# Is my new tracker really better than yours?

Luka Čehovin, Matej Kristan, and Aleš Leonardis
Faculty of Computer and Information Science, University of Ljubljana, Slovenia
Tržaška 25, Ljubljana, Slovenia
{luka.cehovin,matej.kristan,ales.leonardis}@fri.uni-lj.si

## Abstract

*The problem of visual tracking evaluation is sporting an abundance of performance measures, which are used by various authors, and largely suffers from lack of consensus about which measures should be preferred. This is hampering the cross-paper tracker comparison and faster advancement of the field. In this paper we provide an overview of the popular measures and performance visualizations and their critical theoretical and experimental analysis. We show that several measures are equivalent from the point of information they provide for tracker comparison and, crucially, that some are more brittle than the others. Based on our analysis we narrow down the set of potential measures to only two complementary ones that can be intuitively interpreted and visualized, thus pushing towards homogenization of the tracker evaluation methodology.*

## 1. Introduction

Visual tracking is one of the rapidly evolving fields of computer vision. Every year, literally dozens of new tracking algorithms are presented and evaluated in journals and at conferences. When considering the evaluation of these new trackers and comparison to the state-of-the-art, several questions arise. Is there a standard set of sequences that we can use for the evaluation? Is there a standardized evaluation protocol? What kind of performance measures should we use? Unfortunately, there are currently no definite answers to these questions. Unlike some other fields of computer vision, like object detection and classification [11], optical-flow computation [4] and automatic segmentation [2], where widely adopted evaluation protocols are used, visual tracking is still largely lacking these features.

The absence of homogenization of the evaluation protocols makes it difficult to rigorously compare trackers across publications and stands in the way of faster development of the field. The authors of new trackers typically compare their work against a limited set of related algorithms due to the difficulty of adapting these for their own use in the experiments. The issue here is the choice of tracker's performance evaluation measures, which seems to be almost arbitrary in the tracking literature. Worse yet, an abundance of these measures are currently in use. Because of this, experiments in many cases offer a limited insight into tracker's performance, and prohibit comparison across different papers.

In this paper we focus on the problem of performance evaluation in monocular single-target visual tracking and address several challenges therein. We investigate various popular performance evaluation measures, discuss their pitfalls and show that, from a standpoint of tracker comparison, there exist several equivalent measures currently in use. From there on we identify complementary measures that are sensitive to two different aspects of tracker's performance. The goal of our analysis is to homogenize of the tracking performance evaluation methodology and increase the interpretability of results. It is worth noting that our findings have been so far already used as the foundation of the evaluation methodology of a visual tracking challenge, whose results have been presented at a workshop at a major computer vision conference.

### 1.1. Related work

Until recently the majority of papers that address performance evaluation in visual tracking were concerned with multi-target tracking scenarios [32, 17, 8, 9, 10, 25, 6, 24]. One might view the multi-target tracking as a generalization of single-target tracking, however, there is a crucial difference in the focus of the evaluation. In multi-target tracking, the focus is usually on measuring correctness of target labeling assignments coupled with target detection and occlusion handling. The reason is that the algorithms are often focused on a particular tracking domain, which is typically people or vehicle tracking for surveillance [8, 9, 15], animal groups tracking [18] or sports tracking [21], to name a few. A well known PETS workshop (e.g. [6]) has also been organized yearly for more than a decade with the main focus on performance evaluation of surveillance and activity recognition algorithms.

On the other hand, single-target visual tracking evalua-

tion focuses on the tracker's accuracy, robustness and generality. The goal is to demonstrate the tracker's performance on a wide range of challenging scenarios (various types of object, lighting conditions, camera motion, signal noise, etc.). In this respect, the authors for [35] compared several trackers using center error and overlap measures. Their research is focused primarily on investigating strengths and weaknesses of a few trackers. In [36] authors perform an experimental comparison of several trackers. The performance measures in this case are not well chosen which results in a poor qualitative analysis of the results. Nawaz and Cavallaro [26] have presented a system for evaluation of video trackers that aims at addressing the real-world conditions. The system can simulate several real-world sources of noisy input, such as initialization noise, image noise and changes in the frame-rate. They have also proposed a new performance measure to address the tracker's scoring under these simulated conditions that is discussed in Section 2.6. These recent experimental evaluations show the need for a better evaluation of visual trackers, however, none of them addresses an important prerequisite for such evaluation, that is the selection of good performance measures. This selection should be grounded in an analysis of performance measures which is the main focus of this paper.

Recently, Smeulders et al. [31] provided an experimental survey of several recent trackers together with an analysis of several performance measures. Their methodology and the general disposition in this aspect are similar to ours in terms that they search for multiple measures that describe different aspects of tracking performance. However, their selection of measures is from the start biased in favor of detection-based tracking algorithms, which also affects their choice of final measures and the derived conclusions.

Lastly, an interesting idea has been proposed by Pang and Habin [28], who aggregate existing experiments, published in various articles, in a page-rank fashion to form a less biased ranking of trackers. The authors acknowledge that their approach is not appropriate for ranking recently published trackers. Furthermore, their approach does not remove bias that comes from correlation in multiple performance measures, which is one of the goals of our work.

