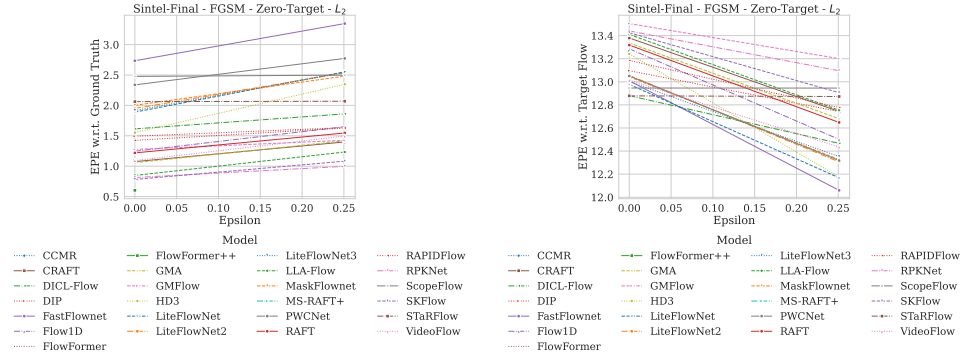


Figure 22: Evaluations for targeted FGSM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

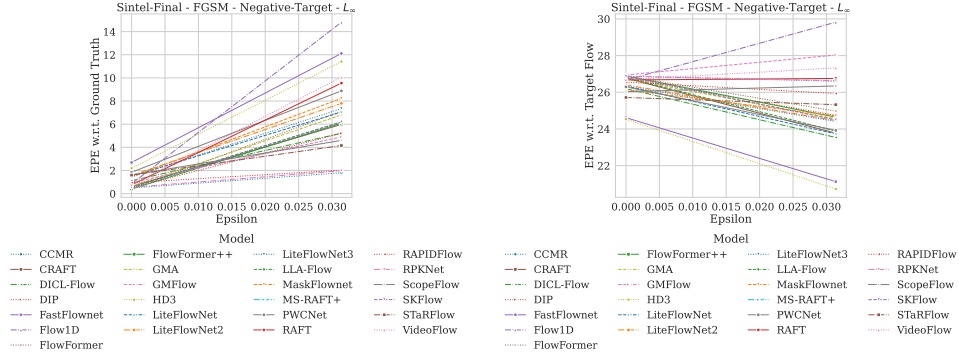


Figure 23: Evaluations for targeted FGSM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

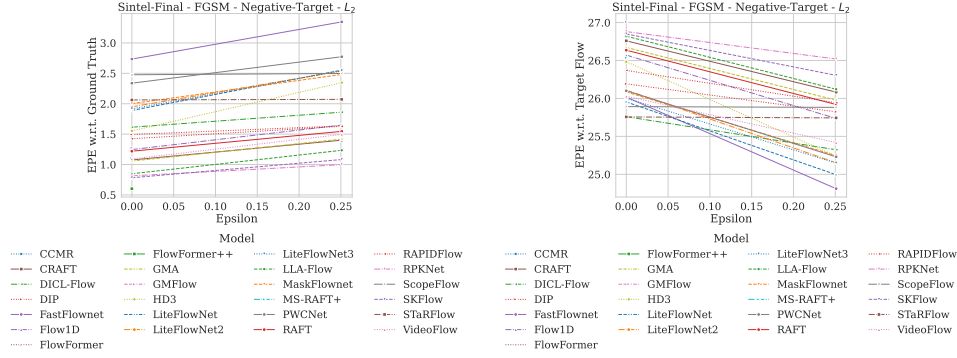


Figure 24: Evaluations for targeted FGSM attack with target  $\vec{-f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

