# Differentiating Objects by Motion:
# Joint Detection and Tracking of Small Flying Objects

Ryota Yoshihashi     Tu Tuan Trinh     Rei Kawakami

The University of Tokyo

Shaodi You     Makoto Iida     Takeshi Naemura

CSIRO-Data61     The University of Tokyo

Australian National University

{yoshi, tu, rei, naemura}@hc.ic.i.u-tokyo.ac.jp
iida@ilab.eco.rcast.u-tokyo.ac.jp
shaodi.you@data61.csiro.au

(a) Poor visual information

(b) Large deformation and background clutter

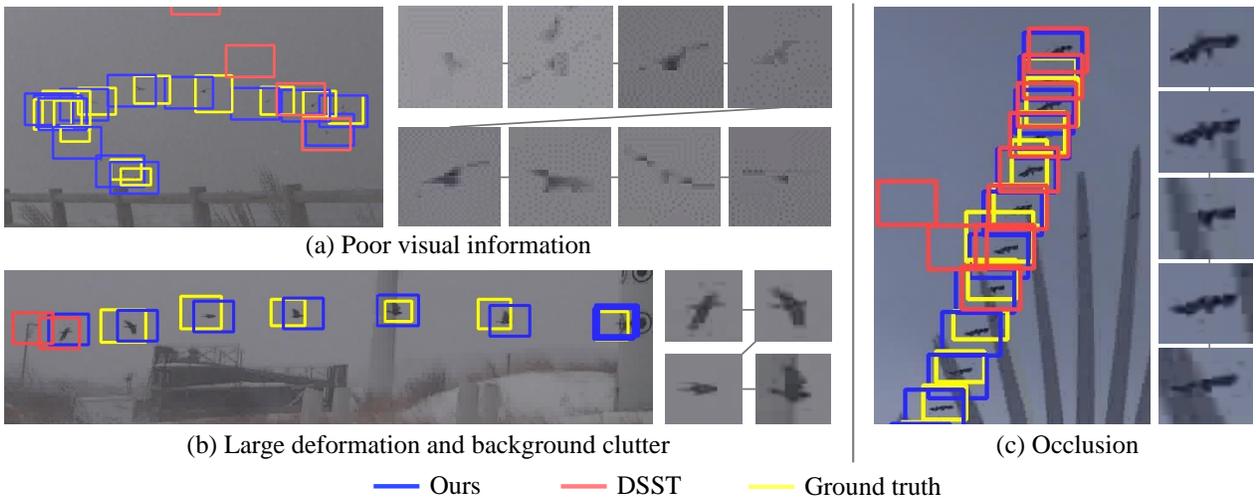(c) Occlusion

— Ours     — DSST     — Ground truth

Figure 1: Importance of multi-frame information for recognizing apparently small flying objects (birds in these examples). While visual features in single frames are vague and limited, multi-frame information, including deformation and pose changes, provides better clues with which to recognize birds. To extract such useful motion patterns, tracking is necessary for compensating translation of objects, but the tracking itself is a challenge due to the limited visual information. The blue boxes are birds tracked by our method that utilizes multi-frame representation for detection, while the red boxes are the results of a single-frame handcrafted-feature-based tracker [11] , which tends to fail when tracking small objects.

## Abstract

*While generic object detection has achieved large improvements with rich feature hierarchies from deep nets, detecting small objects with poor visual cues remains challenging. Motion cues from multiple frames may be more informative for detecting such hard-to-distinguish objects in each frame. However, how to encode discriminative motion patterns, such as deformations and pose changes that characterize objects, has remained an open question. To learn them and thereby realize small object detection, we present a neural model called the Recurrent Correlational Network, where detection and tracking are jointly performed over a multi-frame representation learned through a single, train-able, and end-to-end network. A convolutional long short-term memory network is utilized for learning informative appearance change for detection, while learned representation is shared in tracking for enhancing its performance. In experiments with datasets containing images of scenes with small flying objects, such as birds and unmanned aerial vehicles, the proposed method yielded consistent improvements in detection performance over deep single-frame detectors and existing motion-based detectors. Furthermore, our network performs as well as state-of-the-art generic object trackers when it was evaluated as a tracker on the bird dataset.*