### 1.2. Our approach and contributions

The goal of this paper is not to propose new performance measures. Instead we focus on narrowing the wide variety of existing measures for single-target tracking performance evaluation to only a few complementary ones. This is a crucial step towards the homogenization of the field of is. We claim a three-fold contribution: (1) We provide a detailed survey and experimental analysis of performance measures used in single-target tracking evaluation. (2) We show by experimental analysis that there exist clusters of
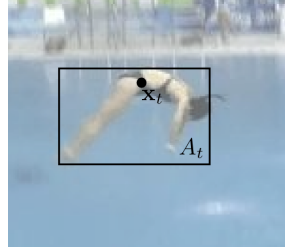


Figure 1. An example of an annotation for a single frame for the *driver* sequence. The center of the object can be estimated using the centroid of $A_t$, which is not true in this case.

performance measures that essentially indicate the same aspect of tracker's performance. (3) By considering the theoretical aspects of existing measures as well as the experimental analysis we identify the two most suitable complementary measures that characterize tracker's performance within the accuracy vs. robustness context and propose an intuitive way to visualize the selected pair of measures.

Our experimental analysis has been carried out in a form of a comparative experiment with 13 state-of-the-art trackers and 25 widely-used video sequences. While the goal of this paper is not to establish the ranking of these trackers, we nevertheless provide detailed results of the experiment as a side-product of our research in the supplementary material [1].

The rest of the paper is organized as follows: Section 2 gives an overview of the current state of performance evaluation techniques. Section 3 describes our experimental setup. We discuss the findings of the experiment in Section 4 where we also propose our selection of good measures together with several insights. Finally, we draw concluding remarks in Section 5.

## 2. Performance measures

There are several measures that have become popular and are widely used in the literature, however, none of them is a *de-facto* standard. As all of these measures assume that manual annotations are given for a sequence, we first establish a general definition of an object state description in a sequence with length $N$ as:

$$\Lambda = \{(A_t, \mathbf{x}_t)\}_{t=1}^{N}, \tag{1}$$

where $\mathbf{x}_t \in \mathcal{R}^2$ denotes a center of the object and $A_t$ denotes the region of the object at time $t$. In practice the region is usually described by a bounding box (that is most commonly axis-aligned), however, a more complex shape could be used for a more accurate description. An example of a single frame annotation can be seen in Figure 1. Performance measures aim at summarizing the extent to which the tracker's predicted annotation $\Lambda_T$ agrees with the ground truth annotation, i.e., $\Lambda_G$.

## 2.1. Center error

Perhaps the oldest means of measuring performance, which has its roots in aeronautics, is the center prediction error. This is still a popular measure [30, 3, 1, 23] and it measures the difference between the target's predicted center from the tracker and the ground-truth center.

$$\Delta(\Lambda^G, \Lambda^T) = \{\delta_t\}_{t=1}^N, \quad \delta_t = \|\mathbf{x}_t^G - \mathbf{x}_t^T\|. \quad (2)$$

The popularity of center prediction measure comes from its minimal annotation effort, i.e., only a single point per frame. The results are usually shown in a plot, as in Figure 9 or summarized as average error (3), or root-mean-square-error (4):

$$\Delta_\mu(\Lambda^G, \Lambda^T) = \frac{1}{N} \sum_{t=1}^N \delta_t, \quad (3)$$

$$\text{RMSE}(\Lambda^G, \Lambda^T) = \sqrt{\frac{1}{N} \sum_{t=1}^N \|\mathbf{x}_t^G - \mathbf{x}_t^T\|^2}. \quad (4)$$

One drawback of this measure is its sensitivity to subjective annotation (i.e., where exactly is the target's center). This sensitivity largely comes from the fact that the measure completely ignores the target's size and does not reflect the apparent tracking failure [26]. To remedy this, a normalized center error $\widehat{\Delta}(\cdot, \cdot)$ is used instead, e.g. [5, 31], in which the center error at each frame is divided by the tacker-predicted visual size of the target, $size(A_t^G)$,

$$\widehat{\Delta}(\Lambda^G, \Lambda^T) = \left\{\widehat{\delta}_t\right\}_{t=1}^N, \quad \widehat{\delta}_t = \|\frac{\mathbf{x}_t^G - \mathbf{x}_t^T}{size(A_t^G)}\|. \quad (5)$$

Nevertheless, despite the normalization, the measure may give misleading results as the center error is reduced proportionally to the estimated target size. Furthermore, when the tracker fails and is drifting over a background, the actual distance between the annotated and reported center, combined with the estimated size (which can be arbitrarily large) influences the averaged score and does not properly reflect the important information that the tracker has failed.