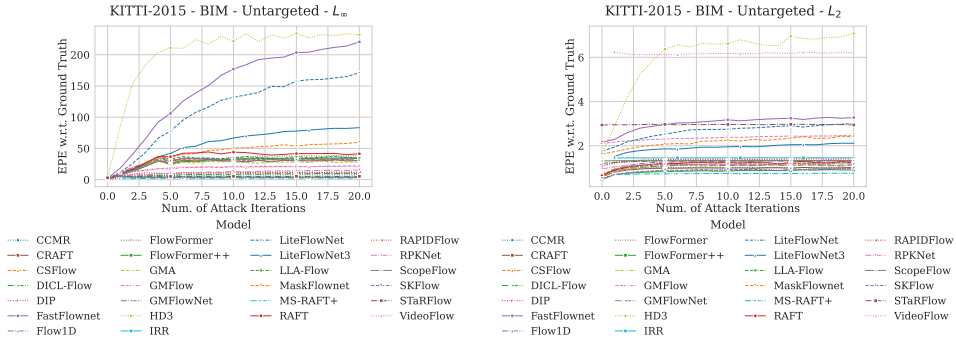


Figure 25: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

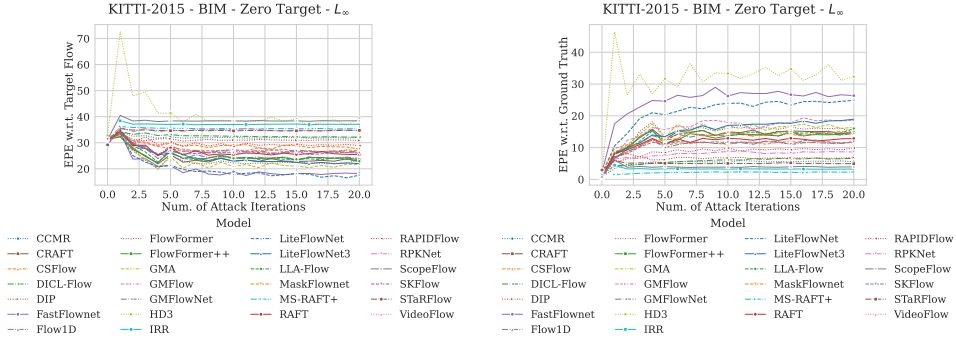


Figure 26: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

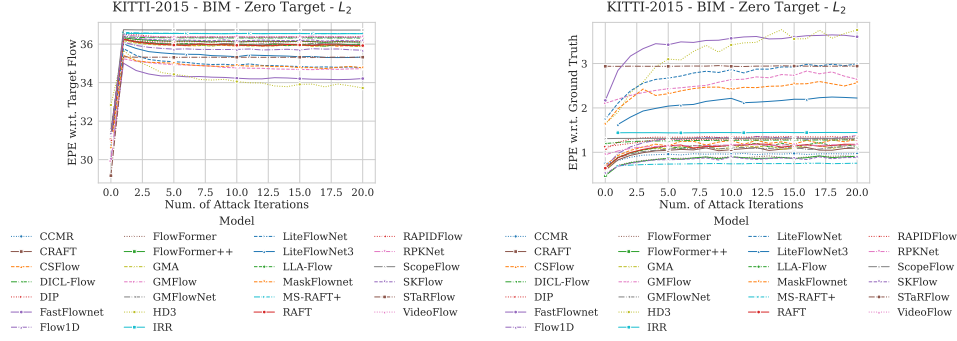


Figure 27: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

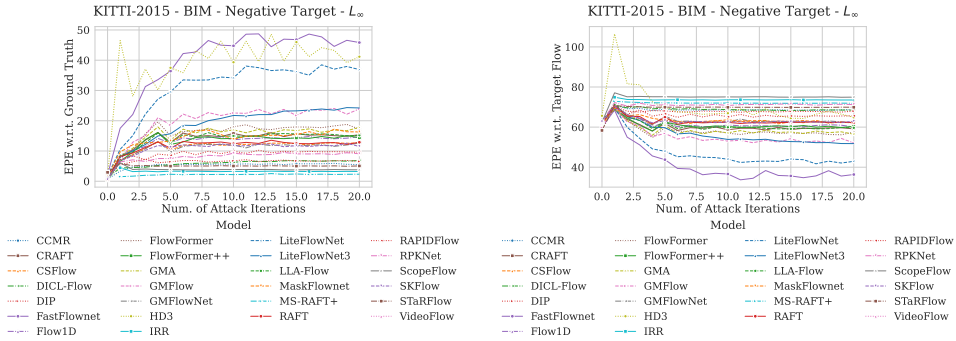


Figure 28: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

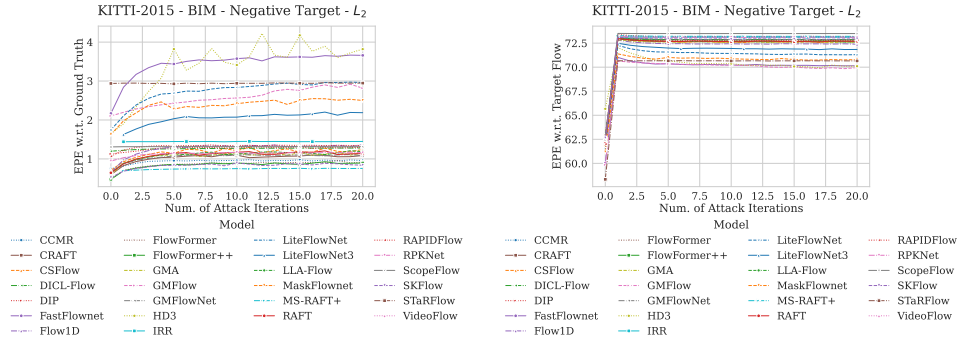


Figure 29: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

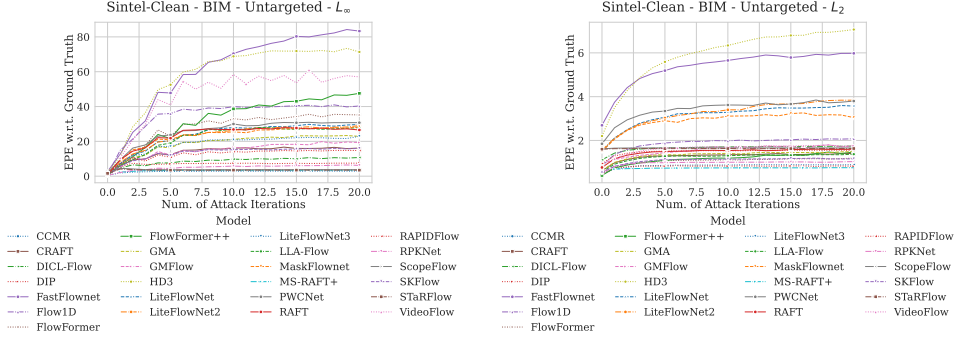


Figure 30: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

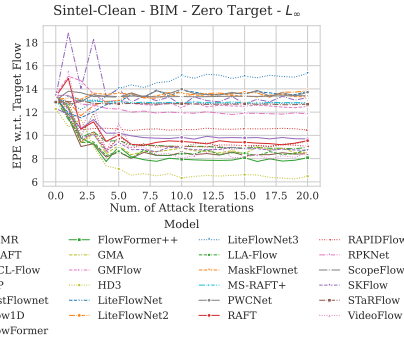


Figure 31: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

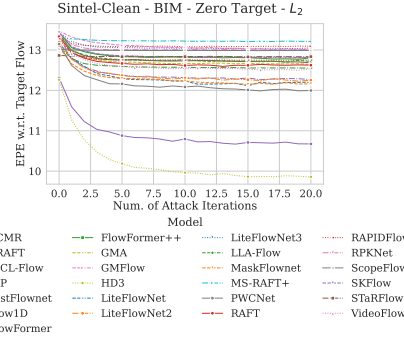


Figure 32: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.



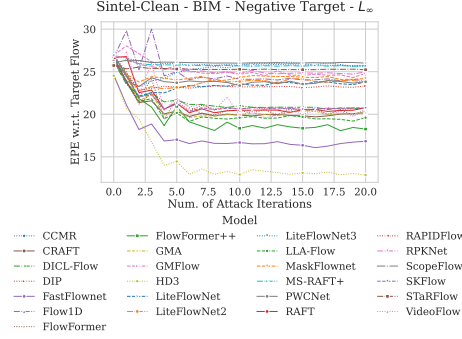


Figure 33: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

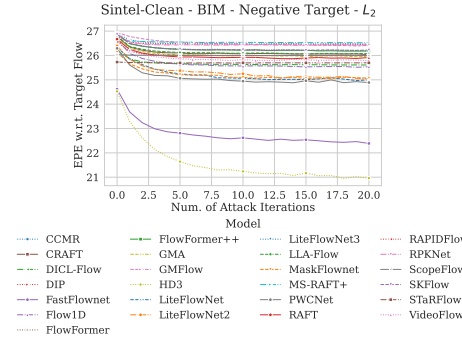


Figure 34: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

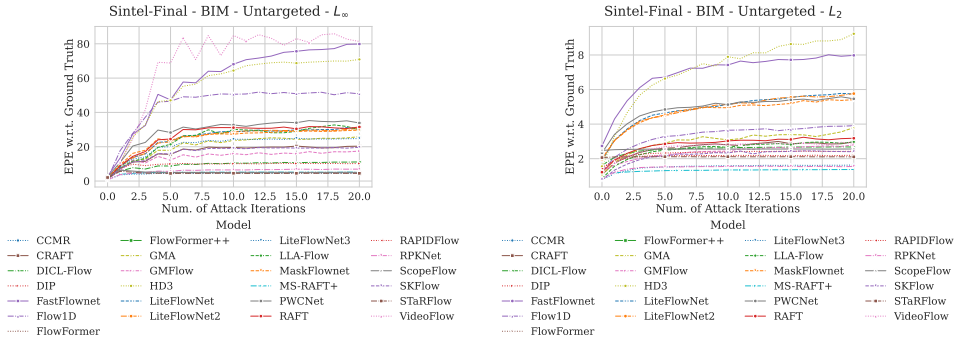


Figure 35: Evaluations for non-targeted BIM attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

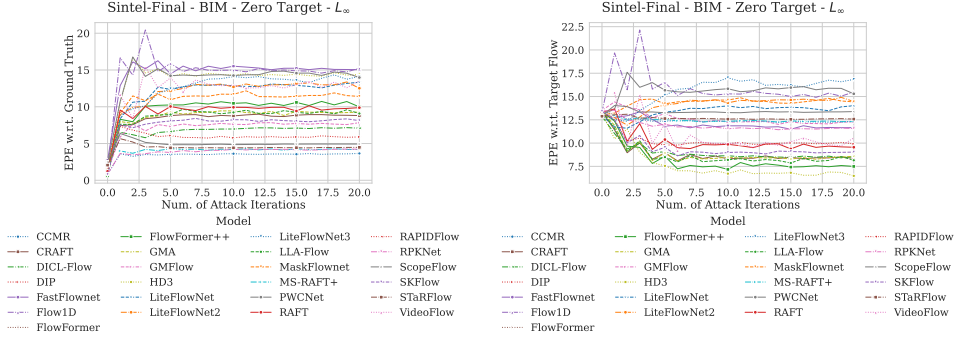


Figure 36: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

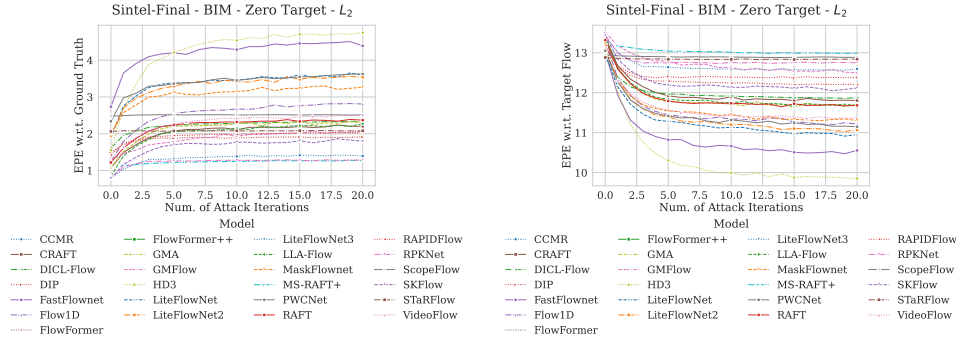


Figure 37: Evaluations for targeted BIM attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

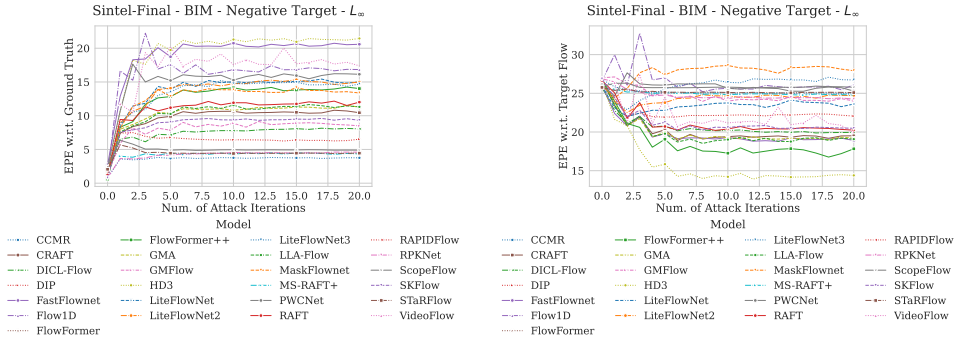


Figure 38: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

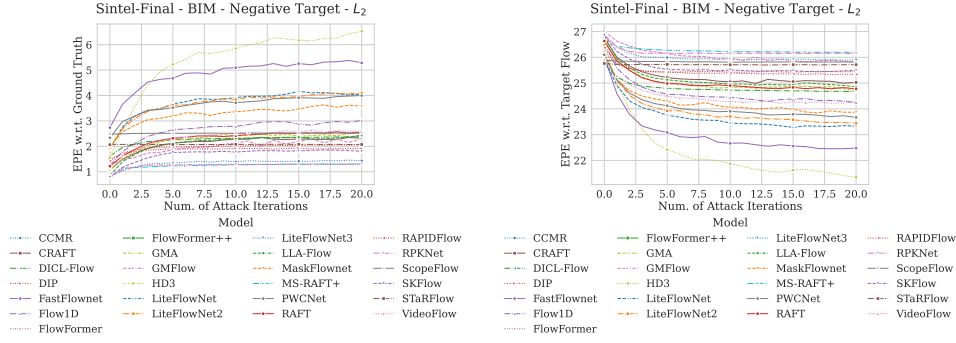


Figure 39: Evaluations for targeted BIM attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

