

A Flow-based Credibility Metric for Safety-critical Pedestrian Detection

Maria Lyssenko^{1,2}, Christoph Gladisch¹, Christian Heinzemann¹, Matthias Woehrle¹, and Rudolph Triebel^{3,4}

¹ Robert Bosch GmbH, Corporate Research, Germany
firstname.lastname@de.bosch.com

² Technical University of Munich, Germany maria.lyssenko@tum.de

³ German Aerospace Center, Germany rudolph.triebel@dlr.de

⁴ Karlsruhe Institute of Technology, Germany
rudolph.triebel@kit.edu

Abstract. Safety is of utmost importance for perception in automated driving (AD). However, a prime safety concern in state-of-the-art object detection is that standard evaluation schemes utilize safety-agnostic metrics to argue sufficient detection performance. Hence, it is imperative to leverage supplementary domain knowledge to accentuate safety-critical misdetections during evaluation tasks. To tackle the underspecification, this paper introduces a novel credibility metric, called *c-flow*, for pedestrian bounding boxes. To this end, *c-flow* relies on a complementary optical flow signal from image sequences and enhances the analyses of safety-critical misdetections without requiring additional labels. We implement and evaluate *c-flow* with a state-of-the-art pedestrian detector on a large AD dataset. Our analysis demonstrates that *c-flow* allows developers to identify safety-critical misdetections.

Keywords: Safe Perception in AD · Optical Flow · Verification & Validation (V&V)

1 Introduction

In automated driving (AD), safety is a crucial aspect of perception systems. Driven by the remarkable performance demonstrated in perception tasks such as camera-based object detection, the demand for systems utilizing deep neural networks (DNN) has surged in the field of AD. In respect thereof, ensuring accurate and reliable perception of vulnerable road users (VRU) like pedestrians is a significant requirement.

Hence, for a systematic safety argumentation the consideration of safety concerns is of utmost importance [1]. Among other things, the concerns encompass the ubiquitous safety-agnostic evaluation that treats all misdetections equally irrespective of their particular relevance for the safe driving task. To tackle safety-agnosticism, recent works by Wolf *et al.* [20], Ceccarelli *et al.* [4], and ourselves [12,13] leverage domain knowledge (in form of, *e.g.* temporal distance) to incorporate a notion of safety into the evaluation of object detectors. This safety-awareness prioritizes pedestrians that are especially relevant to the autonomous vehicle (AV) for further downstream tasks like planning and control.

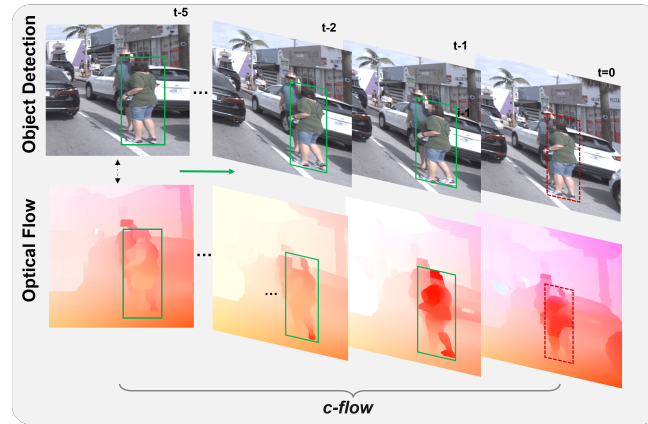


Fig. 1: We present a novel metric, called *c-flow*, to quantify the credibility of pedestrian bounding boxes. Therefore, we exploit temporal information from optical flow (bottom) and a series of pedestrian detections ($\mathcal{B}_{pred,t}$) to rate whether it is credible to have a **misdetection** at $t=0$, *i.e.* a false negative in our case. Thus, *c-flow* provides an supplementary signal that helps to uncover prospectively challenging, safety-critical **misdetctions** ($c-flow \rightarrow 0$). In the optical flow maps, the intense red color illustrates particularly high, relative motion that we leverage for the metric design.

However, in the safety-critical domain of AD, failure cases where the system fails to detect objects (so-called *false negatives*) or where the model produces spurious detections (so-called *ghost objects* or *false positives*) during operation continue to occur [2]. Hence, both of these erroneous detections not only degrade the overall performance but, in the worst case, lead to a hazardous situation. As an example, failing to detect a relevant pedestrian without a warning from the system could result in a safety-critical event such as a collision (as depicted in Fig. 1). As of now, the development-time evaluation in our previous work [13] showed that a state-of-the-art detector tends to produce several misdetections in safety-critical sequences. This, in particular, could potentially lead to a loss of the pedestrian track at runtime without notice. For mitigating such cases, our work is motivated by the following research question:

How can we identify such detection errors in an unsupervised manner and determine whether detections are indeed credible?

Assessing credibility of detections in an unsupervised manner (*e.g.* in an AV) implies that we can no longer leverage information from ground truth annotations. Instead, we need another, supplementary source of information that can complement the input signal. Contemporary works advocate for runtime monitoring of an auxiliary signal to detect potential problems as early as possible, so the AV can adapt his behavior accordingly [7,16]. A concrete implementation is proposed by Geissler *et al.* [6] where the authors present an approach to monitor input faults of an object detector leveraging knowledge from hidden layers and hardware memory. Yet, we want to avoid a tight coupling to a specific hardware, and propose a method to verify detections from any object detector. Hence, in line with works that use optical flow to assess the consistency of detections [14] and segmentation masks [19], we introduce a flow-based method to identify potential errors of our object

detector. As optical flow measures the perceived relative motion of objects between two subsequent images from a sequence, we thereby exploit the time- and space-consistency of objects in the real world and embed it as an auxiliary signal in our proposed method.