## 2.2. Region overlap

The normalization problem is rather well addressed by the overlap-based measures [37, 13, 31]. These measures require region annotations and are computed as an overlap between predicted target's region form the tracker and the ground-truth region:

$$\Phi(\Lambda^G, \Lambda^T) = \{\phi_t\}_{t=1}^N, \quad \phi_t = \frac{A_t^G \cap A_t^T}{A_t^G \cup A_t^T}. \quad (6)$$

A nice property of region overlap measures is that they account for both position and size of the predicted and ground-truth bounding boxes simultaneously, and do not result in arbitrary large errors at tracking failures. In fact, once

the tracker drifts to the background, the measure becomes zero, regardless of how far from the target the tracker is currently located. In terms of pixel classification (see Figure 2), the overlap can be interpreted as

$$\frac{A_t^G \cap A_t^T}{A_t^G \cup A_t^T} = \frac{TP}{TP + FN + FP}, \quad (7)$$

a formulation similar to the F-measure in information retrieval, which can be written as $F = \frac{2TP}{2TP + FN + FP}$. Another closely related measure, used in tracking to account for un-annotated object occlusions is precision [13], or $\frac{TP}{TP+FP}$.
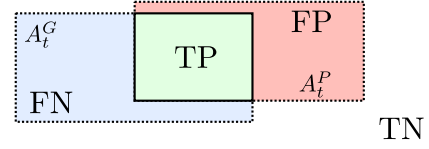


Figure 2. An illustration of the overlap of ground-truth region with the predicted region.
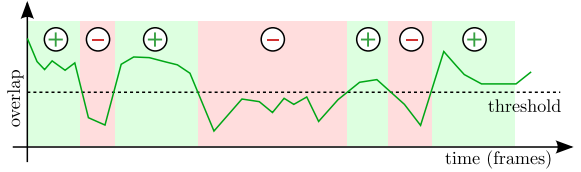


Figure 3. An illustration of overlap being used as a detection measure.

The overlap measure is summarized over an entire sequence by an average overlap, or as a number of correctly tracked frames. The latter approach comes from the object detection community [11], where the overlap threshold for a correctly detected object is set to $0.5$. The same threshold is often used for tracking performance evaluation, e.g. in [37, 35]. To make the final score more comparable across the different sequences, the number of correctly tracked frames is divided by the total number of frames

$$P_\tau(\Lambda^G, \Lambda^T) = \frac{\|\{t|\phi_t > \tau\}_{t=1}^N\|}{N}, \quad (8)$$

where $\tau$ denotes the threshold of the overlap. The $P_\tau$ is a frame-wise definition of the *true-positive* score, an interpretation that has become popular in tracking evaluation with the advent of tracking-by-detection concept. As noted in [31], the F-measure is another score that can be used in this context, however, it is worth noting that the detection based measures favor disregard the sequential nature of the tracking problem. As it is illustrated in Figure 3, these measures do not necessarily account for complete trajectory reconstruction which is an important aspect in many tracking applications.
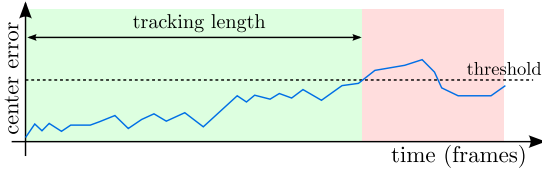
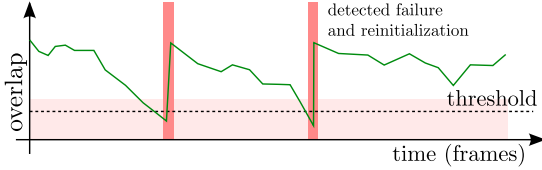Figure 4. An illustration of the tracking length measure for center error.



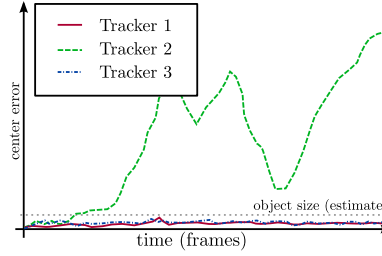Figure 5. An illustration of the failure rate measure for overlap distance.



Figure 6. An example of center-error plot comparison for three trackers. Tracker 2 has clearly failed in the process, yet its large center errors cause the plot to expand its vertical scale, thus reducing the apparent differences of trackers 1 and 3.
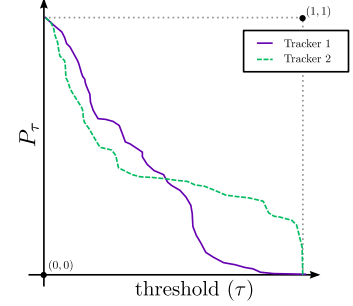


Figure 7. An illustration of the *measure-threshold* plot for two trackers. It is apparent that different values of the threshold would clearly yield different rankings for the trackers.

## 2.3. Tracking length

Another measure that has been used in the literature to compare trackers is *tracking length* [22]. This measure reports the number of successfully tracked frames from tracker's initialization to its (first) failure. A failure criterion can be manual visual inspection (e.g. [13]), which is biased and cannot be repeated reliably even by the same person. A better approach is to automate the failure criterion, e.g., by placing a threshold $\tau$ on the center or overlap measure (see Figure 4). The choice of the criterion may impact the result of comparison. Given the sensitivity of the center-based measures, an overlap criterion makes more sense, and we will denote in the following the tracking length measure with an overlap-based failure criterion by $L_\tau$.