## 1. Introduction

Detection of visually small objects is often required in wide-area surveillance [6, 7, 46]. Rich visual representations by deep convolutional networks (convnets) [18] pretrained on a large-scale still-image dataset [47] are of limited use on such objects, because they appear blurred and textureless in images owing to their small apparent size. For them to be detected, motion, namely the changes in their temporal appearance over a longer time frame, may offer richer information than appearance at a glance. As shown in Fig. 1, a bird is much easier to identify when multiple frames are available. However, it remains unclear how to learn motion features that are powerful enough to differentiate object.

In this paper, we present a method that exploits motion cues for small object detection. Although we utilize learnable pipelines based on convolutional and recurrent networks, our key idea is letting the network focus on informative deformations such as flapping of wings to differentiate target objects for detection, while removing less useful translations [43] by simultaneously tracking them with the learned visual representation. To make this possible, our framework performs joint detection and tracking. It utilizes convolutional long short-term memory (ConvLSTM) [59] to learn a discriminative multi-frame representation for detection, while it also enables correlation-based tracking over its output. Tracking is aided by the shared representation afforded by the training of the detector, and the overall framework is simplified because there are fewer parameters to be learned. We refer to the pipeline as *Recurrent Correlational Network*. Experiments on single-class, fully supervised small object detection in videos targeting birds [52] and unmanned aerial vehicles (UAVs) [46] show consistent improvements by our network over single-frame baselines and previous multi-frame methods. When evaluated as a tracker, ours also outperforms existing hand-crafted-feature-based and deep generic-object trackers in the bird dataset.

Our contribution is three-fold. First, we show that motion patterns learned via ConvLSTM improves detection performance in small object detection. Our network outperforms single-frame baselines, score-averaging baselines, and existing multi-frame methods in flying-object datasets, which indicates the importance of motion cues in these domains. Second, we introduce a novel framework for simultaneous object detection and tracking in video, which efficiently handles motion learning. This is the first recurrent model to achieve joint detection and tracking with deep learning. Third, our network is accurate when evaluated as a separate tracker in the dataset where class-specific detectors can be trained. The proposed network outperforms existing trackers based on various hand-crafted features, and performs slightly better or on par with convnet-based track-

ers. Our results gives a prospect toward domain-specific multi-task representation learning, which should open up application fields that generic detectors or trackers do not directly generalize. The relevant code and data will be published upon acceptance of this paper.

## 2. Related work

**Small object detection** Detection of small objects has been tackled in the surveillance community [6], and recently has attracted much attention since the advent of UAVs [7, 48]. Small pedestrians [3] and faces [26] have also been considered, and some recent studies try to detect small common objects in generic-object detection setting [4, 34]. Studies are more focused on scale-tuned convnets with moderate depths and a wider field of view, and despite of its importance, motion has not yet been incorporated in these domains.

**Object detection in video** Having achieved significant success in generic object detection in still images [18, 17, 45, 37, 8, 44], the research trends have begun moving toward efficient generic object detection in videos [47]. The video detection task poses new challenges, such as how to process voluminous video data efficiently and how to handle appearance of objects differing from still images due to rare poses [16, 61]. Very recent studies have begun improving on detection in videos; examples include T-CNNs [30, 31] that use trackers for propagating high-confidence detection, and deep feature flow [62] and flow-guided feature aggregation [61] that involves feature-level smoothing using optical flow. One of the closest idea to ours is joint detection and bounding-box linking by coordinate regression [16]. These models that have been used in ILSVRC-VID are more like modeling temporal consistency than understanding motion. Thus, it remains unclear whether or how inter-frame information extracted from motion or deformation aids in understanding objects. In addition, they all are based on popular convolutional generic still-image detectors [8, 17, 18, 37, 44, 45] and it is not clear to what extent such generic object detectors, which are designed for and trained in dataset collected from the web, generalize to task-specific datasets [13, 25, 60]. In the datasets for flying objects detection that we use [52, 46], the domain gap is especially large due to differences in the appearance of objects and backgrounds, as well as scale of objects. Thus, we decided to use simpler region proposals and fine-tune our network as region classifiers in each dataset.