An exemplary sequence consisting of RGB images (top) and optical flow maps (bottom) for a safety-critical pedestrian is shown Fig. 1. Here, the flow maps illustrate the created flow within the predicted bounding box (\mathfrak{B}_{pred}) due to the pedestrians’ motion relative to the AV (intense red color). Since the quality of optical flow calculation is better in short ranges, as more pixels exist for objects, using the flow to assess the credibility of detections is also a natural choice for our use case where safety-critical pedestrians occur in short ranges from the AV.

In this paper we present a novel metric, called *c-flow*, for quantifying the credibility of 2D bounding box detections by using information from optical flow. In particular, we (i) apply *c-flow* to a predicted bounding box \mathfrak{B}_{pred} and determine whether \mathfrak{B}_{pred} is indeed a credible true positive (TP) detection. (ii), to detect cases of false negative (FN) candidates after a track of successful pedestrian detections, we propose a technique to infer a *hypothesized* bounding box \mathfrak{B}_{hyp} from past predictions. This allows us to apply *c-flow* to \mathfrak{B}_{hyp} for assessing the credibility of false negative detections.

We provide an thorough experimental evaluation of *c-flow* based on the Argoverse 1.1 dataset [5] and a state-of-the-art RetinaNet [11] for object detection. Accordingly, we (i) demonstrate the validity of our novel metric in a supervised manner where ground truth information is accessible. Our results show that *c-flow* can successfully discriminate between true positives and safety-critical false negatives for most pedestrian samples in vicinity of the AV. For ambiguous cases of true positives by means of *c-flow*, a qualitative analysis shows how the metric is eligible to identify cases with an unreasonable labeling policy (e.g. distant or heavily occluded pedestrians), i.e. aiding as a tool to provide dataset insights. (ii), we evaluate the metric in scope of prospective runtime applications. Hence, as ground truth information is unavailable, we demonstrate how the approach with \mathfrak{B}_{hyp} is effectively applicable to an unsupervised setting.

2 Methodology

In this section, we present our novel credibility metric. *c-flow* leverages changes in optical flow to quantify the credibility of bounding boxes within safety-critical sequences by means of the temporal distance (time-to-collision) between the AV and the respective pedestrian. We motivate our methodology for the *c-flow* metric design in Sec. 2.1, followed by implementation details in Sec. 2.2. In Sec. 2.3, we illustrate the approach to compute so-called *hypothesized* bounding boxes \mathfrak{B}_{hyp} for unsupervised use cases.

2.1 Motivation on *c-flow*

Let us provide a motivating example for using flow to assess the credibility of pedestrian detections. To this end, we select one safety-critical pedestrian track as shown in Fig. 2 (for details see Sec. 3.2 and Fig. 5), with time indicated by the colormap. The respective criticality as time-to-collision (TTC) is on the x-axis. As this is a real pedestrian track with

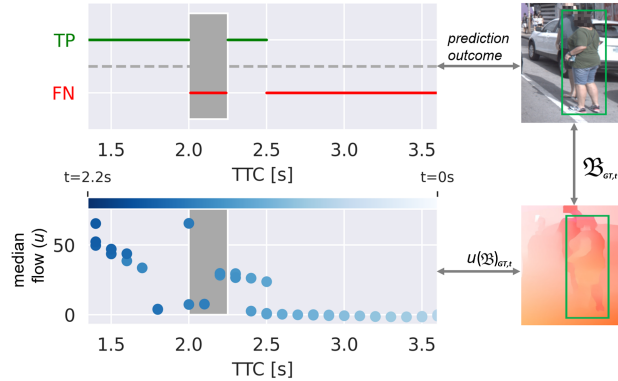


Fig. 2: Selected pedestrian track (red star sequence from Fig. 5) showing evolution over time (in color) for a pedestrian bounding $\mathfrak{B}_{GT,t}$, illustrating (i) the classification outcome (top) and (ii) the median score u gathered from the optical flow map at $\mathfrak{B}_{GT,t}$. The gray window highlights sudden fluctuations in u that correlate with the switch between classification outcomes at, *e.g.*, $TTC \approx 2.25s$.

associated ground truth (GT), the classification outcome that is shown in the top plot for a bounding box $\mathfrak{B}_{GT,t}$ is at different times t either a TP or a FN (*i.e.* not detected).

The upper plot depicts the typical behavior of a longer track: No detections for the pedestrian at larger distances (*i.e.* high TTC) with the first successful detection at $TTC = 2.5s$ after the AV has approached the pedestrian⁵. In the lower plot, we use a pedestrian bounding box $\mathfrak{B}_{GT,t}$ from our selected track at time t in $t \in [1.25s, 3.5s]$ to construct the corresponding data points. Therefore, we apply $\mathfrak{B}_{GT,t}$ to the optical flow map to define the area of interest where we calculate the optical flows' median score u ⁶.

On the one hand, in the temporal plot around the highlighted gray window, we denote significant changes in u , which are a strong indicator for errors concerning the consistency of objects between images [19]. On the other hand, if we look at the detection type, we can derive that the rapid changes in u correlate with the classification outcome switch from TP to FN at $TTC \approx 2.25$. As optical flow estimates the relative motion of the pedestrian between consecutive images, sudden changes in u may refer to abrupt variations in the perceived motion. This may indicate unexpected changes in the scene such as partial or complete vanishing of the pedestrian due to, *e.g.* a sudden occlusion. This conveys a difficult case for our detector-under-test as occlusions are a well-known challenge for object detectors.