While this measure explicitly addresses the tracker's failure cases, which the simple average center-error and overlap measures do not, it suffers from a significant drawback. Namely, it only uses the part of the video up to the first tracking failure. If by some coincidence, the beginning of the video contains a difficult tracking situation, or the target is not visible well, which results in a necessarily poor initialization, the tracker will fail, and the remainder of the video will be discarded. This means that, technically, one would require a significant amount of videos exhibiting the various properties right at its beginning to get a good statistic on this performance measure.

## 2.4. Failure rate

A measure that largely addresses the problem of the tracking length measure is the so-called failure rate measure [20, 18]. The failure rate measure casts the tracking problem as a supervised system in which an operator reinitializes the tracker once it fails. The number of required manual interventions per frame is recorded and used as a comparative score. The approach is illustrated in Figure 5. This measure also reflects the tracker's performance in a real-world situation in which the human operator supervises the tracker and corrects its errors.

Compared to the tracking length measure, the failure rate approach has the advantage that the entire sequence is used in the evaluation process and decreases the importance of the beginning part of the sequence. The question of a failure criterion threshold is even more apparent here as each change in the criterion requires the entire experiment to be repeated. Researchers in [33, 34] consider a failure when the bounding box overlap is lower than $0.1$. This lower threshold reasonable for non-rigid objects, since these are often poorly described by the bounding-box area. An even lower threshold could be used for overlap-based failure criteria if we are interested only in the most apparent failures with no overlap between the regions. We will denote the failure rate measure with an overlap-based failure criterion with threshold $\tau$ as $F_\tau$.

## 2.5. Performance plots

Plots are frequently used to visualize the behavior of a tracker when a single number does not suffice for expressing its performance. The most widely-used plot is a center-error plot that shows the center-error with respect to the frame number [3, 1, 5, 37]. While this kind of plot can be useful for visualizing tracking result of a single tracker, a combined plot for multiple trackers is in many cases misused if applied without caution, because the tracker with an inferior performance steals away the focus from the information that we are interested in with this type of plots, i.e. the tracker accuracy. An illustration of such a problematic plot is shown in Figure 6.

In the previous section we have seen that a failure criterion plays a significant role in visual tracker performance evaluation. Choosing an appropriate value for the threshold may affect the ranking and can also be potentially abused to influence the results of a comparison. However, it is

sometimes better to avoid the use of one specific threshold altogether, especially when the evaluation goal is general and a specific threshold is not a part of the target task. To avoid the choice of a specific threshold, results can be presented as a *measure-threshold* plot. This kind of plots bear some resemblances to a ROC curve [12], like monotony, intuitive visual comparison, and a similar calculation algorithm. Measure-threshold plots were used in [3], where the authors used center-error as a measure as well as in [36], where both center-error and overlap are used.

The percentage of correctly tracked frames, defined in (8) as $P_\tau$, is a good choice for a measure to be used in this scenario, however, other measures could be used as well. The $P_\tau$ measure can be intuitively computed for multiple sequences which makes it useful for summarizing the entire experiment (an example of $P_\tau$ plot is illustrated in Figure 7). Interpretations of such plots have been so far limited to their basic properties which in a way negates the information verbosity of a graphical representation. For example, similarly to ROC curves, we can compute an area-under-the-curve (AUC) summarization score, which is used in [36] to reason about performance of trackers. It can be trivially proven that the AUC score in this case actually matches the average overlap over the entire sequence (proof available in the supplementary material) and therefore adds no additional insight into the performance of the tracker.

A curve that is in shape similar to $P_\tau$ plot is the *survival curve* [31]. In this case the curve summarizes the trackers success (various performance measures can be used) over a dataset of sequences that are ordered from the best performance to worst. While this approach gives a good overview of the overall success, it is unsuitable for sequence-wise comparison as the order of sequences differs from tracker to tracker.

## 2.6. Hybrid measures

Nawaz and Cavallaro [26] propose a threshold-independent overlap-based measure that combines the information on tracking accuracy and tracking failure into a single score. This hybrid measure is called the *Combined Tracking Performance Score* (CoTPS) and is defined as a weighted sum of an accuracy score (based on the frames where the tracker was successful) and a failure score (based on the frames where the tracker failed to predict the position of the object). An appealing property of this measure is that it ranks trackers by accounting for two separate aspects of tracking. However, the authors do not give any justification for this rather complex fusion of measures. This complexity prohibits easy interpretation of the results required for a rigorous scientific analysis.

In terms of interpretation, we therefore believe that a better strategy is to focus on a few complementary performance measures with well-defined meaning, and avoid fus-

ing them into a single measure too soon in the evaluation process.

