

Towards More Flexible and Accurate Object Tracking with Natural Language: Algorithms and Benchmark

Xiao Wang^{1*}, Xiujun Shu^{2,1*}, Zhipeng Zhang³, Bo Jiang⁴, Yaowei Wang¹, Yonghong Tian^{1,5}, Feng Wu^{1,6}

¹Peng Cheng Laboratory, Shenzhen, China

²School of Electronic and Computer Engineering, Peking University, Shenzhen, China

³NLPR, Institute of Automation, Chinese Academy of Sciences

⁴School of Computer Science and Technology, Anhui University, Hefei, China

⁵Department of Computer Science and Technology, Peking University, Beijing, China

⁶University of Science and Technology of China, Hefei, China

<https://sites.google.com/view/langtrackbenchmark/>

Abstract

Tracking by natural language specification is a new rising research topic that aims at locating the target object in the video sequence based on its language description. Compared with traditional bounding box (BBox) based tracking, this setting guides object tracking with high-level semantic information, addresses the ambiguity of BBox, and links local and global search organically together. Those benefits may bring more flexible, robust and accurate tracking performance in practical scenarios. However, existing natural language initialized trackers are developed and compared on benchmark datasets proposed for tracking-by-BBox, which can't reflect the true power of tracking-by-language. In this work, we propose a new benchmark specifically dedicated to the tracking-by-language, including a large scale dataset, strong and diverse baseline methods. Specifically, we collect 2k video sequences (contains a total of 1,244,340 frames, 663 words) and split 1300/700 for the train/testing respectively. We densely annotate one sentence in English and corresponding bounding boxes of the target object for each video. We also introduce two new challenges into TNL2K for the object tracking task, i.e., adversarial samples and modality switch. A strong baseline method based on an adaptive local-global-search scheme is proposed for future works to compare. We believe this benchmark will greatly boost related researches on natural language guided tracking.

1. Introduction

Single object tracking is one of the most important tasks in computer vision and it has been widely used in many applications such as video surveillance, robotics, and autonomous vehicles. Usually, they initialize the target object in the first frame with a bounding box (BBox), as shown in Fig. 1 (a), and adjust the BBox along with the movement of the target object. Most of the existing single object trackers [25–27, 64, 66, 80] are developed based on this setting¹, and many benchmark datasets [19, 28, 44, 50, 58, 70, 71] are proposed for this task.

Although these trackers have been adopted in many applications, however, the setting of tracking-by-BBox still suffers from the following issues. (1) The target object in the first frame with a BBox is inconvenient to initialize in practical scenarios. In another word, the initialization limits the wide applications of existing BBox initialized trackers. (2) The initialized BBox may be not optimal for the representation of target object which may lead to ambiguity. As shown in Fig. 1 (a), the tracker may be confused to track the *bike* or *lower body* of the pedestrian. Similar views can also be found in [21, 43, 65, 78]. (3) Current BBox-based trackers may perform poorly when facing abrupt appearance variation of the target object, like *face/cloth changing* or *species variation* in Fig. 1 (b). Because the appearance feature initialized in the first frame and the object in the tracking procedure are vastly different. Only one sample initialized in the first frame is not enough to handle these challenging scenarios. These observations all inspire us to begin to think about *how can we conduct tracking in a more applicable and accurate way?*

*The first two authors contribute equally to this work. Yaowei Wang is the corresponding author. Email: {wangx03, shuxj, wangyu, tianyh}@pcl.ac.cn, zhangzhipeng2017@ia.ac.cn, jiangbo@ahu.edu.cn, fengwu@ustc.edu.cn.

¹https://github.com/wangxiao5791509/Single_Object_Tracking_Paper_List

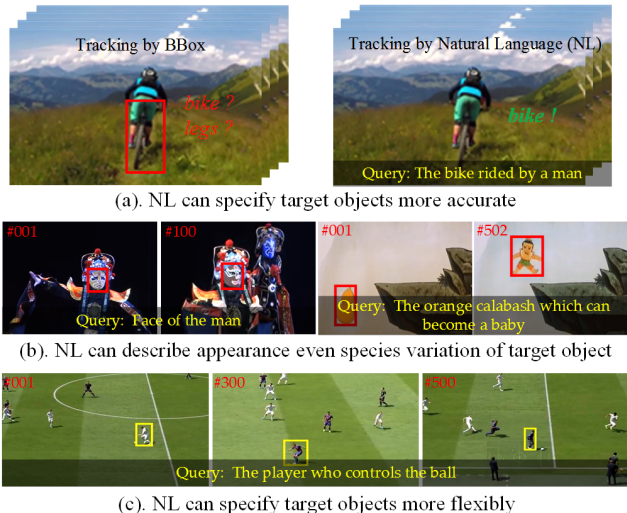


Figure 1. Comparison between the task of tracking-by-BBox and tracking-by-language. We can find that tracking-by-NL can specify target object more accurate and flexibly, and is also good at describing the appearance/species variation.

Recently, some researchers attempt to introduce the natural language description instead of the BBox for tracking [21, 43, 65, 78], termed tracking by natural language. This setting allows for a new type of human-machine interaction in object tracking. For example, it can enhance existing BBox based trackers by helping them against model drift, or simultaneous multiple-video tracking as noted in [43]. More importantly, natural language is more convenient and intuitive to express for human beings compared with BBox. It can provide a more precise expression of the target object from spatial location to high-level semantic information like *attributes*, *category*, *shape*, *properties*, and *structural relationship* with other objects, etc. This information will be beneficial to address the ambiguity issue of BBox and the vast appearance variation of the target object. Meanwhile, the language can also specify target objects more flexibly, for example, “*The player who controls the ball*” in Fig. 1 (c). The intelligent tracker should focus on target players even the ball passed to different persons, without having to re-initialize the target person like the standard setup of visual tracking. However, this research topic has received far less attention than standard target tracking. Only a few works [20, 21, 43, 65, 78] are developed and compared on tracking benchmark datasets specially designed for BBox based tracking. These benchmarks may fail to nicely reflect the true power of tracking-by-language, and this inspires us to design a new and large-scale benchmark for this task.

In this work, we collect a large-scale dataset that contains 2,000 video sequences, named TNL2K. These videos are collected from YouTube², surveillance cameras, and mo-

²<https://www.youtube.com/>

bile. For each video, we densely annotate the location information of the target object for each frame and one sentence in English for the whole video. Specifically, we describe the category, shape, attributes, properties, and spatial location of the target object which will provide rich fine-grained *appearance information* and high-level *semantic information* for tracking. We select 1,300 videos for training and the rest 700 videos for evaluation. Our videos also reflect two attributes for the tracking task, i.e., the adversarial samples and modality switch between RGB and thermal data. To provide a baseline method for other researchers to compare, we design a simple but strong algorithm based on an adaptive local-global-search scheme. Specifically, three kinds of baseline results are provided, i.e., Tracking-by-BBox, Tracking-by-Language, Tracking-by-BBox and Language.

The contributions of this paper can be summarized in the following three aspects:

- We propose the TNL2K dataset for the natural language-based tracking which consists of 2,000 video sequences. It aims at offering a dedicated platform for the development and assessment of natural language-based tracking algorithms.
- We propose a simple but strong baseline approach (termed AdaSwitcher) for future works to compare, which can switch between the local tracking algorithm and global grounding module adaptively.
- To provide extensive baselines for the comparison on TNL2K dataset, we also evaluate more than 40 representative BBox-based trackers and analyze their performance using different evaluation metrics.

2. Related Work

Tracking by Bounding Box The standard trackers begin their tracking procedure based on an initialized BBox in the first frame, including classification based [24, 31, 52, 53], Siamese network based [11, 12, 35, 64, 72], correlation filter based [13, 15, 27, 48], and regression-based [26]. Inspired by the success of neural networks on image classification, most of the recent trackers are developed based on deep learning. Specifically, the Siamese network based trackers achieve state-of-the-art performance on multiple tracking benchmarks. Previous Siamese trackers simply measure the similarity between the static target template with extracted proposals and treat the best-scored proposal as their tracking results. Recently, some researchers begin to collect the tracking results which can be used to dynamically update the target template and attain better results [75, 81]. In addition to learn powerful feature representation and conduct a local search for tracking, some trackers attempt to achieve robust tracking by global search [22, 29, 59, 63, 66, 74]. For more related works on standard visual tracking, please check the following survey papers [39, 42, 49, 56, 79].

Tracking by Natural Language Due to it is a new

rising topic, only a few algorithms are developed and the authors of [43] first validated the effectiveness of natural language for the tracking task by designing three modules (i.e. Lingual Specification Only; Lingual First, then Visual Specification; Lingual and Visual Specification). Wang [65] and Feng [21] also propose to use the language information to generate global proposals for tracking. Yang et al. propose the GTI [78] which decomposes the tracking problem into three sub-tasks, i.e., grounding, tracking and integration, and these modules operate simultaneously and predict the box sequence frame-by-frame. These methods are evaluated on datasets specifically designed for tracking-by-BBox which may fail to reflect the feature of tracking-by-language. To the best of our knowledge, there is still no public benchmark specifically dedicated to the tracking-by-language task. We believe our benchmark will greatly boost the researches on natural language related object tracking.

Benchmarks for Tracking Existing benchmarks for visual tracking can be concluded into two main categories according to whether contains training data. As shown in Table 1, previous benchmarks [32–34, 41, 41, 44, 70, 71] provide test videos only before deep trackers occurred. It is worthy to note that OTB-2013 [70] and OTB-2015 [71] are the first public benchmarks for visual tracking which contain 50 and 100 video sequences, respectively. In the deep learning era, several large scale tracking benchmarks are proposed for the training of deep trackers. For example, GOT-10k [28] contains 10,000 videos which can be categorized into 563 classes. TrackingNet [51] is a subset (31K sequences selected) of video object detection benchmark YT-BB [54] and the ground truth is manually labeled at 1 FPS. OxUvA [58] and LaSOT [19] are two long-term tracking benchmark which consists of 366 and 1400 video sequences respectively.