**Deep trackers** Recent studies have intensively examined convnets and recurrent nets for tracking. Convnet-based trackers learn convolutional layers to acquire rich visual representation. Their localization strategies are diverse, including classification-based [40], similarity-learning-based [35], regression-based [22], and correlation-
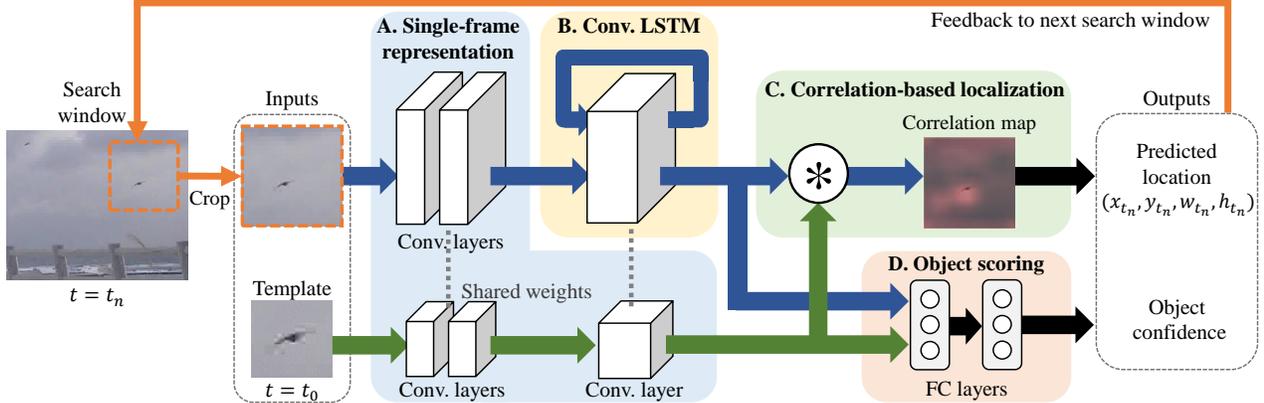
Figure 2: Overview of the proposed network, called *Recurrent Correlation Network* (RCN). It consists of the four modules: Convolutional layers for single-frame representations (A), ConvLSTM layers for multi-frame representations (B), cross-correlation layers for localization (C), and fully-connected layers for object scoring (D). Green arrows show the information stream from templates (the proposals in the first frame at $t = t_0$), and blue arrows show that from search windows.

based [2, 53] approaches. While classification of densely sampled patches [40] is the most accurate in generic benchmarks, its computation is slow and regression-based one [22] and correlation-based ones [2, 53] are used instead in real-time. Our network also incorporates a correlation-based localization mechanism, having its performance enhanced by the representation shared by the detector.

Recurrent nets [57, 23] efficiently handle temporal structures in sequences and thus, they have been used for tracking [41, 20, 39, 55]. However, most utilize separate convolutional and recurrent layers, and have a fully connected recurrent layer, which may lead to a loss of spatial information. Thus, currently recurrent trackers do not perform as well as the best single-frame convolutional trackers in generic benchmarks. One study used ConvLSTM with simulated robotic sensors for handling occlusion [42].

**Joint detection and tracking** The relationship between object detection and tracking is a long-term problem in itself; before the advent of deep learning, it had only been explored with classical tools. In the track–learn–detection (TLD) framework [29], a trained detector enables long-term tracking by re-initializing trackers after temporal disappearance of objects. Andriluka *et al.* uses a single-frame part-based detector and shallow unsupervised learning based on temporal consistency [1]. Tracking by associating detected bounding boxes [27] is another popular approach. However, in this framework, recovering undetected objects is challenging because tracking is more akin to post-processing following detection than to joint detection and tracking.

**Motion feature learning** Motion feature learning, and hence the use of recurrent nets, are more active in video classification [32] and action recognition [51]. Studies have shown that LSTMs yield improvement in accuracy [54, 56, 14]. For example, VideoLSTM [36] uses the idea of inter-frame correlation to recognize actions with attention. How-ever, with action recognition datasets, the networks may not fully utilize human motion features apart from appearance, backgrounds and contexts [21].