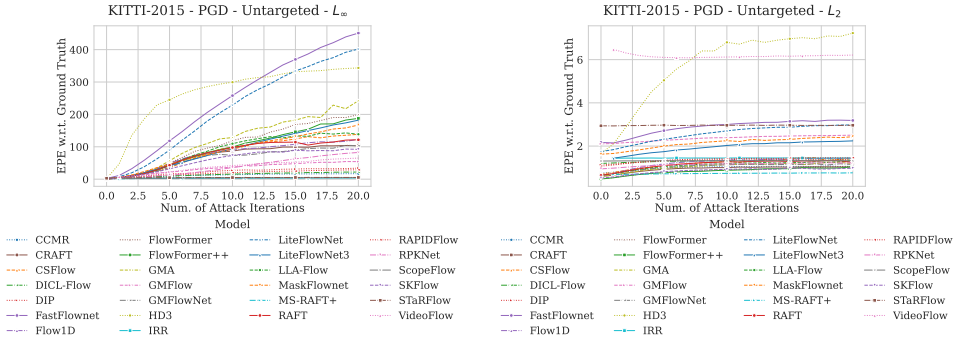


Figure 40: Evaluations for non-targeted PGD attack under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

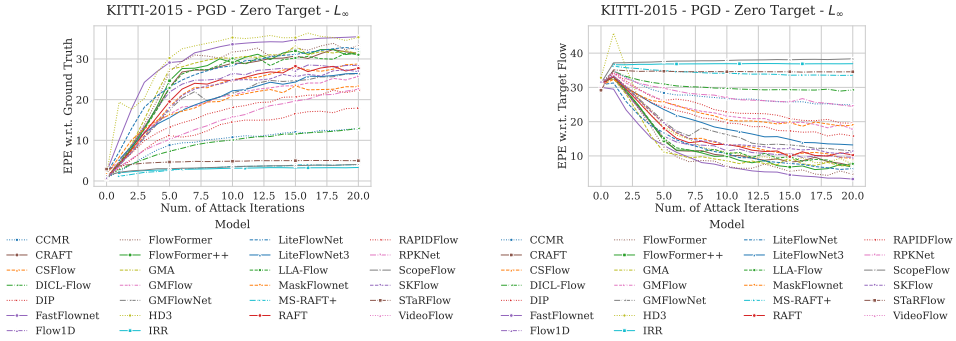


Figure 41: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

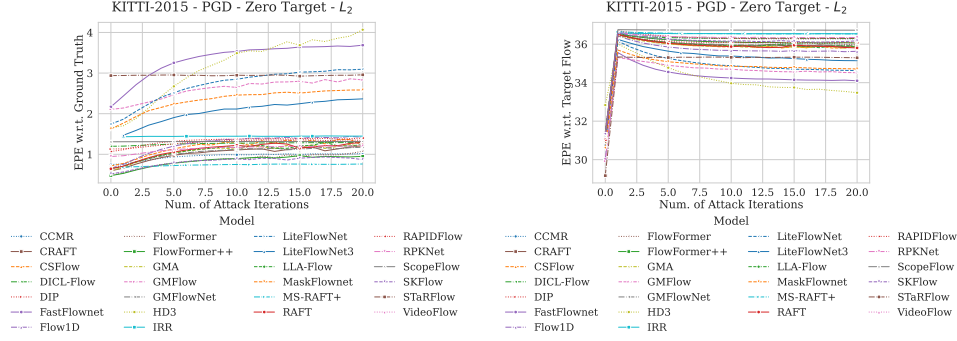


Figure 42: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

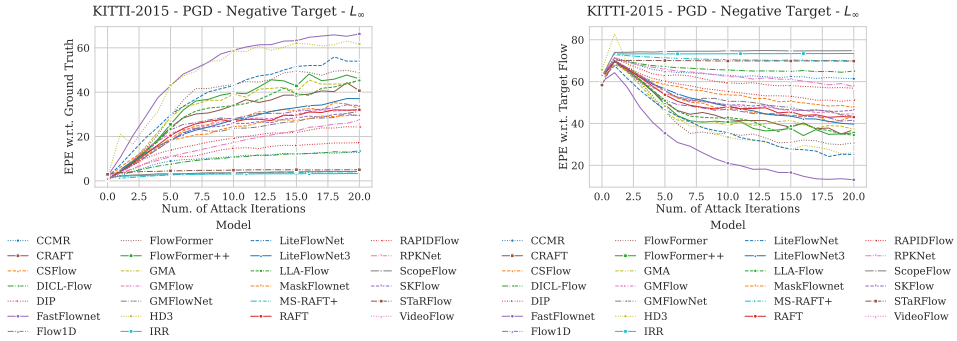


Figure 43: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

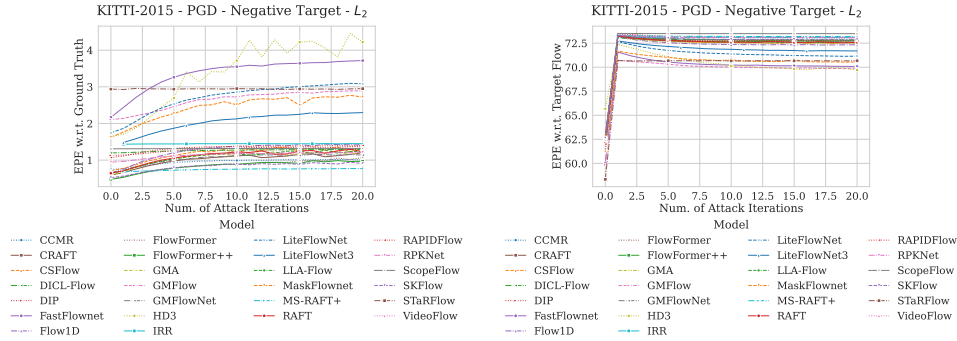


Figure 44: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

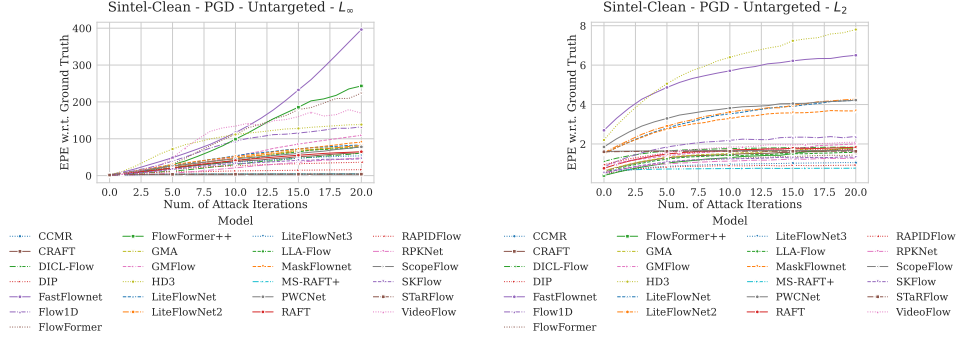


Figure 45: Evaluations for non-targeted PGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