The aforementioned tracking benchmarks are all mainly designed for tracking by BBox, although the LaSOT indeed provides the language specification of the target object. However, they only describe the appearance of the target object but ignore the relative location which may limit the integration of natural language. In another word, their benchmark is suitable for natural language assisted tracking but is not for the task of language initialized tracking. Another issue of the existing benchmark is that these videos do not contain videos with significant appearance variations, such as clothing change for a pedestrian. This also limits the application of existing trackers in practical scenarios. Besides, these benchmarks also ignore the adversarial samples which limit the development of adversarial learning-based trackers [30, 45, 67, 73]. By contrast, our proposed TNL2K is specifically designed for tracking by natural language specification and contains multiple videos with significant appearance variation and adversarial samples. It also contains natural videos, animation videos, in-

fired videos, virtual game videos, which are suitable for the evaluation of domain adaptation of current trackers. We also provide baseline results of three kinds of settings which will be beneficial for future trackers to compare.

3. Tracking by Natural Language

3.1. TNL2K Dataset

Data Collection and Annotation The proposed TNL2K dataset contains 2,000 video sequences, and most of them are downloaded and clipped from YouTube, intelligent surveillance cameras, and mobile phones. We invite seven people for the annotation of these videos. Specifically, we annotate one sentence in English for each video and also one bounding box for each frame in this video. The left corner point (x_1, y_1) , width w and height h of the target’s bounding box are used as the ground truth, i.e., $[x_1, y_1, w, h]$. The annotated natural language description indicates the *spatial position*, *relative location with other objects*, *attribute*, *category* and *property* of target object in the first frame. We also annotate the *absent* label for each frame to enrich the information that is available for more accurate tracking. To construct a rich and heterogeneous benchmark, we also borrow some thermal videos from existing datasets [37, 46] and re-annotate the target object we want to track if necessary. Example sequences and annotations are illustrated in Fig. 2.

Attribute Definition Following popular tracking benchmarks [19, 28, 71], we also define multiple attributes of each video sequence for the evaluation under each challenging factors. As shown in Table 2, our proposed TNL2K dataset has the following 17 attributes: CM (Camera Motion), ROT (Rotate Of Target), DEF (DEFormation), FOC (Fully OCcluded), IV (Illumination Variation), OV (Out of View), POC (Partially OCcluded), VC (Viewpoint Change), SV (Scale Variation), BC (Background Clutter), MB (Motion Blur), ARC (Aspect Ratio Change), LR (Low Resolution), FM (Fast Motion), AS (Adversarial Sample), TC (Thermal Crossover), MS (Modality Switch). It is worthy to note that our dataset contains some thermal videos with challenging factors like TC (target object shares similar intensity with background), MS (the video contains both thermal and RGB images). To provide a good platform for the study of adversarial attack and defense of neural network for tracking, we also generate 100 videos contain adversarial samples as part of the testing subset using attack toolkit [30]. Therefore, these videos contain additional challenging factor, i.e., AS (influence of Adversarial Samples). It is worthy to note that the AS and MS are two new attributes for tracking community first proposed in this work. A more detailed distribution of each challenge is shown in Fig. 3 (c).

Statistical Analysis Our proposed TNL2K contains 663 English words and focuses on expressing the attributes, spa-

Table 1. Comparison of current datasets for object tracking. # denotes the number of corresponding item. Lang-A and Lang-I denote the dataset can be used for language assisted and initialized tracking task. SAV denotes the dataset contains many videos with significant appearance variation. Adv means the dataset contains adversarial samples (i.e., malicious attacks). DA is short for domain adaptation.

Datasets	#Videos	#Min	#Mean	#Max	#Total	#FR	#Attributes	Aim	Absent	Lang-A	Lang-I	SAV	Adv	DA
OTB50 [70]	51	71	578	3,872	29K	30 fps	11	Eval						
OTB100 [71]	100	71	590	3,872	59K	30 fps	11	Eval						
TC-128 [44]	128	71	429	3,872	55K	30 fps	11	Eval						
VOT-2017 [33]	60	41	356	1,500	21K	30 fps	-	Eval						
NUS-PRO [34]	365	146	371	5040	135K	30 fps	-	Eval						
UAV123 [50]	123	109	915	3085	113K	30 fps	12	Eval						
UAV20L [50]	20	1717	2934	5527	59K	30 fps	12	Eval						
NFS [32]	100	169	3830	20665	383K	240 fps	9	Eval						
TrackingNet [51]	30,643	-	480	-	14.43M	30 fps	15	Train/Eval						
OxUvA [58]	366	900	4260	37440	1.55M	30 fps	6	Train/Eval						
GOT-10k [28]	10,000	29	149	1,418	1.5M	10 fps	6	Train/Eval	✓					
LaSOT [19]	1,400	1000	2506	11397	3.52M	30 fps	14	Train/Eval	✓	✓				
TNL2K (Ours)	2,000	21	622	18488	1.24M	30 fps	17	Train/Eval	✓	✓	✓	✓	✓	✓

Table 2. Description of 17 attributes in our TNL2K dataset.

Attributes	Definition
01. CM	Abrupt motion of the camera
02. ROT	Target object rotates in the video
03. DEF	The target is deformable
04. FOC	Target is fully occluded
05. IV	Illumination variation
06. OV	The target completely leaves the video sequence
07. POC	Partially occluded
08. VC	Viewpoint change
09. SV	Scale variation
10. BC	Background clutter
11. MB	Motion blur
12. ARC	The ratio of bounding box aspect ratio is outside the range [0.5, 2]
13. LR	Low resolution
14. FM	The motion of the target is larger than the size of its bounding box
15. AS	Influence of adversarial samples
16. TC	Two targets with similar intensity cross each other
17. MS	Video contain both color and thermal images

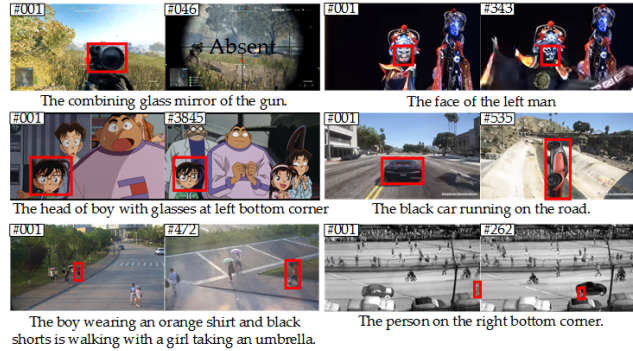


Figure 2. Example sequences and annotations in TNL2K dataset.

tial location of target objects, as shown in Fig. 3 (a). For the distribution of length of all videos, we can see from Fig. 3 (b) that the TNL2K contains [648, 479, 415, 139, 319] videos for the category of 1-300, 300-500, 500-800, 800-1000, and larger than 1000. More details, the number of these five segments for train and evaluation set are [488, 304, 258, 75, 175] and [160, 175, 157, 64, 144] respectively. We can find that our test set contains 144 long-term videos (larger than 1000 frames for each video) which will be suit-

able for the evaluation of long-term trackers. From Fig. 3 (c), we can find that our TNL2K contains many videos with challenging attributes like *background clutter*, *scale variation*, *view change*, *partially occlusion*, *out-of-view* and *rotate*. The videos with these challenging factors will provide a good platform for the evaluation of current trackers.

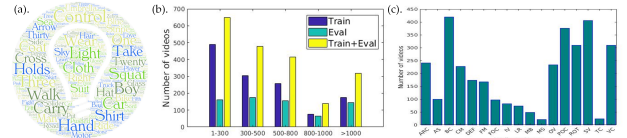


Figure 3. (a) Some words in our language description; (b, c) Distribution of sequences in each attribute and length in our TNL2K. Best viewed by zooming in.

3.2. Our Proposed Approach

In this paper, we propose the adaptive tracking and grounding switch framework for tracking by natural language specification, as shown in Fig. 4. We will first introduce the visual grounding and visual tracking module, then, we will focus on our AdaSwitcher module.

Visual Grounding Module In the tracking by natural language task, we need to first locate the target object only

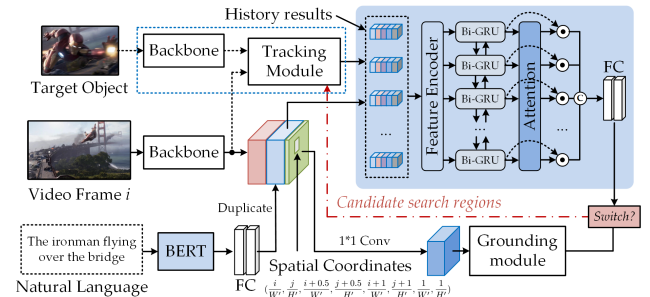


Figure 4. An overview of our proposed adaptive tracking and grounding switch framework. AdaSwitcher is highlighted in blue.

depends on the language description $S = [w_1, w_2, \dots, w_T]$. It is a standard visual grounding task and we follow the algorithm [77] proposed by Yang et al. due to its good performance and efficiency.