## 3. Experimental setup

In order to analyze measures, we have conducted a typical comparative experiment. Our goal is to rank several existing trackers according to the selected measures on a number of typical visual tracking sequences. The selection of measures is based on our theoretical discussion in Section 2. We have therefore selected: average center error, average normalized center error, root-mean-square error, average overlap, percent of correct frames $P_{0.1}$, tracking length $L_{0.1}$, percent of correct frames $P_{0.5}$, tracking length $L_{0.5}$, failure rate $F_0$, average overlap for $F_0$.

We have evaluated 13 trackers that were proposed in recent years: A color-based particle filter (PF) [29], the On-line boosting tracker (OBT) [14], the Flock-of-features tracker (FOF) [19], the Basin-hopping Monte Carlo tracker (BHMC) [22], the Incremental visual tracker (IVT) [30], the Histograms-of-blocks tracker (BH) [27], the Multiple instance tracker (MIL) [3], the Fragment tracker (FRT) [1], the P-N tracker (TLD) [16], the Local-global tracker (LGT) [34], Hough tracker (HT) [13], the L1 Tracker Using Accelerated Proximal Gradient Approach (L1-APG) [5] and the Compressive tracker (CT) [37]. The source code of the trackers was provided the authors and adapted to fit into our framework.

We have run the trackers on 25 different sequences, most of which are already known in the video tracking community [34, 33, 37, 35, 30, 22, 1, 13], and several were acquired additionally. The sequences were annotated with a bounding-box region of the object, as well as the target's central point of the object. To account for stochastic processes that are a part of many trackers, each tracker was executed on each sequence 30 times. The tracker's performance on a particular sequence was then evaluated by averaging these results. Parameters for all trackers were set to their default values and kept constant during the experiment. A separate run was executed for the *failure rate* measure as the re-initialization influences other aspects of tracking performance.

Because of the scale of the experiment, only the most relevant results are presented in Section 4. Additional results, such as the ranking of the trackers according to individual measures, are available in the supplementary material.

## 4. Results and discussion

Different measures may reflect different aspects of tracking performance, so it is impossible to simply establish which measure is the best. We start our analysis by establishing similarities and equivalence between various measures, by experimentally analyzing which measures pro-

duce consistently similar responses when comparing trackers.

## 4.1. Correlation analysis

Similarly to [31], we calculate a correlation matrix from all pairs of measures calculated over all tracker-sequence pairs. Note that we do not calculate the correlation on rankings directly to avoid handling situations where several trackers take the same place (if differences are not statistically significant). The rationale is that strongly correlated measure values will also produce similar ranking for trackers. The obtained correlation matrix is shown in Figure 8. We can see that two clusters emerge, one for measures 1 to 3 and one for measures 4 to 7, where measure 7 is less correlated to the other three. All these correlations are highly statistically significant ($p < 0.001$).

The first cluster of measures consists of the three center-error-based measures. This is expected as these measures all base on *center-error* using different averaging methods. The second cluster of measures contains *average overlap*, *percentage of correctly tracked frames* for two threshold values ($P_{0.1}$ and $P_{0.5}$) and *tracking length* $L_{0.1}$. Measures in the second cluster assume that incorrectly tracked frames do not influence the final score based on the specific (incorrect) position of the tracker. This makes them more representative as a measure for tracking performance than the center-error-based measures. An illustration of this difference for *overlap* and *center-error* is shown as a graph on Figure 9, where we can clearly see that the center-error measure takes into account the exact center distance at frames after the failure has occurred, which depends on the movement of an already failed tracker and does not reflect its true performance.



1. Average center error,
2. Average normalized center error,
3. Root-mean-square error,
4. Average overlap,
5. Percent of correct frames $P_{0.1}$,
6. Tracking length $L_{0.1}$,
7. Percent of correct frames $P_{0.5}$,
8. Tracking length $L_{0.5}$,
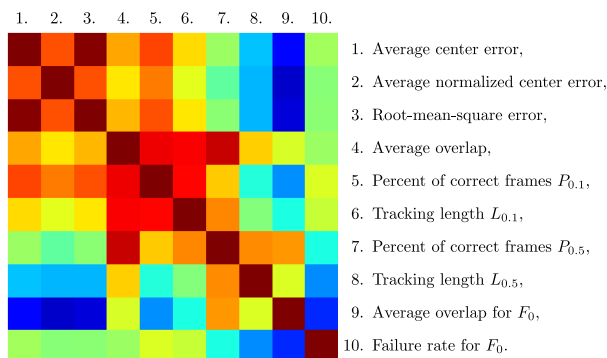9. Average overlap for $F_0$,
10. Failure rate for $F_0$.

Figure 8. Correlation matrix for all measures visualized as a heatmap. The image is best viewed in color.