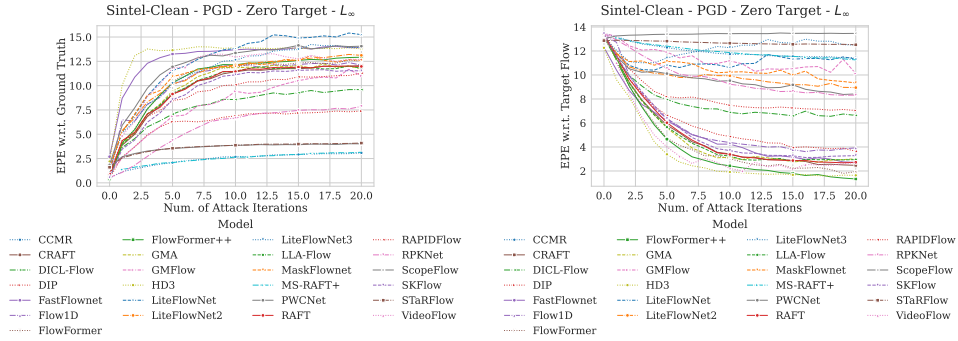


Figure 46: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

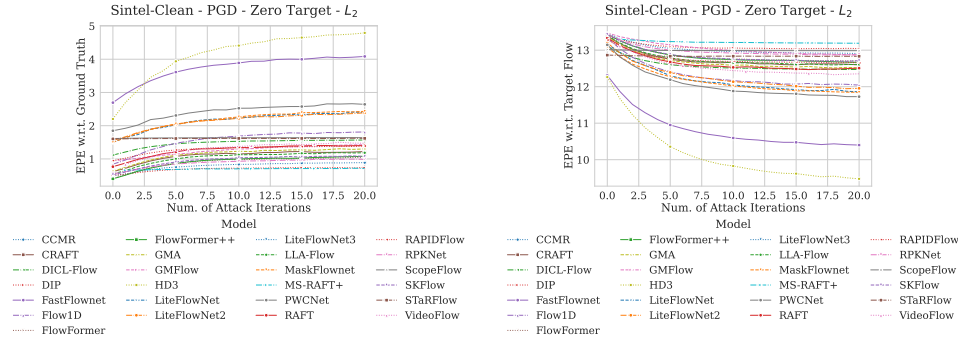


Figure 47: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

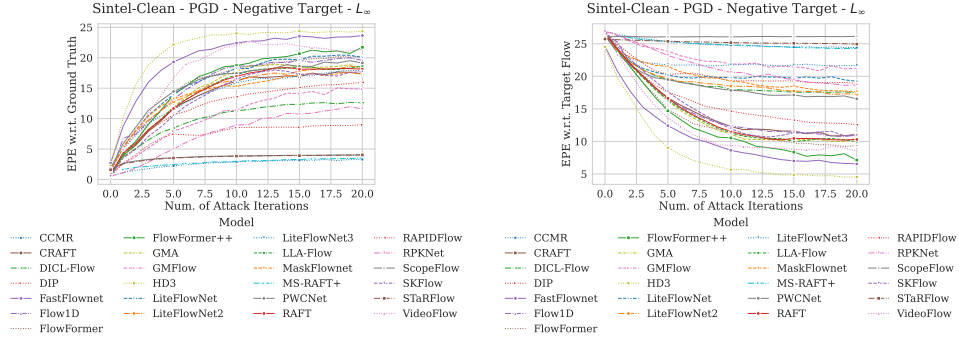


Figure 48: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

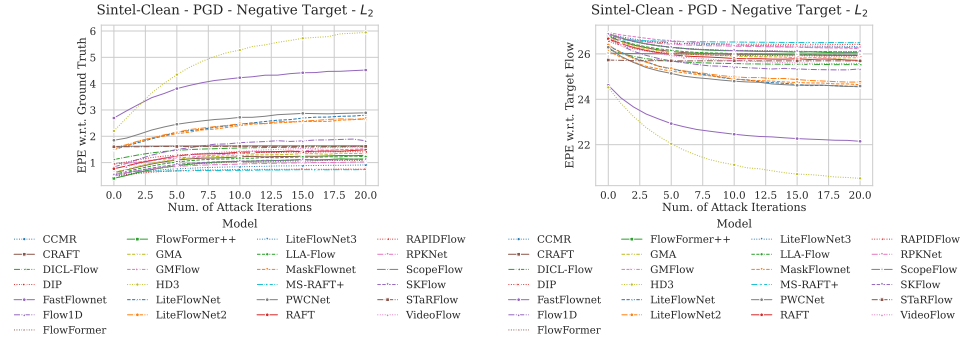


Figure 49: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

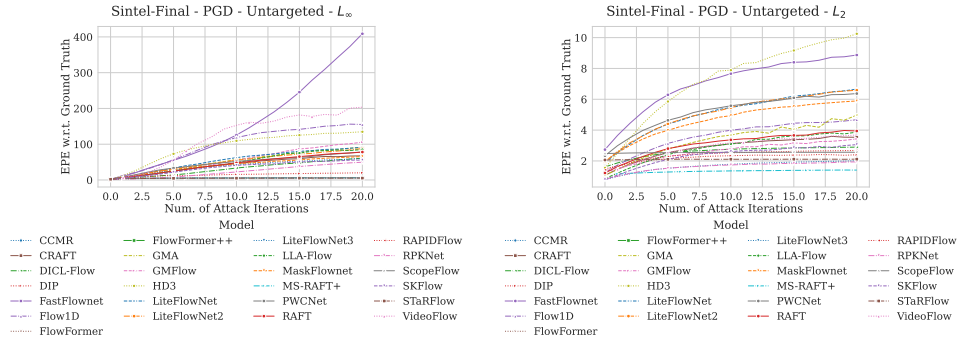


Figure 50: Evaluations for non-targeted PGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.



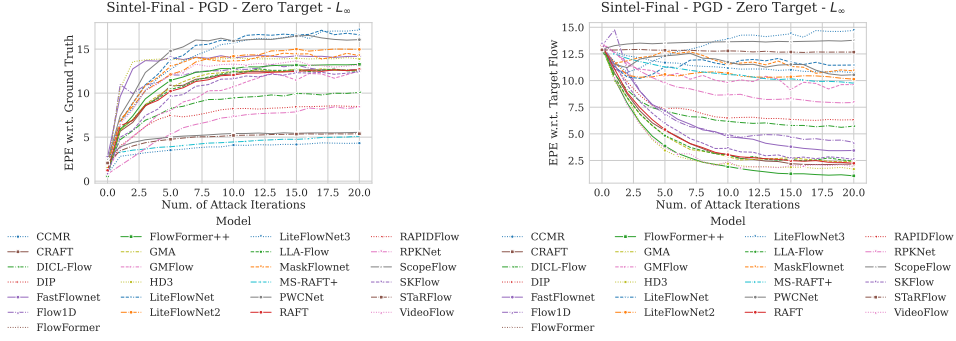


Figure 51: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

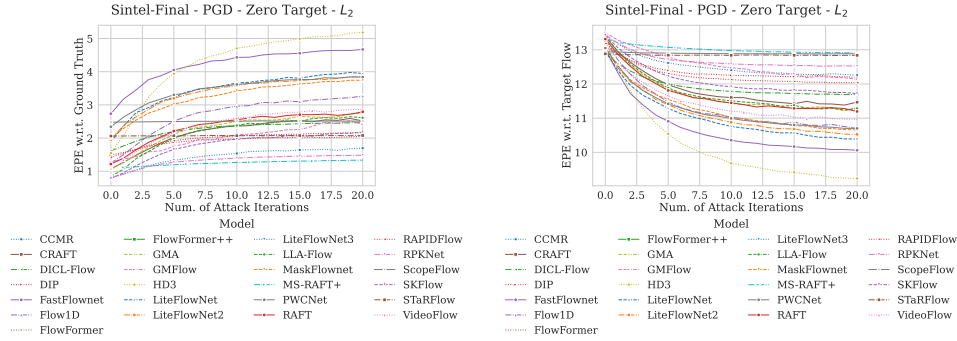


Figure 52: Evaluations for targeted PGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

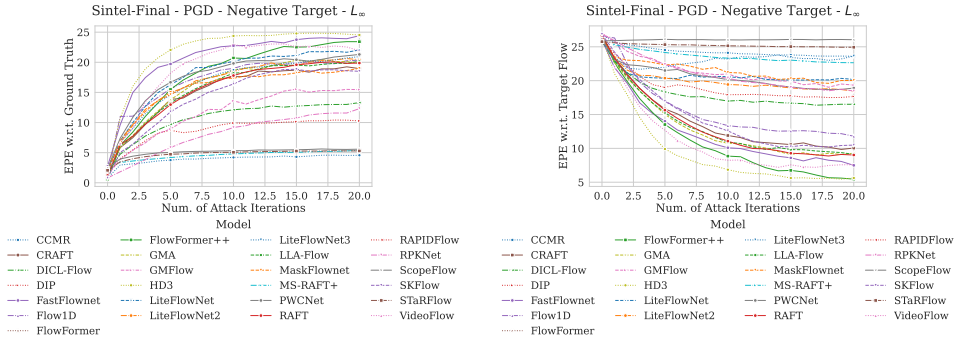


Figure 53: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.