As shown in Fig. 4, the visual grounding module takes video frame and natural language description as input. We use the backbone CNN to obtain the deep feature representation of the i -th video frame F_i . For the natural language, we first embed the words into feature representations $E = [e_1, e_2, \dots, e_T]$ using a pre-trained BERT [16] which is a widely used word embedding model in natural language related tasks. Then, this feature is fed into two fully connected layers for further fine tuning. Following [77], we also duplicate this feature vector into feature maps and concatenate them with visual features of video frame. Another important information for visual grounding is the spatial coordinates encoding due to the spatial configurations are usually adopted to refer to target object. Therefore, the spatial feature for each position is also explicitly encoded in this work by following [77].

The visual feature maps of global frame, duplicated language feature, and the spatial coordinates are concatenated together and fed into convolutional layers with kernel size 1×1 for information fusion. The output feature map is then sent into the grounding module, which will output the predicted location of target object. We treat such visual grounding as a global search procedure for tracking by natural language, which plays an important role at the beginning of the video and when we need to re-detection the target object in tracking procedure. The integration of visual grounding and tracker SiamRPN++ [35] is termed Ours-I in Table 3. Besides, we also explore the target-aware attention (termed TANet) proposed in [65, 66], i.e., Ours-II in Table 3. The TANet takes the feature maps of target object and video image as input, and output corresponding global attention using de-convolutional network which can be used for search target object from global view. We refer the readers to check [65, 66] for further understanding of this module.

Visual Tracking Module Aforementioned visual grounding can help detect the target object at the beginning, however, only grounding is not enough for high performance tracking, since it is easily influenced by background clutter. In this work, we initialize a visual tracker for target object location in a local search manner based on the predicted bounding box from visual grounding in the first frame. The SiamRPN++ [35] is adopted in our experiments due to its good performance.

AdaSwitcher Module Given the visual grounding and visual tracking module, we can capture the target object from global and local view, respectively. One thorny issue that still exists is when we use visual grounding for global search (or visual tracking for local search). One intuitive

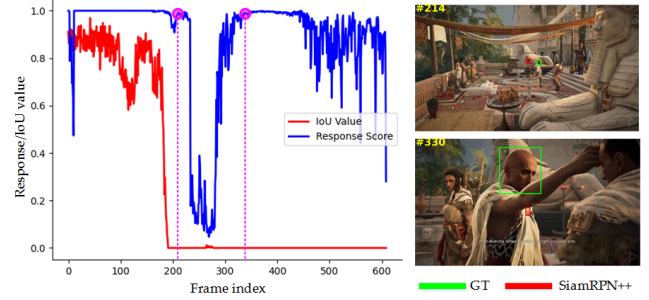


Figure 5. Illustration of current trackers with high response score but low IoU values (take the SiamRPN++ [35] as an example).

approach is to conduct such switch based on the confidence of tracker, however, the confidence score is not always reliable especially in the challenging scenarios. For example, as shown in Fig. 5, the confidence score is very high (larger than 0.9) in some frames, but the model actually locates wrong object. Inspired by anomaly detection (also called outlier detection) whose target is the identification of rare items, events or observations which raise suspicions by differing significantly from the majority of the data. In this work, we take the failure of visual tracking as a kind of anomaly detection and propose a novel AdaSwitcher module to detect such failure. Once the anomaly is detected (the prediction from AdaSwitcher is larger than a pre-defined threshold), we can switch the candidate search regions from visual tracking to visual grounding for more robust and accurate tracking.

In this paper, *confidence score* (1-D), *BBox* (4-D), *result image* ((30 * 30 * 3)-D), *response map* ((23 * 23)-D) and *language embedding* (512-D) are exploited in this work as the input of our AdaSwitcher. This information can be collected from visual tracker easily for each frame. And the historical information of past video frames can also contribute to current anomaly detection. Assume we use the history of past N frames, then, the dimension of these input are $N \times 1$, $N \times 4$, $N \times (23 \times 23)$, $N \times (30 \times 30 \times 3)$, and $N \times 512$, respectively. We use multiple parallel fully connected layers to encode this information and embed them into fixed feature vectors, specifically, we have $F = [F_s, F_b, F_{img}, F_{map}, F_{emb}]$, whose dimension are $N \times 10$, $N \times 10$, $N \times 512$, $N \times 512$, and $N \times 512$, respectively. Then, these features are concatenated and fed into a bi-directional GRUs [7] to learn the temporal information.

Inspired by the fact that various frames may contribute differently, we introduce attention mechanism to encode the inputs differently. The attention weights $\alpha_i (i = 1, \dots, N)$ can be obtained by the multilayer perceptron (MLP):

$$\{\alpha_1, \alpha_2, \dots, \alpha_N\} = MLP([F_s, F_b, F_{img}, F_{map}, F_{emb}]) \quad (1)$$

where $[\cdot]$ denotes concatenate operation. The attention weights $\alpha_i (i = 1, \dots, N)$ are stacked into feature vectors

$\hat{\alpha}_i (i = 1, \dots, N)$ which have same dimension with feature representation $F^i (i = 1, \dots, N)$ of each frame i . Therefore, the attended feature representations can be obtained by:

$$[\bar{F}^1, \bar{F}^2, \dots, \bar{F}^N] = [\hat{\alpha}_1 * F^1, \hat{\alpha}_2 * F^2, \dots, \hat{\alpha}_N * F^N] \quad (2)$$

After that, two fully connected layers are used to determine whether we should switch the candidate search regions from current tracking result to grounding result.

3.3. Implementation Details

Training Phase In our experiments, we directly use the pre-trained weights of baseline tracker for visual tracking. For the visual grounding module, we train it on the training subset of our TNL2K dataset which contains 1,300 video sequences for 40 epochs. The initial learning rate is $1e-4$, batchsize is 5. The YOLO loss function is used for this network by following [55, 77]. For the AdaSwitcher, we first collect the training data by running the baseline tracker on the training subset of our TNL2K dataset. In this process, we treat the video clips whose average IoU (Intersection over Union) score larger than 0.7 as the positive data, and less than 0.5 as the negative data. For the data with average IoU score range from 0.5 to 0.7, we directly discard them due to it may bring confusion to our model. Similar operations can also be found in [31]. The learning rate is $1e-5$, batchsize is 1, the Adagrad [18] is adopted as optimizer and trained for totally 30 epochs. We consider the switch between visual tracking and grounding as a binary classification problem, therefore, the BCE loss function is selected for the training of AdaSwitcher.

Inference Phase In this benchmark, three kinds of baseline methods are studied: 1). *Tracking by Natural Language only*: In this setting, only the natural language is provided for tracking, we need to first locate the target object using visual grounding module. Then, we can conduct adaptive tracking (SiamRPN++ [35] used in this setting) and grounding for high performance object localization. 2). *Tracking by Natural Language and BBox*: We take the natural language as an external modality and conduct robust tracking based on both language and BBox. SiamRPN++ [35] and TANet [66] are used in this setting. 3). *Tracking by BBox only*: To construct a comprehensive benchmark, we also provide baseline results for tracking by BBox only, i.e., the standard setting of visual object tracking. All the evaluated trackers can be found in our supplementary materials.

4. Experiments

4.1. Datasets and Evaluation Protocols

In our experiments, the OTB-Lang [43, 71], LaSOT [19] and our proposed TNL2K dataset are used for the evaluation. The OTB-lang contains 99 videos released from [71], then, the natural language specification is provided by Li et

al. [43]. The LaSOT is a recently released long-term tracking dataset that provides both bounding box and natural language annotations. The test subset of LaSOT contains 280 video sequences.

Two popular metrics are adopted for the evaluation of tracking performance, including **Precision Plot** and **Success Plot**. Specifically, Precision Plot illustrates the percentage of frames where the center location error between the object location and ground truth is smaller than a pre-defined threshold (20-pixels threshold is usually adopted). Success Plot demonstrates the percentage of frames the IoU of the predicted and the ground truth bounding boxes is higher than a given ratio. The evaluation toolkit of this paper can be found at: https://github.com/wangxiao5791509/TNL2K_evaluation_toolkit.

4.2. Benchmark Results

Results of Tracking by Natural Language Only As shown in Table 3, Li et al. [43] attain 0.29|0.25 on the OTB-Lang dataset, while Feng et al. achieve 0.56|0.54 and 0.78|0.54 in [20] and [21] respectively. When we take the result of visual grounding in the first frame as the initialized bbox of visual tracker SiamRPN++, we achieve 0.24|0.19 on the OTB-Lang dataset. On the LaSOT and TNL2K dataset, we attain 0.49|0.51 and 0.06|0.11|0.11 respectively. We can find that our method is comparable with Li et al. on the OTB-Lang dataset. On the larger dataset LaSOT, we attain better results than Feng et al. [20]. These experimental results demonstrate that our baseline method can also achieve good performance on existing LaSOT and our proposed TNL2K dataset.

Results of Tracking by Bounding Box Only This setting is most widely used in existing tracking algorithms, and we provide the results of 43 representative trackers from 2015 to 2021, as shown in Fig. 6. These trackers contain *Classification*-based, *SiameseNet*-based, *Correlation filter*-based, *Reinforcement learning*-based, *Long-term*-based and *Other* trackers. More detailed introductions on these trackers can be found in our supplementary materials due to the limited space in this paper. From Fig. 6, we can find that SiamRCNN [59] achieves the best performance on our benchmark dataset, i.e., 0.528|0.523 on the precision/success plot respectively. Other trackers also attain good performance such as LTMU [9], KYS [22], TACT [6], due to the use of joint local and global search scheme. These experiments fully demonstrate the importance of joint local and global search for visual tracking. We also find that Siamese network based trackers usually achieve better results than other trackers like multi-domain based trackers [31, 52, 53], regression based trackers [26], and correlation filter based trackers [2, 10, 27]. We also notice that GlobalTrack [29] which employs global search

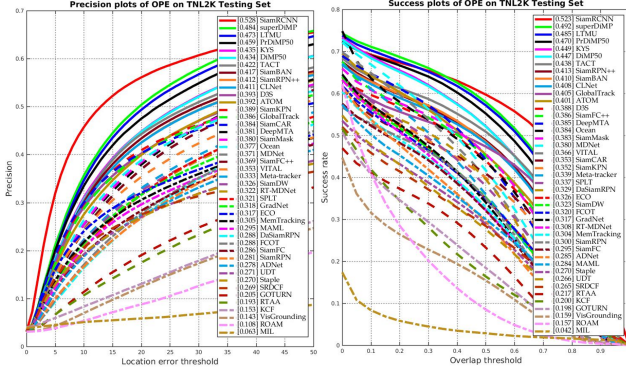


Figure 6. Benchmark results of tracking-by-BBox on TNL2K dataset. Best viewed by zooming in.

scheme only, achieves comparable performance with local search trackers [52, 83], but worse than state-of-the-art. This may demonstrate that only global search is not enough for robust tracking. Overall, the aforementioned observations demonstrate that the structure information mining of global scene, and offline learning indeed contribute to the high-performance visual tracking.