Optical flow [38, 24, 15] is a pixel-level alternative to trackers to describe motion [43, 19, 62, 61]. Accurate flow estimation is, however, challenging in small flying object detection tasks due to the small apparent size of the targets and the large inter-frame disparity by fast motion [46]. While we focus on high-level motion stabilization and motion-pattern learning via tracking, we believe flow-based low-level motion handling is orthogonal and complementary to ours depending on the application areas.

## 3. Recurrent Correlational Networks

To exploit motion information via simultaneous detection and tracking, we present the *Recurrent Correlational Networks* as shown in Fig. 2. The network consists of four modules: (A) convolutional layers, (B) ConvLSTM layers, (C) a cross-correlation layer, and (D) fully connected layers for object scoring. First, the convolutional layers model single-frame appearances of target and non-target regions, including other objects and backgrounds. Second, the ConvLSTM layers encode temporal sequences of single-frame appearances, and extract the discriminative motion patterns. Third, the cross-correlation layer convolves the representation of the template to that of search windows in subsequent frames, and generates correlation maps that are useful for localizing the targets. Finally, the confidence scores of the objects are calculated with fully-connected layers based on the multi-frame representation. The network is supervised by the detection loss, and the tracking gives locational feedback for region of interest in next frames during training and testing.

Our detection pipeline is based on region proposal and classification of the proposal, as in region-based CNNs [18]. The main difference is in that our joint detection and track-

ing network simultaneously track the given proposals in the following frames, and the results of the tracking are reflected in the classification scores, that are used as detectors' confidence scores.

**Convolutional LSTM** In our framework, the ConvLSTM module [59] is used for motion feature extraction (Fig. 2 B). It is a convolutional counterpart of LSTM [23]. It replaces inner products in the LSTM with convolution, and this is more suitable for motion learning, since the network is more sensitive to local spatio-temporal patterns rather than the global patterns. It works as a sequence-to-sequence predictor; specifically, it takes series $(x_1, x_2, x_3, ..., x_t)$ of single-frame representations whose length is $t$ as input, and outputs a merged single representation $h_t$, at each timestep $t = 1, 2, 3, ..., L$.

For the sake of completeness, we show the formulation of ConvLSTM below.

$$
\begin{aligned}
i_t &= \sigma(w_{xi} * x_t + w_{hi} * h_{t-1} + b_i) \\
f_t &= \sigma(w_{xf} * x_t + w_{hf} * h_{t-1} + b_f) \\
c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(w_{xc} * x_t + w_{hc} \circ h_{t-1} + b_c) \\
o_t &= \sigma(w_{xo} * x_t + w_{ho} * h_{t-1} + b_o) \\
h_t &= o_t \circ \tanh(c_t).
\end{aligned}
\tag{1}
$$

Here, $x_t$ and $h_t$ denote the input and output of the layer at timestep $t$, respectively. The states of the memory cells are denoted by $c_t$. $i_t$, $f_t$, and $o_t$ and are called gates, which work for selective memorization. '∘' denotes the Hadamard product. ConvLSTM is also well suited to exploit the spatial correlation for joint tracking, since its output representations are in 2D.

While ConvLSTM is effective at video processing, it inherits the complexity of LSTM. The gated recurrent unit (GRU) is a simpler alternative to LSTM that has fewer gates, and it is empirically easier to train on some datasets [5]. A convolutional version of the GRU (ConvGRU) [49] is as follows:

$$
\begin{aligned}
z_t &= \sigma(w_{xz} * x_t + w_{hz} * h_{t-1} + b_z) \\
r_t &= \sigma(w_{xr} * x_t + w_{hr} * h_{t-1} + b_r) \\
h_t &= z_t \circ h_{t-1} + (1 - z_t) \\
&\quad \circ \tanh(w_{xh} * x_t + w_{hh} * (r_t \circ h_{t-1}) + b_h).
\end{aligned}
\tag{2}
$$

ConvGRU has only two gates, namely an update gate $z_t$ and reset gate $r_t$, while ConvLSTM has three. ConvGRU can also be incorporated into our pipeline; later we provide an empirical comparison between ConvLSTM and ConvGRU.