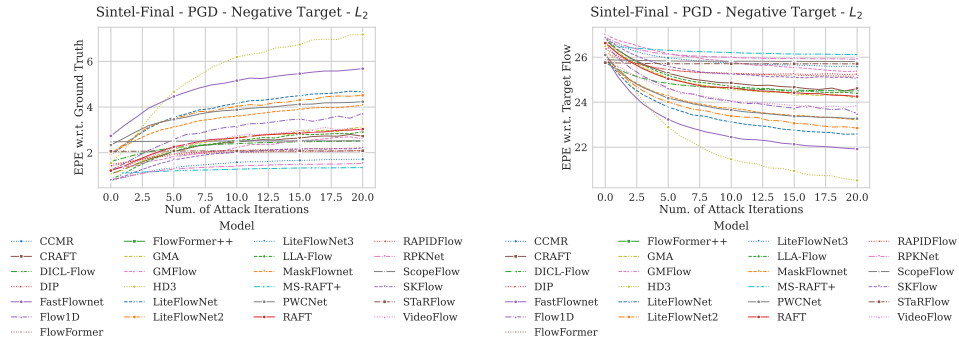


Figure 54: Evaluations for targeted PGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

### G.1.4 CosPGD ATTACK

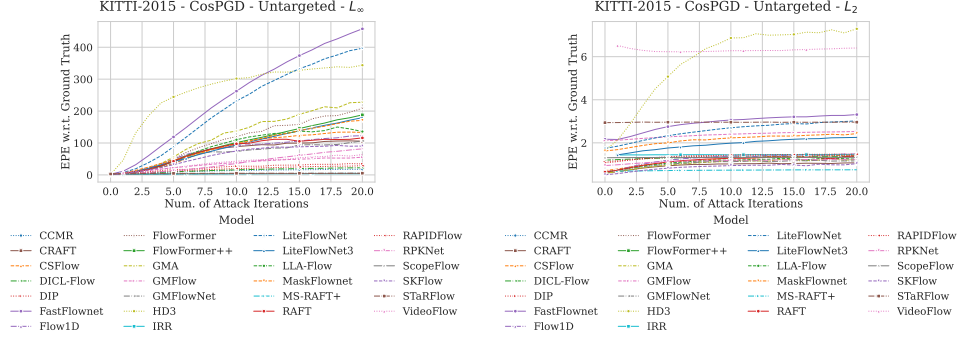


Figure 55: Evaluations for non-targeted CosPGD attack under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

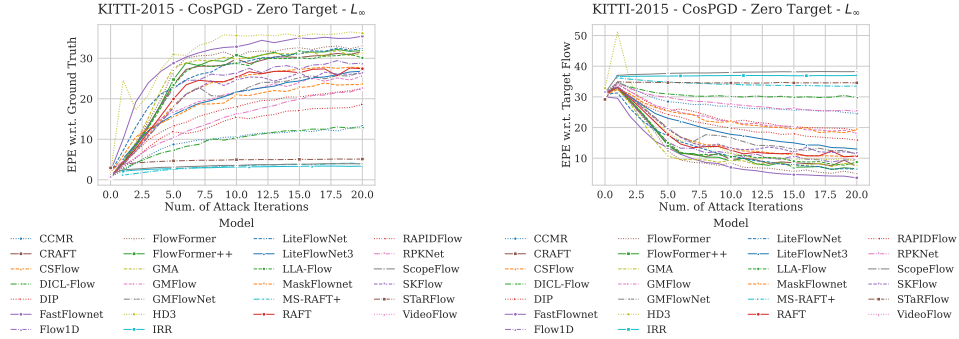


Figure 56: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

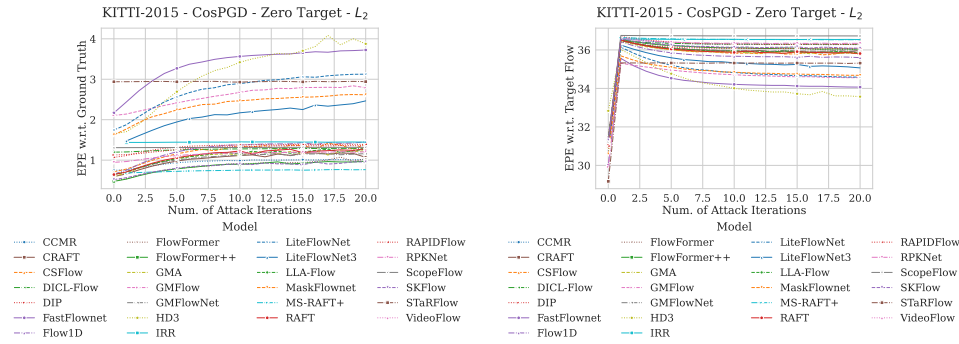


Figure 57: Evaluations for targeted CosPGD attack with target  $\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

Here we report the evaluations using CosPGD 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. 55, Fig. 56, Fig. 57, Fig. 58, Fig. 59, Fig. 60, Fig. 61, Fig. 62, Fig. 63, Fig. 64, Fig. 65, Fig. 66, Fig. 67, Fig. 68, and Fig. 69.

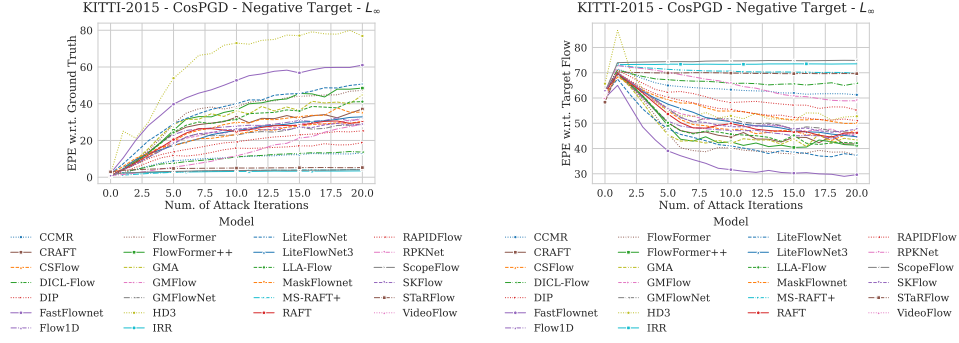


Figure 58: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

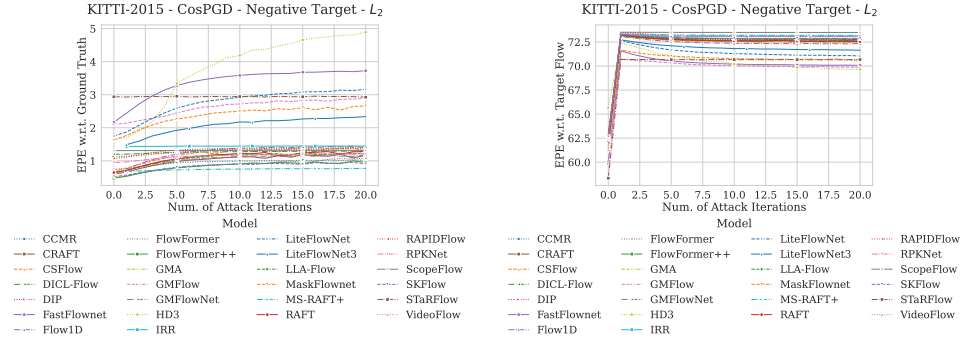


Figure 59: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the KITTI2015 dataset. The attack was optimized w.r.t. the ground truth predictions.

#### G.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.

#### G.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.

### G.2 COMMON CORRUPTIONS OVERVIEW

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

#### G.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.

#### G.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.

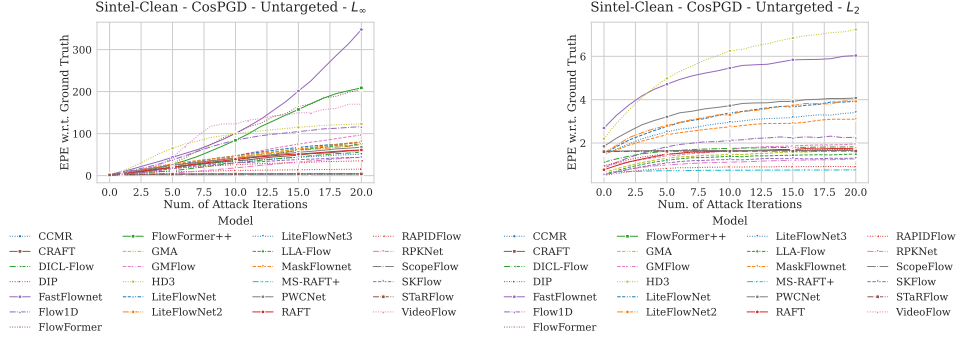


Figure 60: Evaluations for non-targeted CosPGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