Results of Tracking by Joint Language and BBox

As shown in Table 3, there are five trackers designed for this setting [20, 21, 43, 65, 78]. Specifically, Li et al. [43] achieve 0.72|0.55, while Feng et al. [20, 21] attain 0.73|0.67, 0.79|0.61 on the OTB-Lang dataset respectively. GTI [78] combine SiamRPN++ and visual grounding module, and achieves 0.73|0.58, 0.47|0.47 on OTB-Lang and LaSOT dataset. In contrast, we can achieve 0.88|0.68 on the OTB-Lang, 0.55|0.51 on the LaSOT, 0.42|0.50|0.42 on the TNL2K (Ours-II in Table 3), which are significantly better than GTI [78], Wang et al. [65] and Feng et al. [20]. All the experiments on three benchmark datasets validate the effectiveness and advantages of our tracker. Visualization of related tracking results can be found in Fig. 12.

4.3. Ablation Study

In this section, we first analyse the effectiveness of main components in our model. Then, we focus on validating the contributions of each input for AdaSwitcher. Finally, we give the parameter analysis, and attribute analysis.

Effectiveness of AdaSwitcher As shown in Table 5, the baseline tracker SiamRPN++ [35] (AlexNet version) achieves 0.344/0.353 on the precision and success plot, respectively. When integrated with the AdaSwitcher module, the performance can be improved to 0.355/0.370. This result is also better than naive fused method (i.e. 0.347/0.362), which fully demonstrates the effectiveness of our adaptive switch mechanism for robust tracking.

Effectiveness of Frame Attention Due to different frames may contribute differently to our AdaSwitcher, we introduce the frame attention mechanism to achieve this

Table 3. Tracking results on the OTB-Lang, LaSOT, and TNL2K dataset. [Prec.|Norm. Prec. |Succ. Plot] are reported respectively.

Algorithm	Initialize	OTB-Lang	LaSOT	TNL2K
SiamFC [3]	BBox	-	0.40 0.34	0.29 0.35 0.30
MDNet [52]	BBox	-	0.46 0.40	0.37 0.46 0.38
VITAL [57]	BBox	-	0.45 0.39	0.35 0.44 0.37
GradNet [38]	BBox	-	0.35 0.37	0.32 0.40 0.32
ATOM [11]	BBox	-	0.51 0.51	0.39 0.47 0.40
SiamDW [82]	BBox	-	- 0.38	0.33 0.41 0.32
SiamRPN++ [35]	BBox	-	0.50 0.45	0.41 0.48 0.41
GlobalTrack [29]	BBox	-	0.53 0.52	0.39 0.46 0.41
SiamBAN [5]	BBox	-	0.60 0.51	0.42 0.49 0.41
Ocean [83]	BBox	-	0.57 0.56	0.38 0.45 0.38
Li et al. [43]	NL	0.29 0.25	-	-
Li et al. [43]	NL+BBox	0.72 0.55	-	-
Feng et al. [21]	NL	0.56 0.54	-	-
Feng et al. [21]	NL+BBox	0.73 0.67	0.56 0.50	0.27 0.34 0.25
Feng et al. [20]	NL	0.78 0.54	0.28 0.28	-
Feng et al. [20]	NL+BBox	0.79 0.61	0.35 0.35	0.27 0.33 0.25
Wang et al. [65]	NL+BBox	0.89 0.65	0.30 0.27	-
GTI [78]	NL+BBox	0.73 0.58	0.47 0.47	-
Ours-I	NL	0.24 0.19	0.49 0.51	0.06 0.11 0.11
Ours-II	NL+BBox	0.88 0.68	0.55 0.51	0.42 0.50 0.42

goal. As shown in Table 4, with the help of frame attention, the tracking results can be improved from 0.353/0.369 to 0.355/0.370. This fully demonstrates the important role of frame attention in our proposed framework.

Effectiveness of Spatial Coordinates In our visual grounding module, the spatial coordinates are introduced to further improve the final results. As shown in Table 4, our grounding module achieves 0.143/0.159 and 0.103/0.124, respectively, with and without the help of spatial coordinates. This result validates the important role of spatial coordinates for visual grounding.

Table 4. Component analysis of our proposed tracking algorithm. AS is short for AdaSwitcher, FA denotes frame attention in AdaSwitcher, SC is spatial coordinates used in visual grounding. Naive denotes switch method based on response score only.

Track	Ground	SC	TANet	Naive	AS	FA	Results
✓							0.344 0.353
	✓						0.103 0.124
		✓					0.143 0.159
✓			✓	✓			0.347 0.362
✓			✓		✓	✓	0.355 0.370
✓			✓		✓		0.353 0.369

Analysis on History Information Our AdaSwitcher takes multiple inputs for the final decision, in this section, we analyze their contributions by comparing corresponding results in Table 5. Specifically speaking, when the BBox is discarded, we find that the performance is dropped from 0.355/0.370 to 0.350/0.365, this demonstrates that the geometric information of predicted BBox is an important clue for our tracking. Similarly, we attain worse tracking results when the resulting image (i.e. ResImg) is ignored, the results drop from 0.355/0.370 to 0.345/0.362. When all these modules removed, it attains 0.344/0.353 only on the precision and success plot. This demonstrates that this informa-

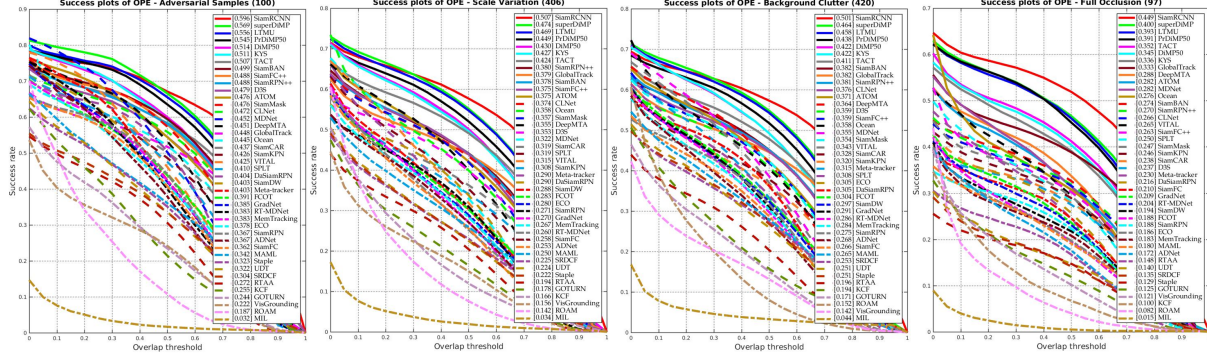


Figure 7. Tracking results under partial attributes of TNL2K dataset. Best viewed by zooming in.

tion are very important for the anomaly (or failure) detection in tracking procedure.

Table 5. Component analysis of history information.

BBox	Score	ResMap	ResImg	Lang	Results
✓	✓	✓	✓	✓	0.355 0.370
✗	✓	✓	✓	✓	0.350 0.365
✓	✗	✓	✓	✓	0.352 0.368
✓	✓	✗	✓	✓	0.352 0.368
✓	✓	✓	✗	✓	0.345 0.362
✓	✓	✓	✓	✗	0.352 0.368
✗	✗	✗	✗	✗	0.344 0.353

Parameter Analysis We report the tracking results with different switch thresholds in Table 6. We can find that the performance is better when switch threshold is set as 0.7.

Table 6. Results with different switch threshold.

Parameter	0.5	0.6	0.7	0.8	0.9	1.0	1.2
Prec. Plot	0.350	0.351	0.355	0.353	0.349	0.352	0.272
Succ. Plot	0.367	0.367	0.370	0.369	0.368	0.368	0.301



Figure 8. Visualization of tracking results on TNL2K dataset.

Attribute Analysis Evaluation under each challenging factors is one of the most important metrics in visual tracking community. In this benchmark, we also report results of evaluated trackers under all the defined 17 attributes.

However, due to limited space in this paper, we select 4 attributes, i.e., *Adversarial Samples*, *Scale Variation*, *Background Clutter*, and *Full Occlusion*, to demonstrate the ability of resistance of these trackers to these challenges. As shown in Fig. 7, we can find that SiamRCNN [59] achieves the best performance which are much better than the second and third ones, i.e., DiMP [4] and LTMU [9] respectively. Interestingly, it is easy to find that the RTAA [30] which is designed for adversarial attack achieves worse results on the challenging factor Adversarial Samples, even compared with their baseline DaSiamRPN [84]. This demonstrates that the detection of adversarial samples is important for high performance tracking. More experimental results on the attribute analysis can be found in our supplementary materials.