In summary, this example shows that optical flow can provide a supplementary signal for analyzing the credibility of predictions: Small changes in u indicate a continuous existence of a pedestrian detection over images, whereas sudden fluctuations in u may help to identify generally difficult cases for the object detector.

⁵ Please note, since near range pedestrians have a better quality of optical flow due to a higher pixel count, the optimal use of optical flow is aligned with our motivation to perform a flow-based credibility assessment on \mathfrak{B}_{GT} for pedestrians close to the AV.

⁶ For further discussion, we only utilize the horizontal flow map, as the vertical flow showed only oscillating displacements between consecutive images, *i.e.* no meaningful vertical shift in \mathfrak{B}_{GT} .

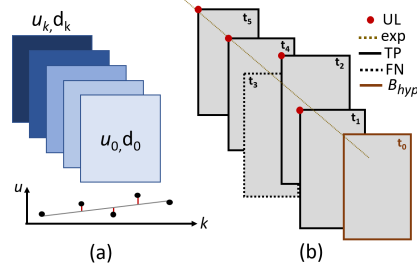


Fig. 3: (a) Leveraging linear regression to determine the variability in u to construct c -flow. (b) For missing detections at t_0 , we apply the methodology of hypothesized bounding boxes (\mathfrak{B}_{hyp}) to infer the required \mathfrak{B}_{hyp} . We extract the upper left corner (UL) of past predictions to extrapolate the hypothesized UL position at t_0 .

Consequently, we want to discriminate correct TP from prospective failure cases where it is non-credible to have a missing detection. This motivates our work to exploit the sudden changes in optical flow to design a metric, which we call c -flow in the following, that can be applied to a bounding box to assess the credibility. With the novel metric we aim to achieve the discrimination between TP detections (high c -flow score) and non-detections, *i.e.* FN (low c -flow score).

2.2 Metric Design: c -flow

Based on this underlying idea, we design the c -flow metric for individual pedestrians and the corresponding bounding box. Therefore, our metric design is guided by the following desiderata: Since we want to assess a credibility score and thus, to discriminate successful detections from cases of FN candidates, we design our novel metric to be c -flow $\rightarrow 0$ for FN and c -flow $\rightarrow 1$ for TP, respectively.

In the first step, we quantify the continuity of pedestrian detections among consecutive images. Therefore, we utilize the given track information from Argoverse 1.1 (see Sec. 3.1), *i.e.* we leverage \mathfrak{B}_{GT} from GT for missing detections (FN). Concretely, we consider a pedestrian track of a certain length with pairs of (i) the pedestrians' bounding box \mathfrak{B}_{GT} and (ii) calculated optical flow median scores u (extracted from flow maps within \mathfrak{B}_{GT}). To include the sudden flow changes in u leveraging a series of successful detections across time, we construct our metric on the basis of a time window w , that encompasses the past k images with $w = [t_k, t_{k-1}, \dots, t_1, t_0]$, *i.e.* $\mathfrak{B}_{GT,t}$ for each image at t in w .

In the second step, we want to quantify the flow change across w up to the current image at t_0 . Therefore, we leverage the calculated u for each $\mathfrak{B}_{GT,t}$, to determine the variability of the flow in $[u_k, u_{k-1}, \dots, u_1, u_0]$ over w . Given the underlying assumption of linear progression of the movement across w , *i.e.* a linear change of flow in w , we now compute a linear regression on $[u_k, u_{k-1}, \dots, u_1, u_0]$ as illustrated in Fig. 3 (a). After applying the linear regression, we estimate the residuals $r = [r_k, r_{k-1}, \dots, r_1, r_0]$ (error bars in red) between regressed u scores and $[u_k, u_{k-1}, \dots, u_1, u_0]$ (black dots), respectively, to quantify the

deviation from the linear progression, *i.e.* a sudden flow change that indicates a difficult case. To define the overall error ε over w , we calculate $\varepsilon = \sum_{i=0}^k r_i$.

In the previous steps, we have identified a strong correlation between rapid changes in u and switches in the classification outcome. However, as flow changes over w also correlate with the change in size of $\mathfrak{B}_{GT,t}$ (*i.e.* a larger bounding box captures more pixels to estimate u), we must consider this change in box size. Therefore, we inject supplementary information on $\mathfrak{B}_{GT,k}$ size into the metric design. To describe the change in size, we utilize the diagonals between the two most recent bounding boxes d_0 and d_1 (at t_0 and t_1). Thus, we introduce $\Delta d_{t_0,1} = d_0 - d_1$ to measure the difference.

Based on the desiderata, $\Delta d_{t_0,1}$, and ε , we can now show the complete formula for the credibility metric in Eq. 1. Please note, as we bring $\Delta d_{t_0,1}$ and ε in relation to each other, we first normalize both parameters. Finally, we propose the sigmoid function

$$c\text{-flow}(\Delta d_{t_0,1}, \varepsilon) = \text{sigmoid}\left(\frac{\Delta d_{t_0,1}}{\varepsilon}\right) \quad (1)$$

that introduces a non-linear, gradual transition from 0 to 1, normalizing *c-flow* to the unit interval.