**Correlation-based localization** The correlation part (Fig. 2 C) aims to stabilize moving objects' appearance by tracking. The localization results are fed back to the next input, as shown in Fig. 3. This feedback allows ConvLSTM
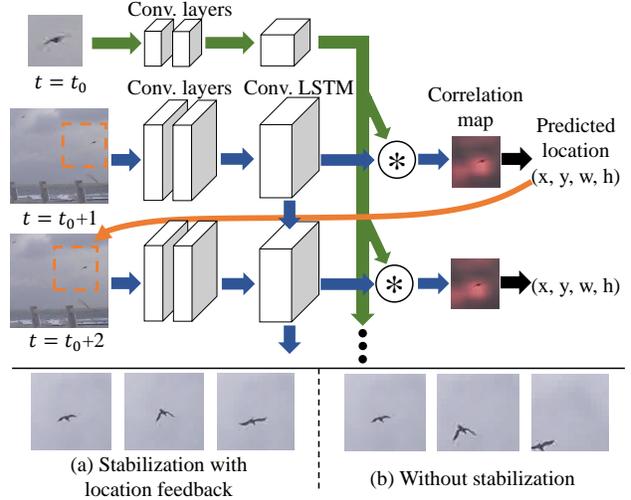


Figure 3: Temporal expansion of the proposed network. The joint tracking is incorporated as part of the feedback in the recurrent cycle. This feedback provides stabilized observation of moving objects (a), while learning from deformation is difficult without stabilization (b).

to learn deformations and pose changes apart from translation (Fig. 3 a), while local motion patterns are invisible due to translation without stabilization (Fig. 3 b).

Cross-correlation is an operation that relates two inputs, and outputs a correlation map that indicates how similar each patch in an image is to another. It is expressed as

$$
C(\boldsymbol{p}) = \boldsymbol{f} * \boldsymbol{h} = \sum_{\boldsymbol{q}} \boldsymbol{f}(\boldsymbol{p} + \boldsymbol{q}) \cdot \boldsymbol{h}(\boldsymbol{q}).
\tag{3}
$$

where $\boldsymbol{f}$ and $\boldsymbol{h}$ denote the multi-dimensional feature representations of the search window and template, respectively. $\boldsymbol{p}$ is for every pixel's coordinates in the domain of $\boldsymbol{f}$, and $\boldsymbol{q}$ is for the same but in the domain of $\boldsymbol{h}$. Two-dimensional (2D) correlation between a target patch and a search window is equivalent to densely comparing the target patch with all possible patches within the search window. The inner product is used here as a similarity measure.

In the context of convolutional neural networks, the cross-correlation layers can be considered to be differentiable layers without learnable parameters; namely, a cross-correlation layer is a variant of the usual convolutional one whose kernels are substituted by the output of another layer. Cross-correlation layers are bilinear with respect to two inputs, and thus are differentiable. The computed correlation maps are used to localize the target by

$$
\boldsymbol{p}_{target} = \text{argmax}_{\boldsymbol{p}} C(\boldsymbol{p})
\tag{4}
$$

**Single-frame representation** A multi-layer convolutional representation is inevitable in natural image recognition, although the original ConvLSTM [59] did not

use non-recurrent convolutional layers in radar-based tasks. Following recent tandem CNN-LSTM models for video recognition [14], we insert non-recurrent convolutional layers before the ConvLSTM layers (Fig. 2 A). Arbitrary covolutional architectures can be incorporated and we should choose the proper ones for each dataset. We experimentally tested two different structures of varying depth.

We need to extract an equivalent representation from the object template for the search windows. For this, we use ConvLSTM, in which the recurrent connection is severed. Specifically, we force the forget gates to be zero and enter zero vectors instead of the previous hidden states. This layer is equivalent to a convolutional layer with $tanh$ nonlinearity and sigmoid gates. It shares weights with $w_{xc}$ in Eq. 1.