5. Conclusion and Future Works

In this paper, we revisit the tracking by natural language, and propose a large-scale benchmark for this task. Specially, a large-scale dataset that contains 2,000 video sequences is proposed, named TNL2K. This dataset is densely annotated with bounding box and natural language description of target object. To construct a sound benchmark, we propose an adaptive switch based tracking algorithm as the baseline approach, i.e., the AdaSwicher, and also test current trackers according to following settings: tracking by natural language only, tracking by bbox, and tracking by joint bbox and language. We believe our benchmark will be greatly boost related researches on the natural language guided tracking. In our future works, we will consider to further extend this benchmark by introducing more videos and baseline trackers. Besides, we will focus on improving the visual grounding module to achieve high performance language initialized tracking.

Acknowledgement: This work is jointly supported by Key-Area Research and Development Program of Guangdong Province 2019B010155002, Postdoctoral Innovative Talent Support Program BX20200174, China Postdoctoral Science Foundation Funded Project 2020M682828, National Natural Science Foundation of China (61976002, 61825101), National Key Research and Development Program of China 2020AAA0106800.

References

- [1] Boris Babenko, Ming-Hsuan Yang, and Serge Belongie. Visual tracking with online multiple instance learning. In *2009 IEEE Conference on computer vision and Pattern Recognition*, pages 983–990. IEEE, 2009.
- [2] Luca Bertinetto, Jack Valmadre, Stuart Golodetz, Ondrej Miksik, and Philip HS Torr. Staple: Complementary learners for real-time tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1401–1409, 2016.
- [3] Luca Bertinetto, Jack Valmadre, Joao F Henriques, Andrea Vedaldi, and Philip HS Torr. Fully-convolutional siamese networks for object tracking. In *European conference on computer vision*, pages 850–865. Springer, 2016.
- [4] Goutam Bhat, Martin Danelljan, Luc Van Gool, and Radu Timofte. Learning discriminative model prediction for tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6182–6191, 2019.
- [5] Zedu Chen, Bineng Zhong, Guorong Li, Shengping Zhang, and Rongrong Ji. Siamese box adaptive network for visual tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6668–6677, 2020.
- [6] Janghoon Choi, Junseok Kwon, and Kyoung Mu Lee. Visual tracking by tridentalign and context embedding. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [7] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [8] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu. Fully convolutional online tracking. *arXiv preprint arXiv:2004.07109*, 2020.
- [9] Kenan Dai, Yunhua Zhang, Dong Wang, Jianhua Li, Huchuan Lu, and Xiaoyun Yang. High-performance long-term tracking with meta-updater. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6298–6307, 2020.
- [10] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. Eco: Efficient convolution operators for tracking. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [11] Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, and Michael Felsberg. Atom: Accurate tracking by overlap maximization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4660–4669, 2019.
- [12] Martin Danelljan, Luc Van Gool, and Radu Timofte. Probabilistic regression for visual tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7183–7192, 2020.
- [13] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4310–4318, 2015.
- [14] Martin Danelljan, Gustav Hager, Fahad Shahbaz Khan, and Michael Felsberg. Learning spatially regularized correlation filters for visual tracking. In *Proceedings of the IEEE international conference on computer vision*, pages 4310–4318, 2015.
- [15] Martin Danelljan, Andreas Robinson, Fahad Shahbaz Khan, and Michael Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *Proceedings of European Conference on Computer Vision*, 2016.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, 2019.
- [17] Xingping Dong, Jianbing Shen, Ling Shao, and Fatih Porikli. Clnet: A compact latent network for fast adjusting siamese trackers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [18] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7):2121–2159, 2011.
- [19] Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5374–5383, 2019.
- [20] Qi Feng, Vitaly Ablavsky, Qinxun Bai, Guorong Li, and Stan Sclaroff. Real-time visual object tracking with natural language description. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 700–709, 2020.
- [21] Qi Feng, Vitaly Ablavsky, Qinxun Bai, and Stan Sclaroff. Robust visual object tracking with natural language region proposal network. *arXiv preprint arXiv:1912.02048*, 2019.
- [22] Bhat Goutam, Danelljan Martin, Van Gool Luc, and Timofte Radu. Know your surroundings: Exploiting scene information for object tracking. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [23] Dongyan Guo, Jun Wang, Ying Cui, Zhenhua Wang, and Shengyong Chen. Siamcar: Siamese fully convolutional classification and regression for visual tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6269–6277, 2020.
- [24] Bohyung Han, Jack Sim, and Hartwig Adam. Branchout: Regularization for online ensemble tracking with convolutional neural networks. In *Proceedings of IEEE International Conference on Computer Vision*, pages 2217–2224, 2017.
- [25] Sam Hare, Stuart Golodetz, Amir Saffari, Vibhav Vineet, Ming-Ming Cheng, Stephen L Hicks, and Philip HS Torr. Struck: Structured output tracking with kernels. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):2096–2109, 2016.
- [26] David Held, Sebastian Thrun, and Silvio Savarese. Learning to track at 100 fps with deep regression networks. In *European Conference on Computer Vision*, pages 749–765. Springer, 2016.

- [27] João F. Henriques, Rui Caseiro, Pedro Martins, and Jorge Batista. High-speed tracking with kernelized correlation filters. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 37(3):583–596, 2015.
- [28] Lianghai Huang, Xin Zhao, and Kaiqi Huang. Got-10k: A large high-diversity benchmark for generic object tracking in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [29] Lianghai Huang, Xin Zhao, and Kaiqi Huang. Globaltrack: A simple and strong baseline for long-term tracking. *AAAI*, 2020.
- [30] Shuai Jia, Chao Ma, Yibing Song, and Xiaokang Yang. Robust tracking against adversarial attacks. In *European Conference on Computer Vision*, pages 69–84. Springer, 2020.
- [31] Ilchae Jung, Jeany Son, Mooyeol Baek, and Bohyung Han. Real-time mdnet. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [32] Hamed Kiani Galoogahi, Ashton Fagg, Chen Huang, Deva Ramanan, and Simon Lucey. Need for speed: A benchmark for higher frame rate object tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1125–1134, 2017.
- [33] Matej Kristan, Jiri Matas, Aleš Leonardis, Tomas Vojir, Roman Pflugfelder, Gustavo Fernandez, Georg Nebehay, Fatih Porikli, and Luka Čehovin. A novel performance evaluation methodology for single-target trackers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(11):2137–2155, Nov 2016.
- [34] A Li, M Lin, Y Wu, MH Yang, and S Yan. NUS-PRO: A New Visual Tracking Challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(2):335–349, 2016.
- [35] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. Siamrpn++: Evolution of siamese visual tracking with very deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4282–4291, 2019.
- [36] Bo Li, Junjie Yan, Wei Wu, Zheng Zhu, and Xiaolin Hu. High performance visual tracking with siamese region proposal network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8971–8980, 2018.
- [37] Chenglong Li, Xinyan Liang, Yijuan Lu, Nan Zhao, and Jin Tang. Rgb-t object tracking: benchmark and baseline. *Pattern Recognition*, 96:106977, 2019.
- [38] Peixia Li, Boyu Chen, Wanli Ouyang, Dong Wang, Xiaoyun Yang, and Huchuan Lu. Gradnet: Gradient-guided network for visual object tracking. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [39] Peixia Li, Dong Wang, Lijun Wang, and Huchuan Lu. Deep visual tracking: Review and experimental comparison. *Pattern Recognition*, 76:323–338, 2018.
- [40] Qiang Li, Zekui Qin, Wenbo Zhang, and Wen Zheng. Siamese keypoint prediction network for visual object tracking. *arXiv preprint arXiv:2006.04078*, 2020.
- [41] Siyi Li and Dit-Yan Yeung. Visual object tracking for unmanned aerial vehicles: A benchmark and new motion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [42] Xi Li, Weiming Hu, Chunhua Shen, Zhongfei Zhang, Anthony Dick, and Anton Van Den Hengel. A survey of appearance models in visual object tracking. *ACM transactions on Intelligent Systems and Technology (TIST)*, 4(4):58, 2013.
- [43] Zhenyang Li, Ran Tao, Efstratios Gavves, Cees GM Snoek, and Arnold WM Smeulders. Tracking by natural language specification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6495–6503, 2017.
- [44] Pengpeng Liang, Erik Blasch, and Haibin Ling. Encoding color information for visual tracking: Algorithms and benchmark. *IEEE Transactions on Image Processing*, 24(12):5630–5644, 2015.
- [45] Siyuan Liang, Xingxing Wei, Siyuan Yao, and Xiaochun Cao. Efficient adversarial attacks for visual object tracking. *arXiv preprint arXiv:2008.00217*, 2020.
- [46] Qiao Liu, Zhenyu He, Xin Li, and Yuan Zheng. Ptb-tir: A thermal infrared pedestrian tracking benchmark. *IEEE Transactions on Multimedia*, 22(3):666–675, 2019.
- [47] Alan Lukezic, Jiri Matas, and Matej Kristan. D3s-a discriminative single shot segmentation tracker. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7133–7142, 2020.
- [48] Alan Lukezic, Tomas Vojir, Luka Čehovin Zajc, Jiri Matas, and Matej Kristan. Discriminative correlation filter with channel and spatial reliability. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4847–4856. IEEE, 2017.
- [49] Seyed Mojtaba Marvasti-Zadeh, Li Cheng, Hossein Ghanei-Yakhdan, and Shohreh Kasaei. Deep learning for visual tracking: A comprehensive survey. *arXiv preprint arXiv:1912.00535*, 2019.
- [50] Matthias Mueller, Neil Smith, and Bernard Ghanem. A benchmark and simulator for uav tracking. In *European conference on computer vision*, pages 445–461. Springer, 2016.
- [51] Matthias Muller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 300–317, 2018.
- [52] Hyeonseob Nam and Bohyung Han. Learning multi-domain convolutional neural networks for visual tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4293–4302, 2016.
- [53] Eunbyung Park and Alexander C. Berg. Meta-tracker: Fast and robust online adaptation for visual object trackers. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [54] Esteban Real, Jonathon Shlens, Stefano Mazzocchi, Xin Pan, and Vincent Vanhoucke. Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5296–5305, 2017.
- [55] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.