2.3 Handling the Absence of Ground Truth in Unsupervised Use Cases

Up to now, we have used GT bounding boxes for analyzing the FN cases from Fig. 2. Thus, we have applied the *c-flow* metric using a window of GT bounding boxes, to assess the credibility that a predicted bounding box \mathfrak{B}_{pred} is missing at t_0 . However, let us assume that we do not have an available \mathfrak{B}_{pred} for the current image at t_0 , as, *e.g.*, when the detector is running in a vehicle or when performing an unsupervised data analysis. Thus, for the current track-under-test for a particular pedestrian, we want to assess whether the lack of \mathfrak{B}_{pred} at t_0 is credible or whether it is a prospective FN. Therefore, to compensate the absence of \mathfrak{B}_{pred} in the *c-flow* computation, we introduce the concept of *hypothesized* bounding boxes, *i.e.* \mathfrak{B}_{hyp} , that are extrapolated from available, past \mathfrak{B}_{pred} for extracting the optical flow crops. In Fig. 3 (b), we can see a sequence of past, predicted \mathfrak{B}_{pred} at t_5, t_4, t_2 , and t_1 .

As a first step, we extract the upper left corner (UL) of previous TP detections. We focus on UL as it is the most probable, visible corner for cases with horizontal or vertical possible occlusions in the front camera perspective *e.g.*, a pedestrian standing behind a car on the right side of the street. Next, we perform an extrapolation (see Fig.3 (b)) of the pixel coordinates to hypothesize UL for t_0 . Please note, we consider the previous five ($k = 5$) frames using the detection window w from Sec. 2.2. However, the algorithm does not require all \mathfrak{B}_{pred} in w to exist (see skipped detection at t_3). It necessitates only a minimum of two existing \mathfrak{B}_{pred} in the detection window (num. TP ≥ 2) to fit the extrapolation line *exp*. Thus, this approach would also work after the first *two* successful detections after the pedestrian track was initially opened.

To hypothesize UL at t_0 , we need to estimate the positional shift from the last available \mathfrak{B}_{pred} (t_1 in our case). Therefore, to determine the mean positional shift between n consecutive detections, we calculate the mean pixel displacement, *i.e.* distance between the first (t_5) and the last available \mathfrak{B}_{pred} (t_1) from our example averaged over $n = 4$. Finally, we apply the determined shift along *exp* to hypothesize UL at t_0 . For the initial proof-of-concept,

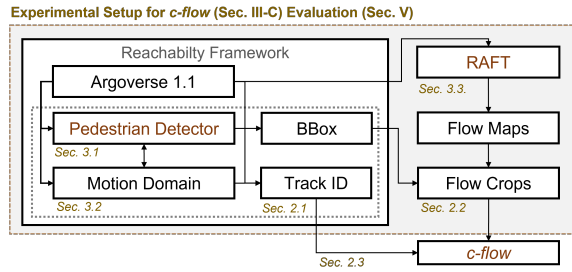


Fig. 4: Experimental setup for *c-flow* evaluation. We extend our reachability framework (RF) from [13] by a DNN-based optical flow estimation using RAFT. We perform our evaluation on the basis of identified, potentially safety-critical pedestrian tracks (Track ID) with respect to TTC derived from the motion domain.

we assume a small change in dimension of \mathfrak{B}_{hyp} to \mathfrak{B}_{pred} from the previous image at t_1 and set the width and height of \mathfrak{B}_{hyp} to be w_{t_1} and h_{t_1} , respectively. As a last step, we extract the estimated crop to calculate $c-flow(t_0)$. Please note, in Fig. 3 (b), t_3 could have been replaced by \mathfrak{B}_{hyp} to boost accuracy for runtime applications. Further approaches to regress the dimensions of \mathfrak{B}_{hyp} from w we leave for further work.

3 Experimental Setup

In the following, we leverage the setup from Fig. 4 to evaluate our novel *c-flow* metric. In our setup we utilize (i) the nuImages dataset to pre-train our RetinaNet from Sec. 3.1, (ii) an adapted sequence selection approach using our reachability framework (RF) [13] in Sec. 3.2, and (iii) the incorporation of RAFT from Sec. 3.3 to generate optical flow maps.

3.1 Datasets and the Pedestrian Detector

For the training of our pedestrian detector we use nuImages [3]⁷. Thereby, we do not train on the Argoverse 1.1 dataset and thus, we can employ all splits from Argoverse 1.1 [5]⁸ with Argoverse HD 2D annotations [10] to evaluate *c-flow*. Please note, in our evaluation, we only use images from the `ring_front_center` camera. Afterwards, we estimate the matching LiDAR GT (different frequency at 10 Hz) using `image_list_sync` to extract physical properties of the AV and the pedestrians, *i.e.* information that is required to calculate the TTC. We implement the RetinaNet⁹ using PyTorch and employ the following training protocol: We utilize the ResNet-50 backbone, the Adam optimizer and the applied learning rate of $1e^{-5}$, the `reduceOnPlateau` scheduler (patience=2) from the `optim` module, and we train RetinaNet for 200 epochs. The trained model achieved a reasonable performance of 0.31 *AP50* for the pedestrian class on the respective nuImages

⁷ <https://github.com/nutonomy/nuscenes-devkit>

⁸ <https://github.com/argoverse/argoverse-api>

⁹ <https://github.com/yhenon/pytorch-retinanet>

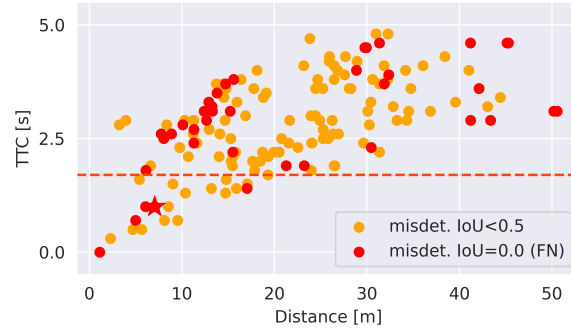


Fig. 5: Extracted pedestrian tracks from Argoverse 1.1 using the reachability framework. Each data point represents an interaction defined by the most critical pedestrian misdetection *w.r.t.* TTC and distance. The dotted line separates critical ($TTC < 2s$) vs. non-critical interactions. Misdetections with a poor detection quality ($0 < IoU < 0.5$) are highlighted in orange and FN ($IoU = 0$) are marked in red.