The first cluster implies that the first three measures are equivalent and it does not matter which one you chose. The second cluster requires some more interpretation. Despite the apparent similarity of overlap-based measures 4 to 7, the correlation is not perfect and the rankings differ in some

cases. One example of such difference can be seen for the TLD tracker on the *woman* sequence, seen in Figure 11. We can see that the tracker loses the target early on in the sequence (during an occlusion), but manages to locate it again later because of its discriminative nature. The *average overlap* measure (number 4) and the *percentage of correct frames* measure (number 5) therefore rank the tracker higher than the *tracking length* measure (number 6). On the general level we can also observe that the choice of a threshold can influence the outcome of the experiment. This can be observed for tracking length measures 6 and 8 and to some extent for the percentage of correct frames measures 5 and 7. There, the scores for a higher threshold (0.5) result in a different ranking of trackers compared to the lower threshold (0.1). This means that care must be taken when choosing the thresholds at they may affect the outcome of the evaluation.

We can in fact observe a slight overlap between the first two clusters in the correlation matrix, implying similarity in their information content. Based on the above analysis and discussion in Section 2 we conclude that the *average overlap* measure is the most appropriate to be used in tracker comparison, as it is simple to compute, it is scale and threshold invariant, exploits the entire sequence, and it is easy to interpret. Note also that it is highly correlated with a more complex *percentage-of-correctly-tracked-frames* measure.

## 4.2. Failure rate

As mentioned before, the *failure rate* measure influences the tracker's entire trajectory, because of the re-initializations, therefore the data for measures 9 and 10 was acquired separately. The advantage of the *failure rate* (measure 10) is that the entire sequence is used, which makes the results statistically significant at smaller number of sequences. It does not matter that much if one tracker fails at the "difficult" beginning of the sequence, while the other one barely survives and then tracks the rest successfully. In Figure 12 we can see the performance of the LGT tracker on the *bicycle* sequence. Because of the occlusion near frame 175 the tracker fails, although it is clearly capable of tracking the rest of the sequence reliably if re-initialized.

## 4.3. Accuracy vs. robustness

The *failure rate* measure itself measures the robustness of the tracker, however, it tells us nothing about its accuracy. We therefore propose to use the *average overlap* measure on the same (re-initialized) data to take into account this aspect of tracking. We define a new, A-R measure, as a pair of scores

$$\text{A-R}(\Lambda^G, \Lambda^T) = \left( \overline{\Phi}(\Lambda^G, \Lambda^T), F_0 \right), \qquad (9)$$

where $\overline{\Phi}$ denotes *average overlap* and $F_0$ denotes the *failure rate* for $\tau = 0$. Note that the value of failure threshold $\tau$ can
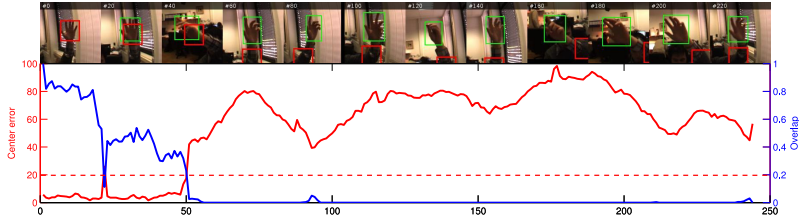
Figure 9. A comparison of overlap and center error distance measures for tracker CT on sequence *hand* [33]. The dashed line shows the estimated threshold above which the center error is greater than the size of the object. The tracker fails around frame 50.



Figure 10. An accuracy-reliability data visualization for all trackers over all sequences.



Figure 11. An overlap plot for tracker TLD on sequence *woman* [1]. The dashed line shows the threshold below which the tracking length detects failure (for threshold 0.1), which happens around frame 120.



Figure 12. An overlap plot for tracker LGT on sequence *bicycle* [34]. The green plot shows the unsupervised overlap, and the blue plot shows the overlap for supervised tracking, where the failure is recorded and the tracker re-initialized.

influence the final results. If the value is set to a high value (i.e. close to 1) the tracker is restarted frequently even for small errors and the final score is hard to interpret. Instead, we propose to use the lowest theoretical threshold $\tau = 0$ to only measure complete failures where the regions have no overlap at all and a reinitialization is clearly justified. In theory a tracker can also report an extremely large region as the position of the target and avoids failures, however, the accuracy will be very low in this case. This is how the two measures complement each other.

The A-R pair of each can be visualized as a 2-D point plot as seen in Figure 10, where we show the average scores for all sequences, from which one can read the tracker's performance in terms of accuracy (the tracker is more accurate if it is higher along vertical axis) and reliability (the tracker fails less if it is further to the right on horizontal axis). Because the robustness does not have an upper bound we propose to interpret it as a reliability for visualization purposes. The reliability of a tracker is defined as $e^{-S^{F_0}/N}$ and is interpreted as a probability that the tracker will still successfully track the object up to $S$ frames since the last failure. Here the failure probability is modeled using an exponential failure distribution based on $F_0$. Note that the choice of $S$ does not influence the order of the trackers, however, changing its value can be useful for visualization and interpretation of results.