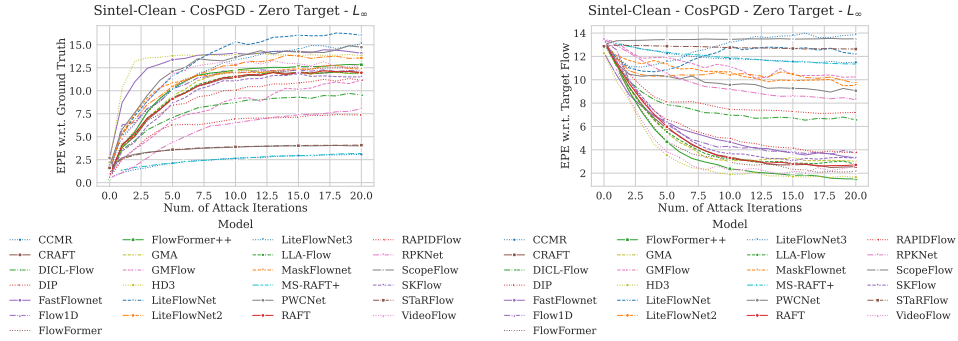


Figure 61: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

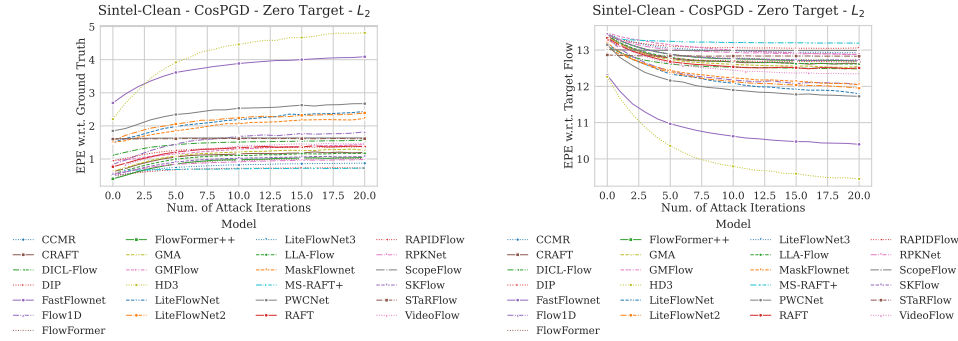


Figure 62: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

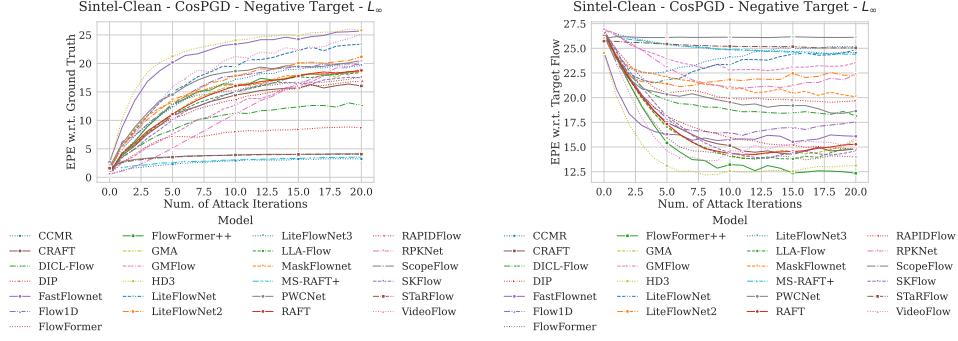


Figure 63: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

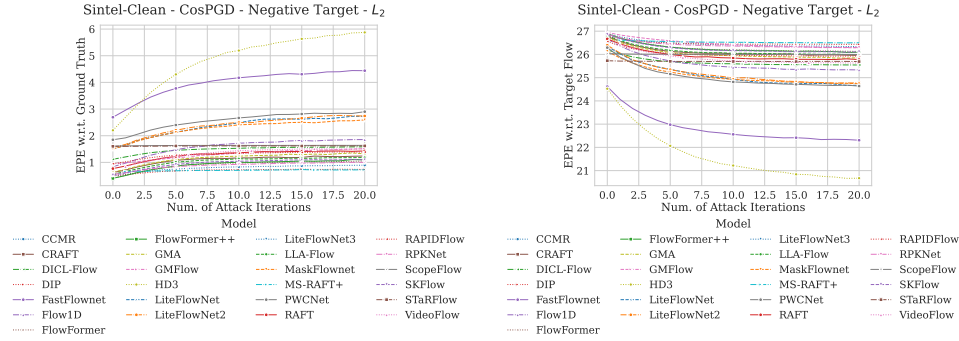


Figure 64: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (clean) dataset. The attack was optimized w.r.t. the ground truth predictions.

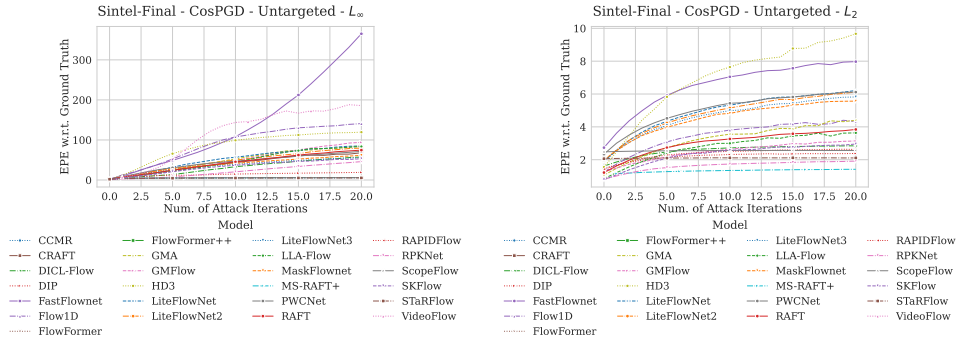


Figure 65: Evaluations for non-targeted CosPGD attack under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

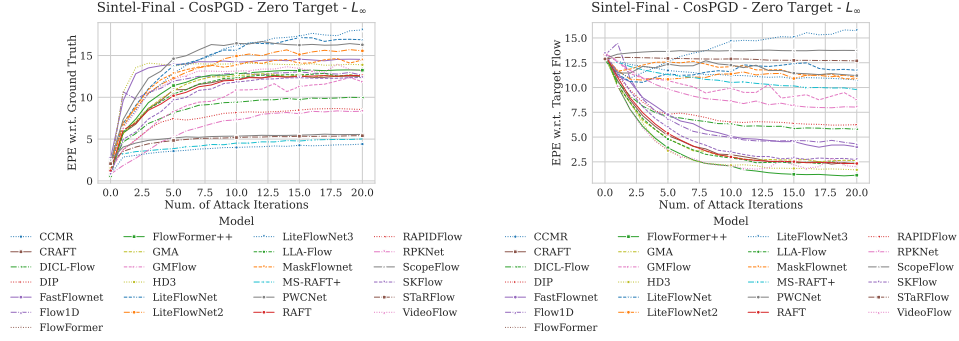


Figure 66: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

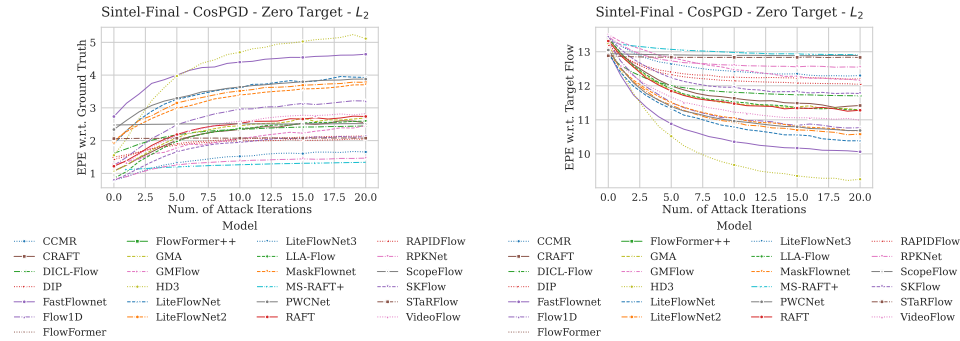


Figure 67: Evaluations for targeted CosPGD attack with target  $\vec{0}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

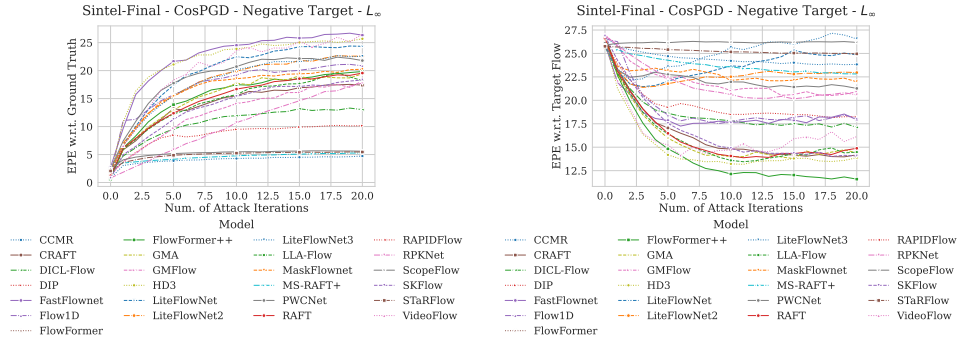


Figure 68: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_\infty$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.