- [56] A. W. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah. Visual tracking: An experimental survey. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 36(7):1442–68, 2014.
- [57] Yibing Song, Chao Ma, Xiaohe Wu, Lijun Gong, Linchao Bao, Wangmeng Zuo, Chunhua Shen, Rynson W.H. Lau, and Ming-Hsuan Yang. Vital: Visual tracking via adversarial learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [58] Jack Valmadre, Luca Bertinetto, Joao F Henriques, Ran Tao, Andrea Vedaldi, Arnold WM Smeulders, Philip HS Torr, and Efstratios Gavves. Long-term tracking in the wild: A benchmark. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 670–685, 2018.
- [59] Paul Voigtlaender, Jonathon Luiten, Philip HS Torr, and Bastian Leibe. Siam r-cnn: Visual tracking by re-detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6578–6588, 2020.
- [60] Guangting Wang, Chong Luo, Xiaoyan Sun, Zhiwei Xiong, and Wenjun Zeng. Tracking by instance detection: A meta-learning approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6288–6297, 2020.
- [61] Ning Wang, Yibing Song, Chao Ma, Wengang Zhou, Wei Liu, and Houqiang Li. Unsupervised deep tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1308–1317, 2019.
- [62] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip HS Torr. Fast online object tracking and segmentation: A unifying approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1328–1338, 2019.
- [63] Xiao Wang, Zhe Chen, Jin Tang, Bin Luo, Yaowei Wang, Yonghong Tian, and Feng Wu. Dynamic attention guided multi-trajectory analysis for single object tracking. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [64] Xiao Wang, Chenglong Li, Bin Luo, and Jin Tang. Sint++: Robust visual tracking via adversarial positive instance generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4864–4873, 2018.
- [65] Xiao Wang, Chenglong Li, Rui Yang, Tianzhu Zhang, Jin Tang, and Bin Luo. Describe and attend to track: Learning natural language guided structural representation and visual attention for object tracking. *arXiv preprint arXiv:1811.10014*, 2018.
- [66] Xiao Wang, Rui Yang, Tao Sun, and Bin Luo. Learning target-aware attention for robust tracking with conditional adversarial network. In *30TH British Machine Vision Conference (BMVC)*, page 131, 2019.
- [67] Rey Reza Wiyatno and Anqi Xu. Physical adversarial textures that fool visual object tracking. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4822–4831, 2019.
- [68] Ancong Wu, Wei-Shi Zheng, Shaogang Gong, and Jianhuang Lai. Rgb-ir person re-identification by cross-modality similarity preservation. *International journal of computer vision*, pages 1–21, 2020.
- [69] Ancong Wu, Wei-Shi Zheng, Hong-Xing Yu, Shaogang Gong, and Jianhuang Lai. Rgb-infrared cross-modality person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 5380–5389, 2017.
- [70] Yi Wu, Jongwoo Lim, and Ming Hsuan Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 37(9):1834, 2015.
- [71] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Object tracking benchmark. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1834–1848, 2015.
- [72] Yinda Xu, Zeyu Wang, Zuoxin Li, Ye Yuan, and Gang Yu. Siamfc++: Towards robust and accurate visual tracking with target estimation guidelines. In *AAAI*, pages 12549–12556, 2020.
- [73] Bin Yan, Dong Wang, Huchuan Lu, and Xiaoyun Yang. Cooling-shrinking attack: Blinding the tracker with imperceptible noises. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 990–999, 2020.
- [74] Bin Yan, Haojie Zhao, Dong Wang, Huchuan Lu, and Xiaoyun Yang. ‘skimming-perusal’ tracking: A framework for real-time and robust long-term tracking. *ICCV*, 2019.
- [75] Tianyu Yang and Antoni B. Chan. Learning dynamic memory networks for object tracking. In *The European Conference on Computer Vision (ECCV)*, September 2018.
- [76] Tianyu Yang, Pengfei Xu, Runbo Hu, Hua Chai, and Antoni B Chan. Roam: Recurrently optimizing tracking model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6718–6727, 2020.
- [77] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. A fast and accurate one-stage approach to visual grounding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4683–4693, 2019.
- [78] Zhengyuan Yang, Tushar Kumar, Tianlang Chen, Jingsong Su, and Jiebo Luo. Grounding-tracking-integration. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [79] Alper Yilmaz, Omar Javed, and Mubarak Shah. Object tracking: A survey. *ACM Comput. Surv.*, 38(4):13–es, dec 2006.
- [80] Sangdoo Yun, Jongwon Choi, Youngjoon Yoo, Kimin Yun, and Young Choi Jin. Action-decision networks for visual tracking with deep reinforcement learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1349–1358, 2017.
- [81] Lichao Zhang, Abel Gonzalez-Garcia, Joost van de Weijer, Martin Danelljan, and Fahad Shahbaz Khan. Learning the model update for siamese trackers. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4010–4019, 2019.
- [82] Zhipeng Zhang and Houwen Peng. Deeper and wider siamese networks for real-time visual tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4591–4600, 2019.
- [83] Zhipeng Zhang, Houwen Peng, Jianlong Fu, Bing Li, and Weiming Hu. Ocean: Object-aware anchor-free tracking. In

European Conference on Computer Vision, volume 12366, pages 771–787. Springer, 2020.

- [84] Zheng Zhu, Qiang Wang, Bo Li, Wei Wu, Junjie Yan, and Weiming Hu. Distractor-aware siamese networks for visual object tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 101–117, 2018.

A. The TNL2K Benchmark

A.1. Motivation and Protocols

Motivation: Directly extending existing datasets like GOT-10k [28] is an intuitive and good idea for this task, but GOT-10k contains few videos with special properties as mentioned in Fig. 1 in our paper. Also, its videos are all short-term which can't reflect performance gain of re-detection with language. As for LaSOT [19], many of its language annotations can not point out target object clearly, as shown in Fig. 9. Thus, LaSOT is not suitable for tracking-by-language only. Similar views can also be found in GTI [78]. Therefore, we build the TNL2K (from video collection, dense bbox and language annotation, to diverse baseline construction) to better reflect the characteristics (see below) of tracking by natural language. The target of this work is not to construct the largest tracking dataset, but to build the first benchmark specifically designed for tracking-by-language task. Compared with GOT-10k and LaSOT, the data collection of TNL2K is a compromise between length and quantity.

Protocol: When collecting the videos, we attempt to search the target object is *severely occluded in the first frame*, with *significant appearance variation* (e.g., cloth changing for human), *can only be located with reasoning*, which correspond to Fig. 1 in our paper. Also, we collect videos from other thermal tracking datasets and annotate language descriptions only to check the robustness to certain challenging factors like domain adaptation, modality switch, etc.

A.2. Why add Attribute Modality Switch (MS) ?

In the proposed TNL2K dataset, we design a new attribute termed Modality Switch (MS) for object tracking. This is mainly motivated by the fact that the RGB cameras work well in the daytime but nearly ineffective at night, meanwhile, the thermal cameras work well in the night time. If we track a target for an extremely long-term (e.g., several days or weeks), collaboration between RGB and thermal cameras are needed. Therefore, the connections between the two modalities need to be set up. Similar views can be found in cross-modality person re-identification [68, 69]. There are still no works on object tracking try to build such connections and they usually study these two cameras separately (i.e., RGB tracking [19, 64, 70], Thermal Tracking [46]) or in an integrated approach (i.e., RGB-T tracking [37]). In this work, we propose the modality switch and attempt to encourage researches on such cross-modality object tracking.

A.3. Highlights of TNL2K Dataset

Generally speaking, our proposed benchmark TNL2K have the following features as shown in Table 1:

- **TNL2K is the first benchmark specifically designed for tracking-by-natural language.** Different from regular tracking benchmarks like OTB, GOT10k, and TrackingNet, we provide both language annotation and dense bounding box annotation for each video sequence which will be a good platform for natural language-related tracking. Different from the recently released long-term tracking dataset LaSOT which also provides language annotation, their annotation only describes the attribute of target object, but ignores

the spatial position. Therefore, this benchmark can be only used for the task of *tracking by joint language and bbox*. Our language annotations not only embody the attribute, category, shape, properties, and structural relationship with other objects, therefore, our dataset can also be used for the task of *tracking by natural language only*. Some video sequences and corresponding annotations are provided in Figure 9 to give an intuitive understanding of the difference between our TNL2K and LaSOT.

- **TNL2K is the first benchmark to provides videos with actively introduced adversarial samples** which will be beneficial for the development of adversarial training for tracking.
- **TNL2K is the first benchmark to provides videos with significant appearance variation**, such as *cloth/face changing*. We believe our benchmark will greatly boost related research on abrupt appearance variation based tracking.
- **TNL2K provides a heterogeneous dataset** that contains RGB video, Thermal video³, Cartoon, and Synthetic data (i.e., videos from games). It can be used for the study of domain adaptation, e.g., train the tracker on RGB data and test it on Thermal videos.
- **TNL2K provides three kinds of baseline methods for future works to compare**, including Tracking-by-BBox, Tracking-by-Language, Tracking-by-Joint-BBox-Language.