validation split. For the pedestrian class in Argoverse 1.1, the model reached a performance of 0.35 AP_{50} and 0.34 AP_{50} on the train and validation split, respectively. Please note that we do not strive for highest performance as we are interested in a variety of cases for our *c-flow* metric.

3.2 Pedestrian Track Selection on Argoverse 1.1

As we want to evaluate *c-flow* on individual pedestrian tracks of different criticality, we first need to extract the respective sequences *w.r.t.* TTC. Therefore, we leverage our RF as described in previous work [13] to calculate reachable sets and its intersection (between AV and pedestrians) utilizing physical properties as, *e.g.*, velocity and position and semantic map information with explicit lanelets.

Fig. 5 shows the total sequence count that encompasses 142 of such tracks with 17 highlighted interactions that include critical scenarios, *i.e.* completely missed pedestrians (red) or insufficient performance with $0 < IoU < 0.5$ (orange). Note that we introduce an improved approach that increases the sequence count from 32 (using [13]) to 142 available sequences. The novel implementation incorporates a matching refinement to find corresponding tracks between Argoverse 1.1 and Argoverse HD. For instance, Argoverse HD creates a new ID as soon as the tracked pedestrian reappears after a full occlusion. Thus, multiple tracks in Argoverse-HD may describe sub-parts from a single Argoverse 1.1 sequence. Our adaptation alleviates this inconsistency and summarizes multiple tracks from Argoverse-HD to one unique sequence that is matched with its counterpart from Argoverse 1.1.

3.3 RAFT: Optical Flow-based Motion Estimation

We employ the RAFT [17]¹⁰ model to determine the optical flow on the basis of two consecutive images. The model was already trained on KITTI¹¹ with a performance of 7.51 *epe* (endpoint error) and a *F1* score of 0.269. Thereafter, for each pedestrian sequence from 3.2, we infer image pairs to calculate a flow map as shown in Fig. 1 (bottom). For visualization purposes, we use `flow_uv_to_colors` [17].

4 Experimental Results

In this section, we conduct an evaluation of the *c-flow* metric. In Sec. 4.1, we utilize the safety-critical pedestrian tracks to perform a thorough validation of the metric in supervised manner where detailed GT information is available. In Sec. 4.2, we evaluate whether *c-flow* can provide accurate information for FN detections in an unsupervised setting when using \mathfrak{B}_{hyp} instead of actual GT information.

4.1 *c-flow* Evaluation: TP vs. FN

We start by computing *c-flow* values for all pedestrians in the selected pedestrian tracks of different criticality using the predicted \mathfrak{B}_{pred} for TP, and GT annotations for FN. As discussed in Sec. 2.2, the goal is to detect cases of FN candidates after a series of successful detections (TP) by means of a low *c-flow*, *i.e.* $c-flow \rightarrow 0$. At the same time, cases of TP should have a high *c-flow* $\rightarrow 1$. Ideally, depending on a set threshold ξ , *c-flow* would be able to discriminate between (i) successful detections with $c-flow > \xi$ and (ii) cases of FN candidates with $c-flow \leq \xi$.

As a first step, we evaluate the *c-flow* scores using Fig. 6. The histogram (top) shows the distribution of the *c-flow* scores conditioned on classification GT (TP or FN). In the histogram, we can see that all FN samples received a $c-flow < \xi = 0.4$ and that almost all TP have a *c-flow* score $\geq \xi = 0.5$. However, in comparison to the occurrence of FN with a *c-flow* mostly in the interval $\in [0.0, 0.3]$, we can also see a spread of TP over lower *c-flow* scores.

Let us now investigate, how the selection of a concrete threshold on *c-flow* affects the identification of FN in safety-critical pedestrian tracks *w.r.t.* TTC. The heatmap on the bottom of Fig. 6 visualizes the amount of *c-flow* scores (*x*-axis) in relation to the TTC value (*y*-axis). The figure also shows the critical zone by the TTC threshold of $2s$ (horizontally dotted line) and a *c-flow* threshold of $\xi = 0.3$ (vertically dotted line). From the heatmap, we can derive that for high criticality ($TTC \leq 2s$), the *c-flow* value is either close to zero or shows *c-flow* scores $\geq \xi = 0.5$. Only for larger TTC values, we can see that the *c-flow* scores are less discriminating.