For a better understanding of the complementing nature of the two measures we introduce two theoretical trackers. The first one, denoted by *T0*, always reports the region of the object to equal the image size of the sequence. This tracker provides too loose regions, but does not fail and is
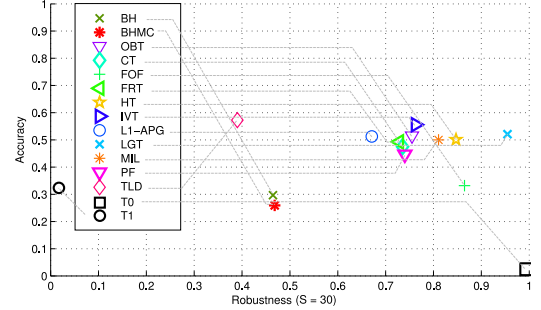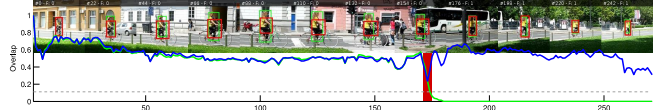
therefore displayed closer to the bottom-right corner. Another tracker, denoted by *T1*, reports its initial position for the entire sequence. This tracker will likely fail a lot, and will achieve decent accuracy because of frequent manual interventions. However, because of numerous failures it would be displayed near the left edge of the plot, while an ideal tracker would reach the top-right corner.

The A-R plot, seen in Figure 10, shows that the LGT tracker is on average the most robust one in the evaluated set, while it is not necessary the most accurate one. As the averaged results can convey only a limited amount of information, we invite the reader to look at the supplementary material for more detailed results that include also persequence A-R plots.

## 5. Conclusion

In this paper we have provided an analysis of the measures for single-target video tracking performance evaluation with the goal of determining the most appropriate ones. Besides a theoretical description and categorization of frequently used measures, we have also performed a large-scale experimental evaluation on a diverse set of 13 well established or recently presented trackers and 25 available evaluation sequences, in order to determine the properties and relationships of measures. From the experiment we can conclude that the selection of a single measure can influence the result of the comparison. Results, summarized by any single measure cannot sufficiently describe the performance of such a complex system as a visual tracker, nor can they help discover the true meaning of its failure, especially on a large set of testing sequences. As a main conclusion of our

experiment we propose that a pair of two complementary measures. This pair takes into account the accuracy (using *average overlap*) and the robustness (using *failure rate*) of each tracker.

While narrowing down the abundance of performance measures is a big step toward homogenizing the tracking evaluation methodology this is only of the requirements for a consistent evaluation methodology for video trackers. The measures that were proposed in this paper have been already adopted as the foundation of the evaluation methodology of a recently organized visual tracking challenge, where a rigorous analysis in terms of accuracy and robustness has provided multiple interesting insights into performance of individual trackers. In our future work we will also address the quality and reliability of manual ground-truth annotations by various people [7].

# References

[1] A. Adam, E. Rivlin, and I. Shimshoni. Robust Fragments-based Tracking using the Integral Histogram. In *CVPR*, volume 1, pages 798–805. IEEE Computer Society, June 2006. 3, 4, 5, 7

[2] S. Alpert, M. Galun, R. Basri, and A. Brandt. Image Segmentation by Probabilistic Bottom-Up Aggregation and Cue Integration. In *CVPR*, pages 1–8. IEEE, June 2007. 1

[3] B. Babenko, M.-H. Yang, and S. Belongie. Robust Object Tracking with Online Multiple Instance Learning. *TPAMI*, 33(8):1619–1632, Aug. 2011. 3, 4, 5

[4] S. Baker, D. Scharstein, J. P. Lewis, S. Roth, M. J. Black, and R. Szeliski. A Database and Evaluation Methodology for Optical Flow. *IJCV*, 92(1):1–31, Nov. 2010. 1

[5] C. Bao, Y. Wu, H. Ling, and H. Ji. Real time robust L1 tracker using accelerated proximal gradient approach. In *CVPR*, pages 1830–1837. IEEE, June 2012. 3, 4, 5

[6] F. Bashir and F. Porikli. Performance Evaluation of Object Detection and Tracking Systems. 2006. 1

[7] A. Baumann, M. Boltz, J. Ebling, M. Koenig, H. Loos, M. Merkel, W. Niem, J. Warzelhan, and J. Yu. A Review and Comparison of Measures for Automatic Video Surveillance Systems. *EURASIP Journal on Image and Video Processing*, 2008. 8

[8] J. Black, T. Ellis, and P. Rosin. A novel method for video tracking performance evaluation. In *VS-PETS*, pages 125–132, 2003. 1

[9] L. M. Brown, A. W. Senior, Y.-l. Tian, J. Connel, A. Hampapur, C.-F. Shu, H. Merkl, and M. Lu. Performance evaluation of surveillance systems under varying conditions. In *PETS*, pages 1–8, 2005. 1

[10] K. K. Edward, P. D. Matthew, and B. H. Michael. An information theoretic approach for tracker performance evaluation. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1523–1529. IEEE, Sept. 2009. 1

[11] M. Everingham, L. Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The Pascal Visual Object Classes (VOC) Challenge. *IJCV*, 88(2):303–338, Sept. 2009. 1, 3