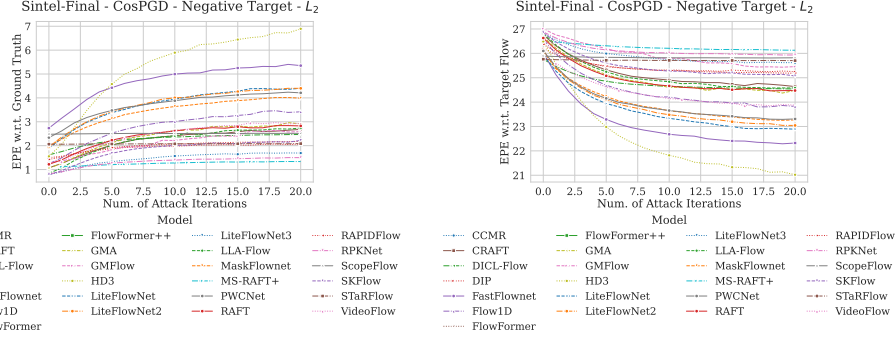


Figure 69: Evaluations for targeted CosPGD attack with target  $-\vec{f}$  under  $\ell_2$ -norm bound using the MPI Sintel (final) dataset. The attack was optimized w.r.t. the ground truth predictions.

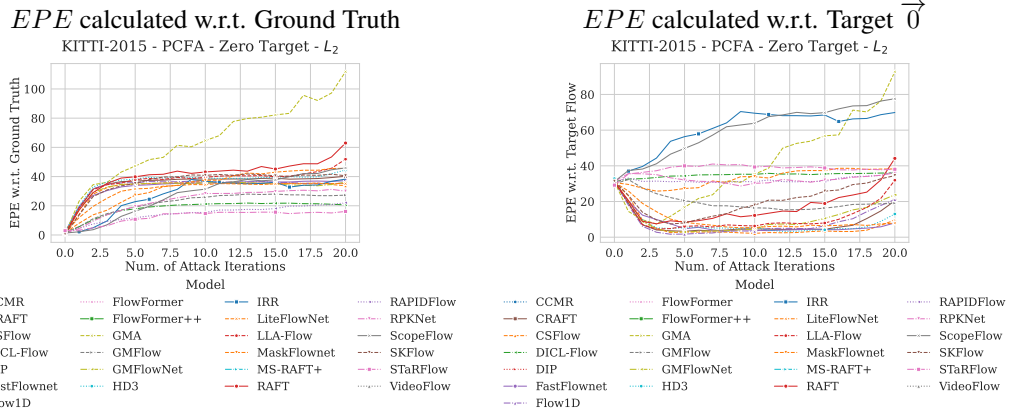


Figure 70: Evaluating all optical flow estimation methods against PCFA attack with target  $\vec{0}$  over multiple attack iterations.

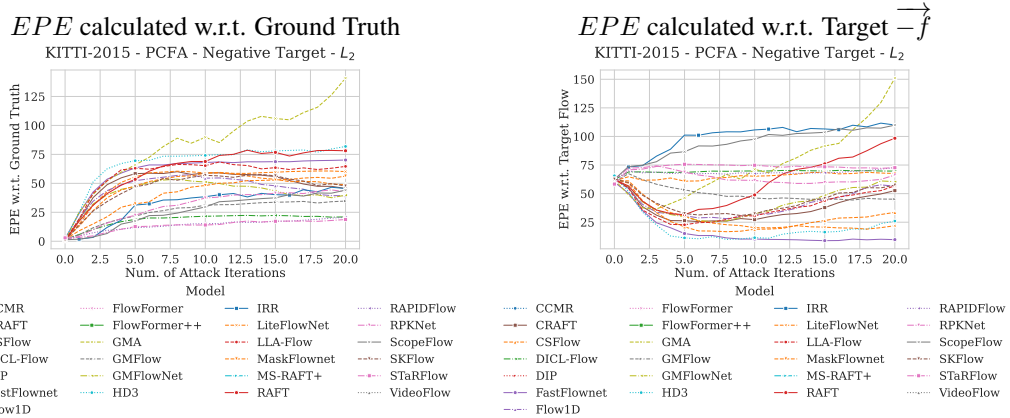


Figure 71: Evaluating all optical flow estimation methods against PCFA attack with target  $-\vec{f}$  over multiple attack iterations.



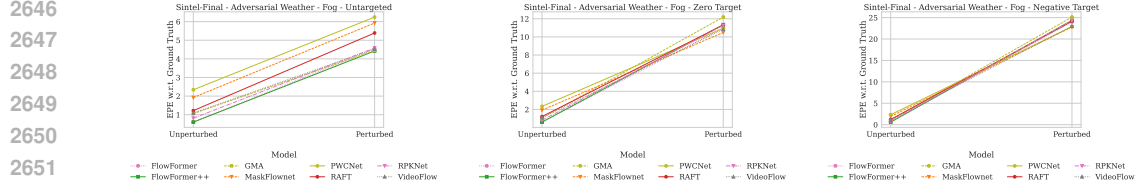


Figure 72: Evaluations for Adversarial Weather attack with Fog optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

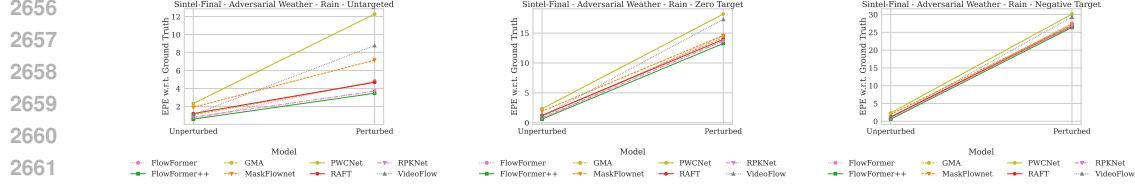


Figure 73: Evaluations for Adversarial Weather attack with Rain optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

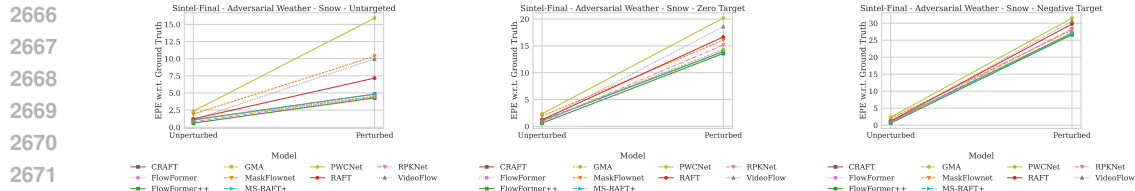


Figure 74: Evaluations for Adversarial Weather attack with Snow optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

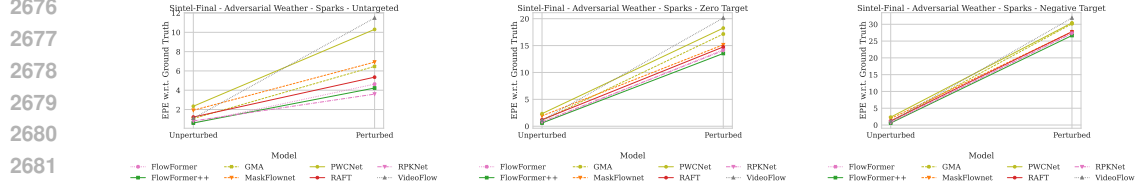


Figure 75: Evaluations for Adversarial Weather attack with Sparks optimized as a non-targeted attack (left), and targeted attack with targets  $\vec{0}$  (center) and  $-\vec{f}$  (right).

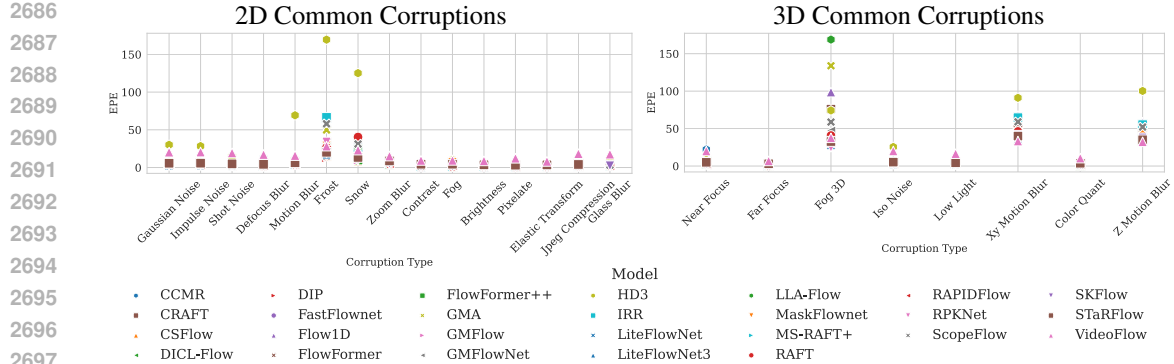


Figure 76: Performance of various optical flow estimation methods after corruptions on the KITTI2015 dataset.