As a second step, in Fig. 7, we study the relative distribution of identified TP and FN samples (*y*-axis) with respect to TTC intervals of different criticality (*x*-axis). Therefore, we select two concrete thresholds ξ_1 and ξ_2 to study how the selected threshold affects the trade-off between correctly identified TP and FN on the basis of *c-flow*. From the result in

¹⁰ <https://github.com/princeton-vl/RAFT>

¹¹ <https://www.cvlibs.net/datasets/kitti/>

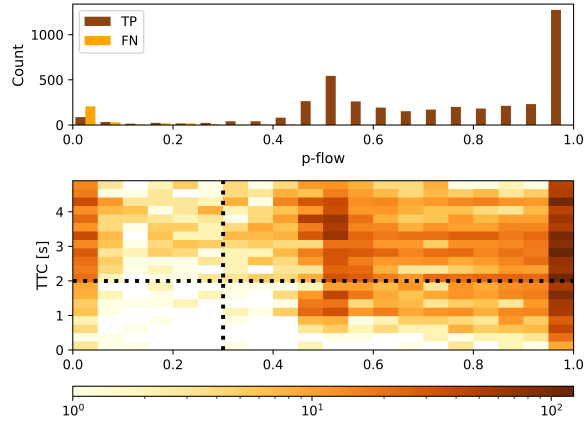


Fig. 6: Histogram (top) shows the distribution of the *c-flow* scores conditioned on the prediction types (FN or TP). The heatmap (bottom) depicts the *c-flow* scores with respect to their TTC criticality. The dotted lines divide the heatmap into regions of different criticality, constrained by the $TTC = 2s$.

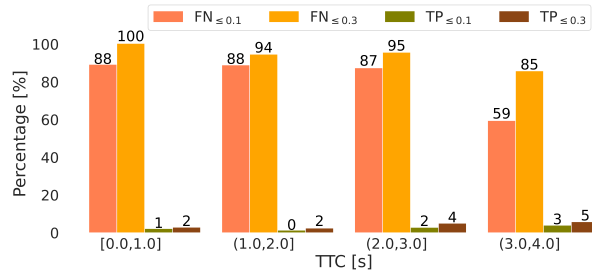


Fig. 7: Bar plot depicts the percentage of identified prediction types (TP and FN, out of annotated GT) among TTC intervals of different criticality for two *c-flow* thresholds: $\xi = 0.1$ and $\xi = 0.3$.

Fig. 6, we are particularly interested in: (i) $\xi_1 = 0.1$ as it defines the interval with almost all FN samples by means of a *c-flow* $\rightarrow 0$ but also TP occurrences and (ii) $\xi_2 = 0.3$ as for $0.1 < c\text{-flow} < 0.3$ there is yet a small number of TP and FN cases.

Now, we apply $\xi_1 = 0.1$ and $\xi_2 = 0.3$, respectively, to investigate the impact on the assigned classification outcome among specified TTC intervals. Therefore, we focus on the two most critical intervals ($TTC \leq 1s$ and $1s < TTC \leq 2s$), *i.e.* the first two sets of bars on the left, for the two most critical TTC intervals. We can see that with ξ_1 , 88% of all FN are correctly identified by means of a low *c-flow*, whereas 1% and 0% of TP samples are falsely marked as non-credible, respectively. Raising the *c-flow* threshold to ξ_2 results in 100% and 94%, *i.e.* all and almost all FN are marked as non-credible with a low *c-flow* for the two most critical TTC intervals. However, this comes at the cost of also more TP with a low *c-flow*. For safety-critical zones we can see a maximum of 2% for this threshold, which is



Fig. 8: TP samples with a low c -flow score (<0.3), to illustrate possible causes for misleading c -flow estimates. Here, we can see pairs of RGB (top) and optical flow (bottom) with predicted bounding boxes (green) and annotations (blue), showing problematic cases of TP including inaccurate predictions and no perceived flow.

still low overall, yet a four-fold increase in the critical region compared to the threshold of ξ_1 .

To investigate the underlying reason behind problematic TP with a low c -flow, we analyzed TP detections by means of c -flow <0.3 for 35 test images that we gathered from different tracks for variability. The qualitative results illustrate (i) 9 cases of occlusion, (ii) 12 cases of distant pedestrians (*i.e.*, minimally perceived flow), and (iii) 14 cases of inaccurate detections. As illustrated in Fig. 8 (left and right), multiple effects as, *e.g.* inaccurate predictions and minimal perceived flow (visualized in blue) may occur simultaneously. Thus, artifacts such as occlusion and truncation, no or minimally perceived flow for distant pedestrians, and inaccurate detections (low IoU) may produce misleading estimates due to *e.g.*, poor flow aggregation in \mathcal{B}_{pred} , *i.e.* TP with low c -flow. Thereby, we can see that low c -flow scores for TP may help to identify difficult cases for the object detector during runtime applications in an AV, especially *w.r.t.* to safety-critical pedestrians (left and mid). We leave elaborate analyses on concrete corner cases for future work.

4.2 FN in an Unsupervised Setting

Further, we evaluate c -flow’s ability to determine FN candidates in pedestrian tracks in an unsupervised setting where no GT information is available. As described in Sec. 2.3, we exploit existing \mathcal{B}_{pred} for the pedestrian from previous images for estimating a \mathcal{B}_{hyp} that we use for c -flow. In this experiment, we compare the resulting c -flow values for \mathcal{B}_{hyp} to the c -flow values obtained based on ground truth \mathcal{B}_{GT} to validate that the use of \mathcal{B}_{hyp} indeed works as intended. For differentiating the different values in the following, we use c -flow_{*h*} to denote a c -flow obtained based on \mathcal{B}_{hyp} and c -flow_{*GT*} to denote a c -flow obtained based on GT. The heatmap in Fig. 9 depicts the count of occurring FN samples conditioned on c -flow_{*GT*} and c -flow_{*hyp*}. Ideally, most of the samples should lie on the bisecting line that represents a perfect linear correlation between c -flow_{*hyp*} and c -flow_{*GT*}. In the heatmap plot, we can see a high correlation (pearson’s coefficient $\rho = 0.83$) between c -flow_{*GT*} and c -flow_{*hyp*}. Especially for c -flow $\leq \xi_1$, the bar plot illustrates that the majority of FN