[12] T. Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874, June 2006. 5

[13] M. Godec, P. M. Roth, and H. Bischof. Hough-based tracking of non-rigid objects. In *ICCV*, pages 81–88. IEEE, Nov. 2011. 3, 4, 5

[14] H. Grabner, M. Grabner, and H. Bischof. Real-time tracking via on-line boosting. In *BMVC*, pages 47–56, 2006. 5

[15] C. Jaynes, S. Webb, R. Steele, and Q. Xiong. An open development environment for evaluation of video surveillance systems. In *PETS*, 2002. 1

[16] Z. Kalal, J. Matas, and K. Mikolajczyk. P-N learning: Bootstrapping binary classifiers by structural constraints. In *CVPR*, pages 49–56. IEEE, June 2010. 5

[17] R. Kasturi, D. Goldgof, P. Soundararajan, V. Manohar, J. Garofolo, R. Bowers, M. Boonstra, V. Korzhova, and J. Zhang. Framework for performance evaluation of face, text, and vehicle detection and tracking in video: data, metrics, and protocol. *TPAMI*, 31(2):319–36, Feb. 2009. 1

[18] Z. Khan, T. Balch, and F. Dellaert. MCMC-based particle filtering for tracking a variable number of interacting targets. *TPAMI*, 27:1805–1819, 2005. 1, 4

[19] M. Kölsch and M. Turk. Fast 2d hand tracking with flocks of features and multi-cue integration. In *CVPR Workshops*, volume 10, page 158, Washington, DC, USA, 2004. IEEE Computer Society. 5

[20] M. Kristan, S. Kovačič, A. Leonardis, and J. Perš. A two-stage dynamic model for visual tracking. *Trans. Sys. Man Cyber. Part B*, 40(6):1505–1520, Dec. 2010. 4

[21] M. Kristan, J. Perš, M. Perše, and S. Kovačič. Closed-world tracking of multiple interacting targets for indoor-sports applications. *CVIU*, 113(5):598–611, May 2009. 1

[22] J. Kwon and K. M. Lee. Tracking of a non-rigid object via patch-based dynamic appearance modeling and adaptive basin hopping monte carlo sampling. In *CVPR*, pages 1208–1215, 2009. 4, 5

[23] J. Kwon and K. M. Lee. Visual tracking decomposition. In *CVPR*, pages 1269–1276. IEEE, June 2010. 3

[24] I. Leichter and E. Krupka. Monotonicity and Error Type Differentiability in Performance Measures for Target Detection and Tracking in Video. *TPAMI*, 35(10):2553–2560, 2013. 1

[25] V. Manohar, P. Soundararajan, H. Raju, D. Goldgof, R. Kasturi, and J. Garofolo. *Computer Vision ACCV 2006*, volume 3852 of *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, Berlin, Heidelberg, Jan. 2006. 1

[26] T. Nawaz and A. Cavallaro. A protocol for evaluating video trackers under real-world conditions. *Tran. Image Proc.*, 2012. 2, 3, 5

[27] S. S. Nejhum, J. Ho, and M.-H. Yang. Online visual tracking with histograms and articulating blocks. *CVIU*, 114(8):901–914, Aug. 2010. 5

[28] Y. Pang and H. Ling. Finding the Best from the Second Bests Inhibiting Subjective Bias in Evaluation of Visual Tracking Algorithms. In *ICCV*, pages 2784–2791, 2013. 2

[29] P. Pérez, C. Hue, J. Vermaak, and M. Gangnet. Color-based probabilistic tracking. In *ECCV*, volume 1, pages 661–675. Springer-Verlag, 2002. 5

[30] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang. Incremental Learning for Robust Visual Tracking. *IJCV*, 77(1-3):125–141, May 2008. 3, 5

[31] A. W. M. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah. Visual Tracking: an Experimental Survey. *TPAMI*, 2013. 2, 3, 5, 6

[32] K. Smith, D. Gatica-Perez, and J. Odobez. Evaluating Multi-Object Tracking. In *CVPR Workshops*, volume 3, pages 36–36. IEEE, 2005. 1

[33] L. Čehovin, M. Kristan, and A. Leonardis. An adaptive coupled-layer visual model for robust visual tracking. In *ICCV*, Barcelona, 2011. 4, 5, 7

[34] L. Čehovin, M. Kristan, and A. Leonardis. Robust Visual Tracking using an Adaptive Coupled-layer Visual Model. *TPAMI*, 35(4):941–953, Apr. 2013. 4, 5, 7

[35] Q. Wang, F. Chen, W. Xu, and M.-H. Yang. An experimental comparison of online object-tracking algorithms. In *SPIE*, pages 81381A–81381A–11, San Diego, 2011. 2, 3, 5

[36] Y. Wu, J. Lim, and M.-h. Yang. Online Object Tracking: A Benchmark. In *CVPR*, 2013. 2, 5

[37] K. Zhang, L. Zhang, and M.-H. Yang. Real-time Compressive Tracking. In *ECCV*, 2012. 3, 4, 5