B. The Proposed Method

B.1. YOLO Loss and BCE Loss Functions

In the training phase, we use the YOLO loss function for the optimization of the visual grounding module by following [77]. This loss is first proposed in YOLOv3 [55] which attempt to predict the five quantities of each anchor box by shifting its *center*, *width*, *height*, and the *confidence* on this shifted box. To better use it for visual grounding, the following two changes are modified by Yang et al.: 1). recalibrate its anchor boxes; 2). change its sigmoid layer to a softmax function. Due to the object detection is designed for output multiple locations, while visual grounding only needs to predict one bbox which best fit the language description. Therefore, the sigmoid function in YOLOv3 is replaced by softmax function. The cross-entropy is used for the measurement of confidence scores, and the regions with maximum IoU with ground truth are labeled as 1, other regions are set as 0. More details can be found in [55, 77]. For the training of TANet, we adopt Binary Cross-Entropy (BCE) loss to measure the distance between the ground truth mask and the prediction.

B.2. Details of Evaluated Trackers

In this section, we provide the details of evaluated BBox-based trackers on our TNL2K dataset. As shown in Table 7, the publication, feature representation, update or not, need pre-train or not, search scheme, tracking efficiency, and results (Precision Plot and Success Plot) on the TNL2K are all reported. These tracking algorithms are ranked according to the results.

³There are 518 videos totally borrowed from existing RGB-T dataset [37] and infrared tracking dataset [46].

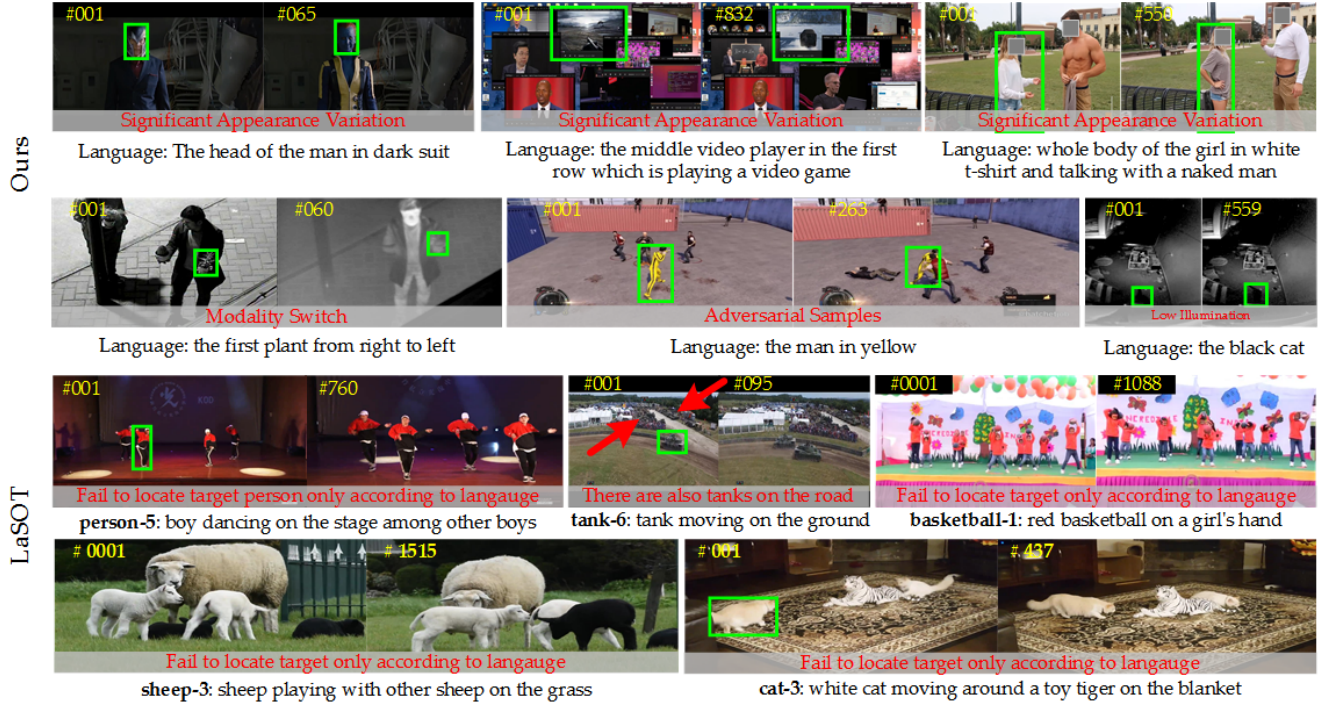


Figure 9. Comparison between our proposed TNL2K dataset and existing LaSOT dataset. Best viewed by zooming in.

B.3. Introduction to TANet

Inspired by [63, 65, 66], we introduce the TANet for the global search to replace the Grounding module [77] in the setting of *tracking-by-joint language and BBox*, termed Ours-II. Generally speaking, the TANet is inspired by semantic segmentation, which takes the target object and video frames as input and output an attention map using a decoder network. The estimated attention maps can highlight the possible search regions from a global view. Therefore, it can be seen as a kind of global search scheme and can be integrated with the baseline tracker and our proposed AdaSwitcher module for robust and accurate tracking. Our experimental results also demonstrate that we can attain good performance on three used datasets, i.e., the OTB-Lang [43], LaSOT [19], and TNL2K. This will be a strong baseline method for future works to compare on the language guided visual tracking. The implementation of our all networks will be released for other researchers to follow.

C. Experimental Results

C.1. Attribute Analysis

As shown in Figure 10, we provide experimental results of all the defined 17 attributes of our TNL2K dataset. Generally speaking, we can find that the SiamRCNN [59] achieves the best performance on most of the attributes, like *Scale Variation*, *Rotation*, *Background Clutter*, *Partial Occlusion*, *Adversarial Samples*, *Deformation*, *Fast Motion*, *Out-of-view*, *Motion Blur*, *Aspect Ratio Change*, *Illumination Variation*, *Camera Motion*, and *Viewpoint Change*. Meanwhile, the SuperDiMP [4], LTMU [9],

PrDiMP [12] and KYS [22] also attains good performance on these attributes, and the KYS also achieves top-1 results on the *Low Resolution*. These results all demonstrate the strong performance of Siamese network based trackers with the help of pre-training and joint local and global search scheme. Interestingly, we can also find that on the attribute *Thermal Crossover* which are all thermal videos, the MDNet [52] which is an online learned tracker attain the best results. Even the Staple and SRDCF are better than most of the other Siamese trackers, such as SiamKPN, SiamCAR, SiamRPN++, SiamRCNN, KYS, etc. The huge contrast demonstrates that online learning is very important for the tracker which is trained on one domain and tested on another domain (for example, the tracker trained on RGB videos and tested on Thermal videos).

C.2. Efficiency Analysis

In this work, two baseline methods are proposed for the *natural language initialized tracking* (Our-I) and *natural language guided tracking* (Our-II). For Our-I, the overall running efficiency is 24.39 FPS on the OTB-Lang, tested on a laptop with Intel Core I7, RTX2070. For Our-II, the overall efficiency on the OTB-Lang is 12.44 FPS.

C.3. More Visualization

In this section, more visualization on the tracking results is given to better understand our proposed method. As shown in Figure 11, 20 video sequences from OTB-Lang are selected to demonstrate the results of the visual grounding module. From the first three rows, we can find that the grounding module can locate the target object accurately when the background is relatively clean.

Table 7. Summary of evaluated trackers on TNL2K dataset.