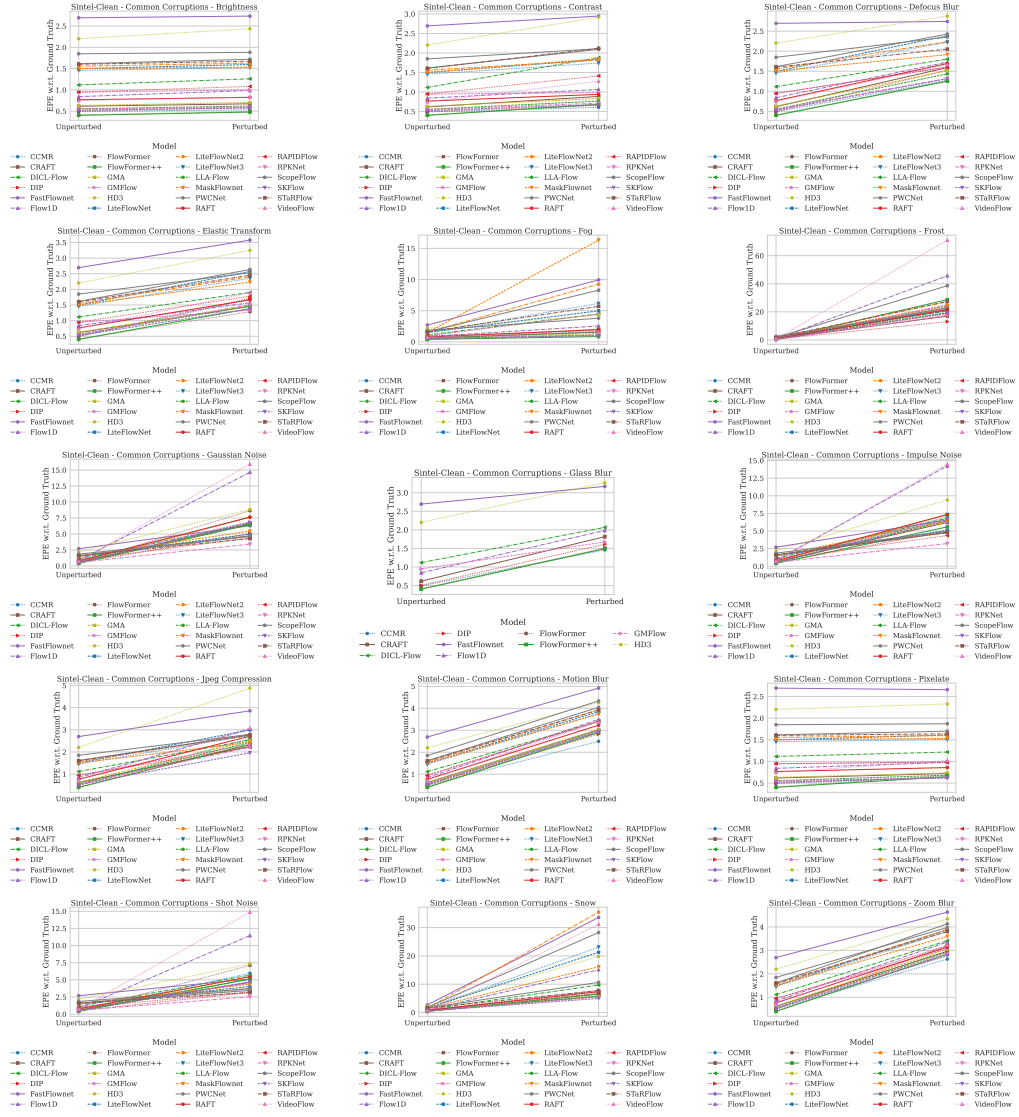


Figure 78: Evaluating optical flow estimation methods against all 2D Common Corruptions on the MPI Sintel (clean) dataset.

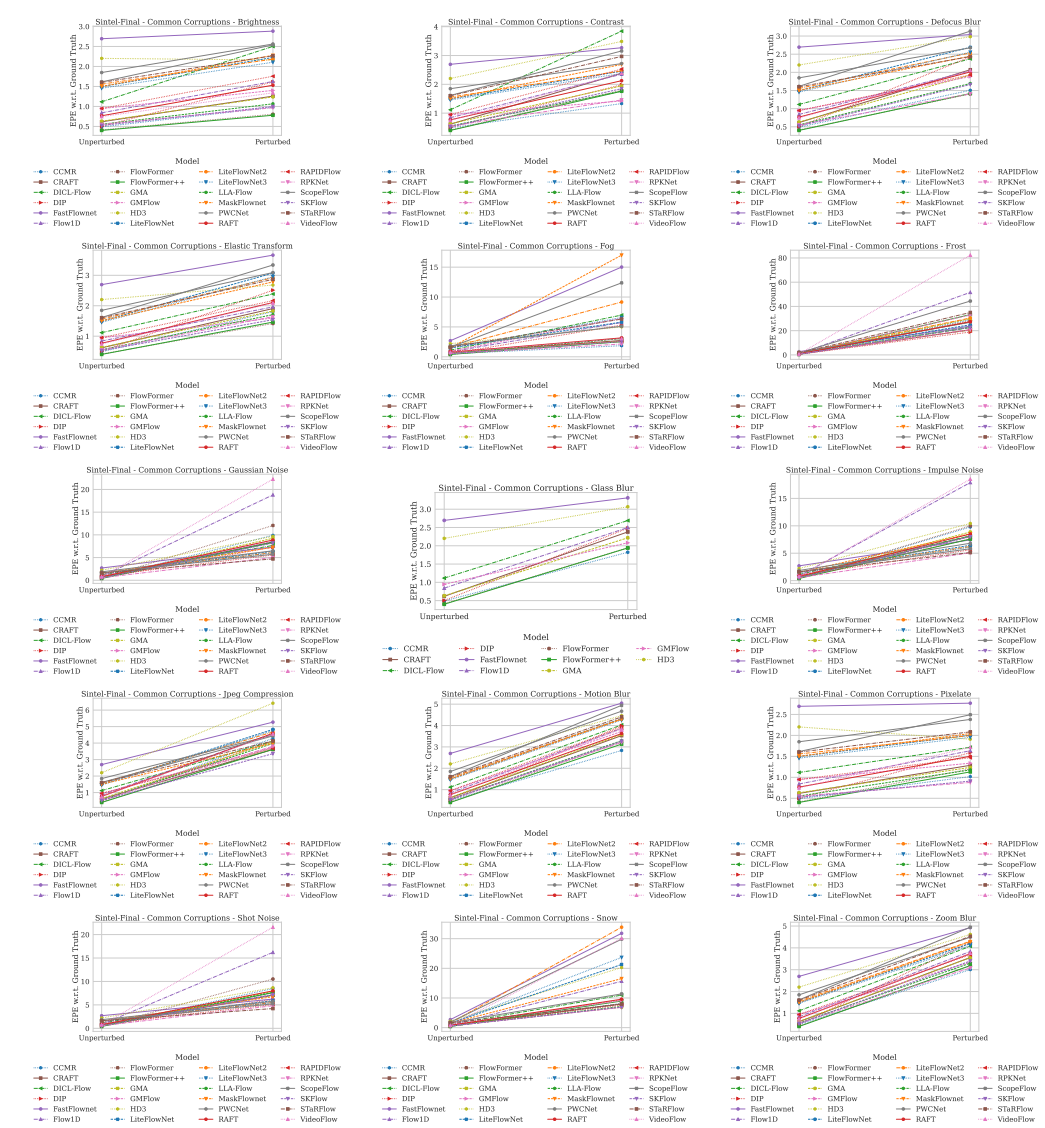


Figure 79: Evaluating optical flow estimation methods against all 2D Common Corruptions on the MPI Sintel (final) dataset.

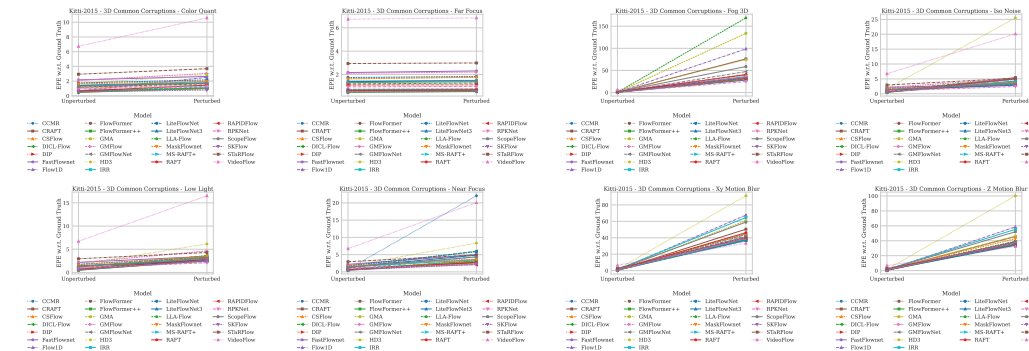


Figure 80: Evaluating optical flow estimation methods against the considered 3D Common Corruptions on the KITTI2015 dataset.