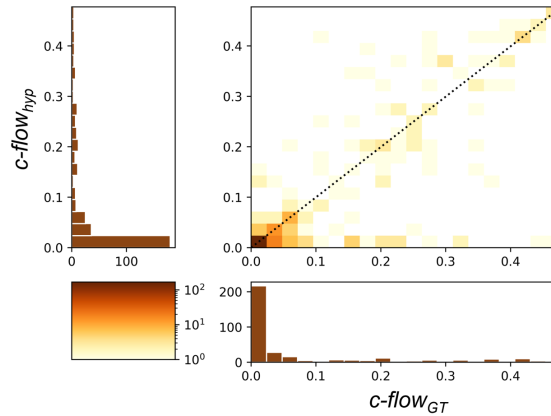


Fig. 9: The heatmap indicates the correlation ($\rho = 0.83$) of hypothesized $c\text{-flow}_{hyp}$ and annotation-based $c\text{-flow}_{GT}$, dependent on the count of FN samples. The bar plot on the left and at the bottom, show the distribution of FN samples with respect to $c\text{-flow}_{hyp}$ and $c\text{-flow}_{GT}$, respectively.

samples fall within that interval, *i.e.* we have a high FN sample count, *w.r.t.* $c\text{-flow}_{hyp}$ (in the top left) and $c\text{-flow}_{GT}$ (at the bottom). Concretely, most of the FN samples correctly received a $c\text{-flow}_{hyp} \rightarrow 0$. The plot shows, that $c\text{-flow}_{hyp}$ introduces some errors especially for higher $c\text{-flow}$ estimates ($c\text{-flow} > \xi_1$), as there are some outliers around the bisecting line. However, the corresponding bar plots depict that the errors are produced for a minority of samples. Thus, our evaluation indicates that the concept of \mathcal{B}_{hyp} generates accurate results for most FN as the results for $c\text{-flow}_{hyp}$ identified $\approx 97\%$ of all FN that were also labeled as FN by means of $c\text{-flow}_{GT}$.

5 Related Work

Although optical flow has been widely researched to enhance detection performance [15,8,9], the scope of our work is to exploit optical flow V&V in the context of AD, *i.e.* we want to determine safety-critical faults in the DNN’s output without performing particular architectural changes.

Two papers that focus on monitoring perception correctness for segmentation are [19,21]. The most similar approach to ours is proposed by Varghese *et al.* [19] which presents a temporal consistency metric to measure the stability of consecutive semantic segmentation predictions utilizing optical flow. Varghese *et al.* [19] argue that their approach may be used as a supplementary observer (online-monitor) to support safety requirements which also motivates our work. Yet, the proposed metric weights all pixels equally whereas we focus on object detection, *i.e.* we are particularly interested in investigating the credibility of detections in safety-critical sequences.

Zhang *et al.*[21] propose a perceptual consistency measure between segmentation maps on consecutive image pairs to capture the temporal consistency of a video segmentation. Unlike optical flow, the perceptual consistency measure does not seek for exact pixel correspondences across two images but finds pairs of maximally correlated pixels to

mitigate the impact of occlusion. Our approach, however, utilizes optical flow to measure motion-based inconsistencies within bounding boxes. This enables us to identify generally difficult cases for object detection like occlusion.

6 Conclusion

Safety-agnostic evaluation is a great safety concern for AD perception. To consider safety in the evaluation of a pedestrian detector, this work leverages additional information extracted from sequences of images to inject domain-knowledge. Particularly, it leverages temporal consistency information from optical flow to estimate the credibility of detections in a time-coherent pedestrian track. Therefore, we introduce a novel metric *c-flow* to identify false negative predictions of a pedestrian detector. To demonstrate the validity of *c-flow*, (i) we perform controlled experiments with ground truth bounding boxes on a large AD dataset to show that *c-flow* achieves high accuracy with a low number of false alarms. Furthermore, (ii) we perform the computation of *c-flow* without ground truth based on the concept of hypothesized bounding boxes and demonstrate that even in the unsupervised case *c-flow* provides valid and very promising results. This qualifies *c-flow* for perspective runtime applications such as an observer for safety-critical misdetections or active learning.

As future work, we see that the detection of non-credible samples can be further optimized, e.g. using optical flow vectors to refine the estimated hypothesized bounding boxes [18].

References

1. Abrecht, S., Hirsch, A., Raafatnia, S., Woehle, M.: Deep learning safety concerns in automated driving perception. ArXiv **abs/2309.03774** (2023), <https://api.semanticscholar.org/CorpusID:261582607> 1
2. Board, N.T.S.: Collision between vehicle controlled by developmental automated driving system and pedestrian. <https://www.nts.gov/investigations/Pages/HWY18MH010.aspx> (18032018), accessed: 17.11.2023 2
3. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuScenes: A multimodal dataset for autonomous driving. CoRR **abs/1903.11027** (2019), <http://arxiv.org/abs/1903.11027> 7
4. Ceccarelli, A., Montecchi, L.: Evaluating object (mis)detection from a safety and reliability perspective: Discussion and measures. IEEE Access **11**, 44952–44963 (2022), <https://api.semanticscholar.org/CorpusID:252762677> 1
5. Chang, M., Lambert, J., Sangkloy, P., Singh, J., Bak, S., Hartnett, A., Wang, D., Carr, P., Lucey, S., Ramanan, D., Hays, J.: Argoverse: 3D tracking and forecasting with rich maps. CoRR **abs/1911.02620** (2019), <http://arxiv.org/abs/1911.02620> 3, 7
6. Geissler, F., Qutub, S., Paulitsch, M., Pattabiraman, K.: A low-cost strategic monitoring approach for scalable and interpretable error detection in deep neural networks. In: Computer Safety, Reliability, and Security: 42nd International Conference, SAFECOMP 2023, Toulouse, France, September 20–22, 2023, Proceedings. pp. 75–88. Springer-Verlag, Berlin, Heidelberg (2023). https://doi.org/10.1007/978-3-031-40923-3_7, https://doi.org/10.1007/978-3-031-40923-3_7 2
7. Hawkins, R., Conmy, P.R.: Identifying run-time monitoring requirements for autonomous systems through the analysis of safety arguments. In: International Conference on Computer Safety, Reliability, and Security (2023), <https://api.semanticscholar.org/CorpusID:261894104> 2