Index	Tracker	Publication	Feature	Update	Pre-train	Search Scheme	FPS	Results
001	SiamRCNN [59]	CVPR-2020	ResNet-101	X	✓	Local + Global	5@GPU	0.528 0.523
002	SuperDiMP [4]	ICCV-2019	ResNet-50	✓	✓	Local	40@GPU	0.484 0.492
003	LTMU [9]	CVPR-2020	ResNet-50	✓	✓	Local	13@GPU	0.473 0.485
004	PrDiMP50 [12]	CVPR-2020	ResNet-50	✓	✓	Local	30@GPU	0.459 0.470
005	KYS [22]	ECCV-2020	ResNet-50	✓	✓	Local + Global	20@GPU	0.435 0.449
006	DiMP50 [4]	ICCV-2019	ResNet-50	✓	✓	Local	40@GPU	0.434 0.447
007	TACT [6]	ACCV-2020	ResNet-50	X	✓	Local + Global	42@GPU	0.422 0.438
008	SiamBAN [5]	CVPR-2020	ResNet-50	X	✓	Local	40@GPU	0.417 0.410
009	SiamRPN++ [35]	CVPR-2019	ResNet-50	X	✓	Local	35@GPU	0.412 0.413
010	CLNet [17]	ECCV-2020	ResNet-50	X	✓	Local	45@GPU	0.411 0.408
011	D3S [47]	CVPR-2020	ResNet-50	X	✓	Local	25@GPU	0.393 0.388
012	ATOM [11]	CVPR-2019	ResNet-50	X	✓	Local	30@GPU	0.392 0.401
013	SiamKPN [40]	arXiv-2020	ResNet-50	X	✓	Local	24@GPU	0.389 0.352
014	GlobalTrack [29]	AAAI-2020	ResNet-50	X	✓	Global	6@GPU	0.386 0.405
015	SiamCAR [23]	CVPR-2020	ResNet-50	X	✓	Local	52@GPU	0.384 0.353
016	DeepMTA [63]	TCSVT-2021	ResNet-50	✓	✓	Local + Global	12@CPU	0.381 0.385
017	SiamMask [62]	CVPR-2019	ResNet-50	X	✓	Local	55@GPU	0.380 0.383
018	Ocean [83]	ECCV-2020	ResNet-50	X	✓	Local	58@GPU	0.377 0.384
019	MDNet [52]	CVPR-2016	CNN-3	✓	✓	Local	1@GPU	0.371 0.384
020	SiamFC++ [72]	AAAI-2020	GoogLeNet	✓	✓	Local	90@GPU	0.369 0.386
021	VITAL [57]	CVPR-2018	CNN-3	✓	✓	Local	1.5@GPU	0.353 0.366
022	Meta-Tracker [53]	ECCV-2018	CNN-3	✓	✓	Local	1@GPU	0.333 0.339
023	SiamDW [82]	CVPR-2019	Res22W	X	✓	Local	150@GPU	0.326 0.323
024	RT-MDNet [31]	ECCV-2018	CNN-3	✓	✓	Local	46@GPU	0.322 0.308
025	SPLT [74]	ICCV-2019	ResNet-50	X	✓	Local + Global	25@GPU	0.321 0.337
026	GradNet [38]	ICCV-2019	CNN-5	✓	✓	Local	80@GPU	0.318 0.317
027	ECO [10]	CVPR-2017	VGG	✓	X	Local	8@CPU	0.317 0.326
028	MemTracking [75]	ECCV-2018	CNN-5	✓	✓	Local	50@GPU	0.305 0.304
029	MAML [60]	CVPR-2020	ResNet-50	X	✓	Local	40@GPU	0.295 0.284
030	DaSiamRPN [84]	ECCV-2018	ResNet-50	X	✓	Local	110@GPU	0.288 0.329
031	FCOT [8]	arXiv-2020	ResNet-50	X	✓	Local	45@GPU	0.288 0.320
032	SiamFC [3]	ECCVW-2016	CNN-5	X	✓	Local	58@GPU	0.286 0.295
033	SiamRPN [36]	CVPR-2018	ResNet-50	X	✓	Local	160@GPU	0.281 0.300
034	ADNet [80]	CVPR-2017	CNN-3	✓	✓	Local	3@GPU	0.278 0.285
035	UDT [61]	CVPR-2019	CNN-5	X	✓	Local	70@GPU	0.271 0.266
036	Staple [2]	CVPR-2016	HOG	✓	X	Local	80@CPU	0.270 0.270
037	SRDCF [14]	ICCV-2015	HOG	✓	X	Local	6@CPU	0.269 0.265
038	GOTURN [26]	ECCV-2016	CaffeNet-5	X	✓	Local	100@GPU	0.205 0.198
039	RTAA [30]	ECCV-2018	ResNet-50	X	✓	Local	2.2@GPU	0.193 0.217
040	KCF [27]	TPAMI-2015	HOG	✓	X	Local	172@CPU	0.153 0.200
041	VisGround [77]	ICCV-2019	DarkNet-53	X	✓	Global	147@GPU	0.143 0.159
042	ROAM [76]	CVPR-2020	ResNet-50	X	✓	Local	13@GPU	0.108 0.157
043	MIL [1]	CVPR-2009	HOG	✓	X	Local	25@CPU	0.063 0.042

Also, it works well in some challenge videos, like *car*, and *human head*. For the fourth row, the grounding is not accurate enough for tracking, including the central location and scale. We can find that the performance of visual grounding is needed to be further improved for more accurate tracking. More experimental results of our proposed baseline and other trackers can be found in Figure 12.

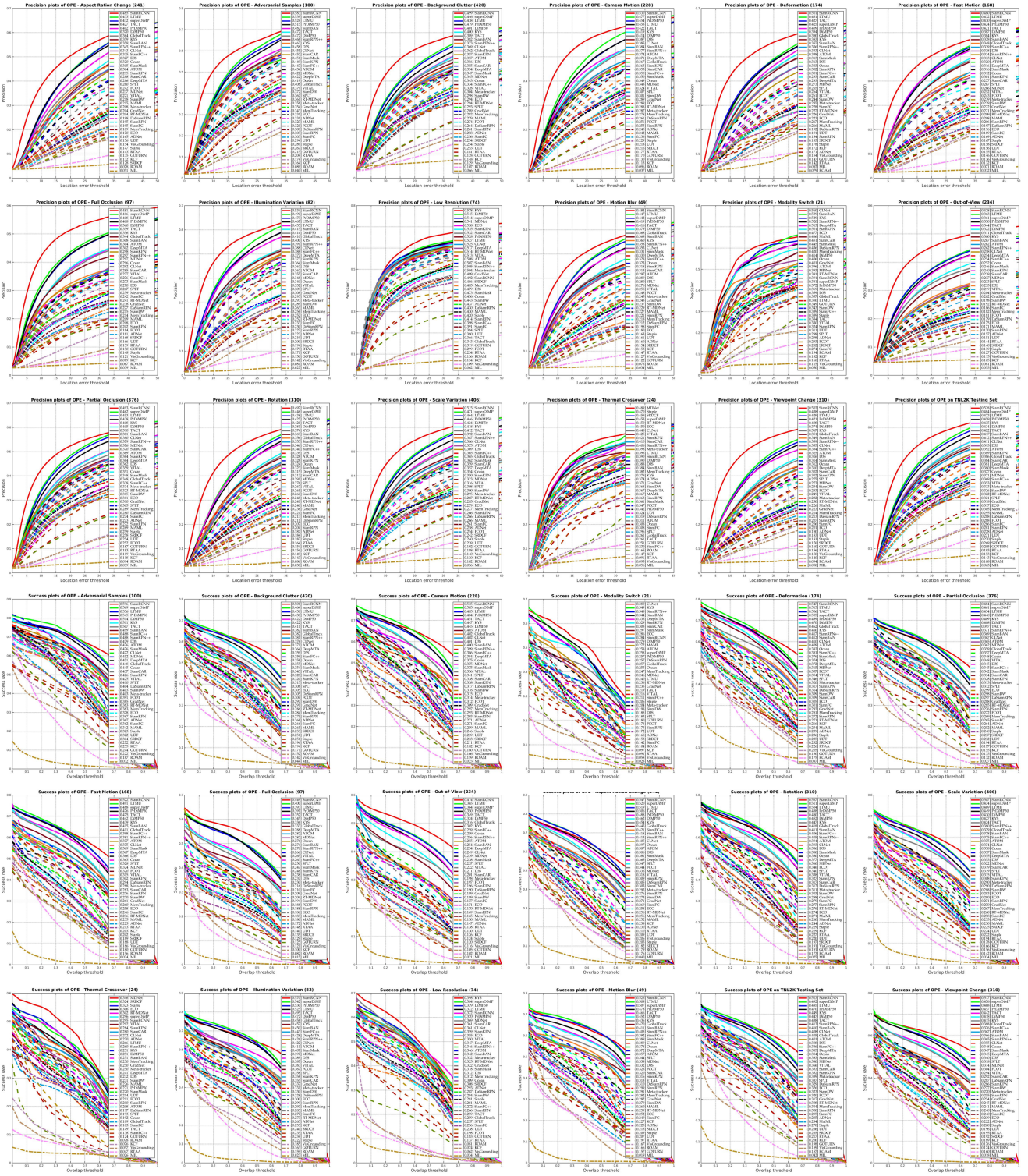


Figure 10. Tracking results under each challenging factors on TNL2K dataset (Tracking-by-BBox). Best viewed by zooming in.

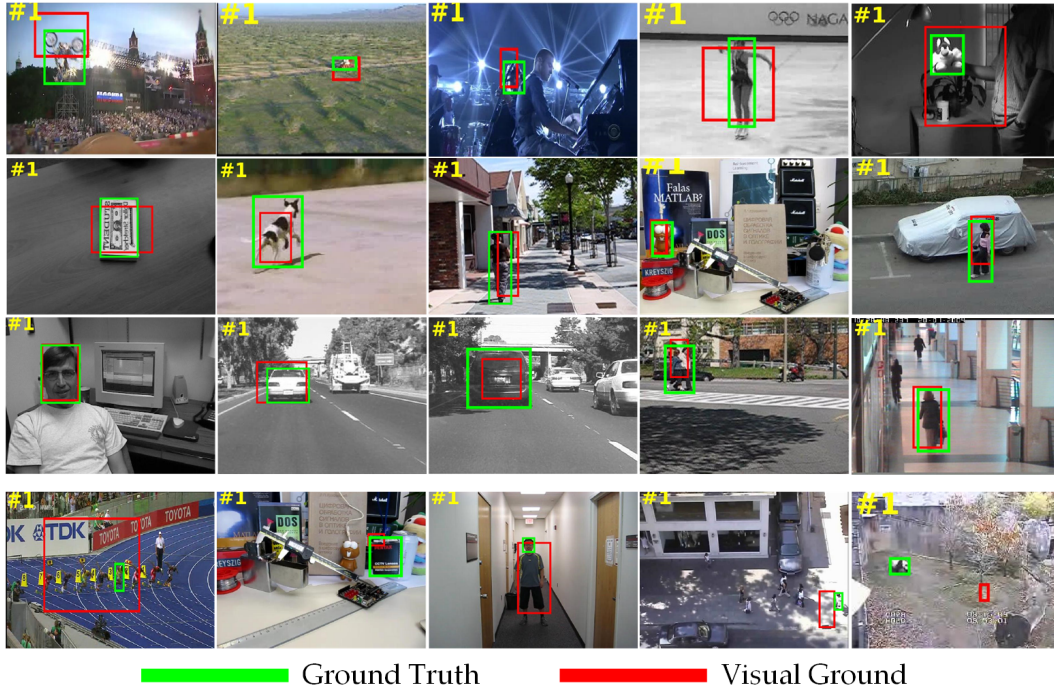


Figure 11. Results of the first frame of visual grounding module.



Figure 12. Tracking results of our method and other state-of-the-art tracking algorithms.