8. Hu, Q., Wang, P., Shen, C., van den Hengel, A., Porikli, F.: Pushing the limits of deep cnns for pedestrian detection. *IEEE Transactions on Circuits and Systems for Video Technology* **28**(6), 1358–1368 (2017) [12](#)
9. Kang, K., Li, H., Yan, J., Zeng, X., Yang, B., Xiao, T., Zhang, C., Wang, Z., Wang, R., Wang, X., et al.: T-cnn: Tubelets with convolutional neural networks for object detection from videos. *IEEE Transactions on Circuits and Systems for Video Technology* **28**(10), 2896–2907 (2017) [12](#)
10. Li, M., Wang, Y.X., Ramanan, D.: Towards streaming perception. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (eds.) *Computer Vision – ECCV 2020*. pp. 473–488. Springer International Publishing, Cham (2020) [7](#)
11. Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. In: *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22–29, 2017*. pp. 2999–3007. IEEE Computer Society (2017). <https://doi.org/10.1109/ICCV.2017.324>, <https://doi.org/10.1109/ICCV.2017.324> [3](#)
12. Lyssenko, M., Gladisch, C., Heinzemann, C., Woehrle, M., Triebel, R.: From evaluation to verification: Towards task-oriented relevance metrics for pedestrian detection in safety-critical domains. In: *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. pp. 38–45 (2021). <https://doi.org/10.1109/CVPRW53098.2021.00013> [1](#)
13. Lyssenko, M., Gladisch, C., Heinzemann, C., Woehrle, M., Triebel, R.: Towards safety-aware pedestrian detection in autonomous systems. In: *International Conference on Intelligent Robots and Systems (IROS 2022)*. pp. 293–300. IEEE (2022). <https://doi.org/10.1109/IROS47612.2022.9981309> [1](#), [2](#), [7](#), [8](#)
14. Nishimura, H., Komorita, S., Kawanishi, Y., Murase, H.: SDOF-Tracker: fast and accurate multiple human tracking by skipped-detection and optical-flow. *CoRR* **abs/2106.14259** (2021), <https://arxiv.org/abs/2106.14259> [2](#)
15. Ramzan, H., Fatima, B., Shahid, A., Ziauddin, S., Ali, A.: Intelligent pedestrian detection using optical flow and hog. *International Journal of Advanced Computer Science and Applications* **7** (01 2016). <https://doi.org/10.14569/IJACSA.2016.070955> [12](#)
16. Siefke, L., Sommer, V., Baylan, M.C., Grunske, L.: Probabilistic spatial relations for monitoring behavior of road users. In: *Computer Safety, Reliability, and Security: 42nd International Conference, SAFECOMP 2023, Toulouse, France, September 20–22, 2023, Proceedings*. pp. 151–164. Springer-Verlag, Berlin, Heidelberg (2023). https://doi.org/10.1007/978-3-031-40923-3_12, https://doi.org/10.1007/978-3-031-40923-3_12 [2](#)
17. Teed, Z., Deng, J.: RAFT: recurrent all-pairs field transforms for optical flow. *CoRR* **abs/2003.12039** (2020), <https://arxiv.org/abs/2003.12039> [9](#)
18. True, J., Khan, N.: Motion vector extrapolation for video object detection. *CoRR* **abs/2104.08918** (2021), <https://arxiv.org/abs/2104.08918> [13](#)
19. Varghese, S., Bayzidi, Y., Bär, A., Kapoor, N., Lahiri, S., Schneider, J.D., Schmidt, N., Schlicht, P., Hüger, F., Fingscheidt, T.: Unsupervised temporal consistency metric for video segmentation in highly-automated driving. In: *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. pp. 1369–1378 (2020). <https://doi.org/10.1109/CVPRW50498.2020.00176> [2](#), [4](#), [12](#)
20. Wolf, M., Douat, L.R., Erz, M.: Safety-aware metric for people detection. In: *24th IEEE International Intelligent Transportation Systems Conference, ITSC 2021, Indianapolis, IN, USA, September 19–22, 2021*. pp. 2759–2765. IEEE (2021). <https://doi.org/10.1109/ITSC48978.2021.9564734>, <https://doi.org/10.1109/ITSC48978.2021.9564734> [1](#)
21. Zhang, Y., Borse, S., Cai, H., Wang, Y., Bi, N., Jiang, X., Porikli, F.: Perceptual consistency in video segmentation. In: *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022, Waikoloa, HI, USA, January 3–8, 2022*. pp. 2623–2632. IEEE

(2022). <https://doi.org/10.1109/WACV51458.2022.00268>, <https://doi.org/10.1109/WACV51458.2022.00268> **12**