**Search window strategy**  In object tracking, as the speed of the target objects is physically limited, limiting the area of the search windows, where the correlations are computed, is a natural way to reduce computational costs. We place windows the centers of which are at the previous locations of the objects, with a radius $R = \alpha \max(W, H)$, where $W$ and $H$ are the width and height of the bounding box of the candidate object. We then compute the correlation map for windows around each candidate object. We empirically set the size of the search windows to $\alpha = 1.0$. The representation extracted from the search windows is also fed to the object scoring part of the network, which yields large field-of-view features and provides contextual information for detection.

**Object scoring**  For object detection, the tracked candidates need to be scored according to likeness. We use fully connected (FC) layers for this purpose (Fig. 2 D). We feed both the representations from the templates (green lines in Fig. 2) and the search windows (blue lines in Fig. 2) into the FC layers by concatenation. We use two FC layers, where the number of dimensions in the hidden vector was 1,000.

We feed the output of each timestep of ConvLSTM into the FC layers and average the scores. In theory, the representation of the final timestep after feeding the last frame of the sequence should provide the maximum information. However, we found that the average scores are more robust in case of tracking failures or the disappearance of targets.

**Training**  Our network is trainable with ordinary gradient-based optimizers in an end-to-end manner because all layers are differentiable. We separately train convolutional parts and ConvLSTM to ensure fast convergence and avoid overfitting. We first initialize single-frame-based convnets by pre-trained weights in the ILSVRC2012-CLS dataset, the popular and largest generic image dataset. We then fine-tune single-frame convnets in the target datasets (birds and drones) without ConvLSTM. Finally, we add the convolutional LSTM, correlation layer, and FC layers to the networks and fine-tune them again. For optimization, we use

Table 1: Statistics of the datasets.

|  | Bird [52] | UAV [46] |
|---|---|---|
| Frame resolution | $3840 \times 2160$ | $752 \times 480$ |
| Ave. object resolution | 55 pixels | 18 pixels |
| #Test frames | 2,222 | 5,800 |
| #Training boxes | 10,000 | 8,000 |

the SGD solver of Caffe [28]. In the case reported here, the total number of iterations was 40,000 and the batch size was five. The original learning rate was 0.01, and was reduced by a factor of 0.1 per 10,000 iterations. The loss was the usual sigmoid cross-entropy for detection. We freeze the weights in the pre-trained convolutional layers after connecting to the convolutional LSTM to avoid overfitting.

During training of ConvLSTM, we use pre-computed trajectories predicted by a single-frame convolutional tracker, which consists of the final convolutional layers of the pre-trained single-frame convnet and a correlation layer. They are slightly inaccurate but have similar trajectories to those of our final network. Then, we store cropped search windows in the disk during training for efficiency, to reduce disk access by avoiding the re-cropping of the regions of interest out of the 4K-resolution frames during training. During the testing phase, the network observes trajectories estimated by itself, which are different from the ground truths that are used in the training phase. This training scheme is often referred to as teacher forcing [58]. Negative samples also need trajectories in training, but we do not have their ground truth trajectories because only the positives are annotated in the detection datasets.

## 4. Experiments

The main purpose of the experiments was to investigate the performance gain owing to the learned motion patterns with joint tracking in small object detection tasks. We also investigated the tracking performance of our method and compared it with that of trackers with a variety of features as well as convolutional trackers.

We first used a recently constructed video-based bird dataset [52]. This dataset involves detecting birds around a wind farm. The resolution is 4K and the frame rate is 30 fps, which made processing the dataset a challenge due to its large volume. The most frequent size of the birds is 55 pixel. Although the dataset consists of images taken from a fixed-point camera, it has changes in illumination owing to the weather, changing background patterns owing to clouds, and variation in the appearance of birds due to occlusion and deformation. We also tested our method on a UAV dataset [46] to see whether it can be applied to other flying objects. This dataset consists of 20 sequences of hand-captured videos. It consists of approximately 8,000 bounding boxes of flying UAVs. All the UAVs in this dataset are multi-copters. We followed the training/testing split provided by the authors of [46]. The properties of the dataset
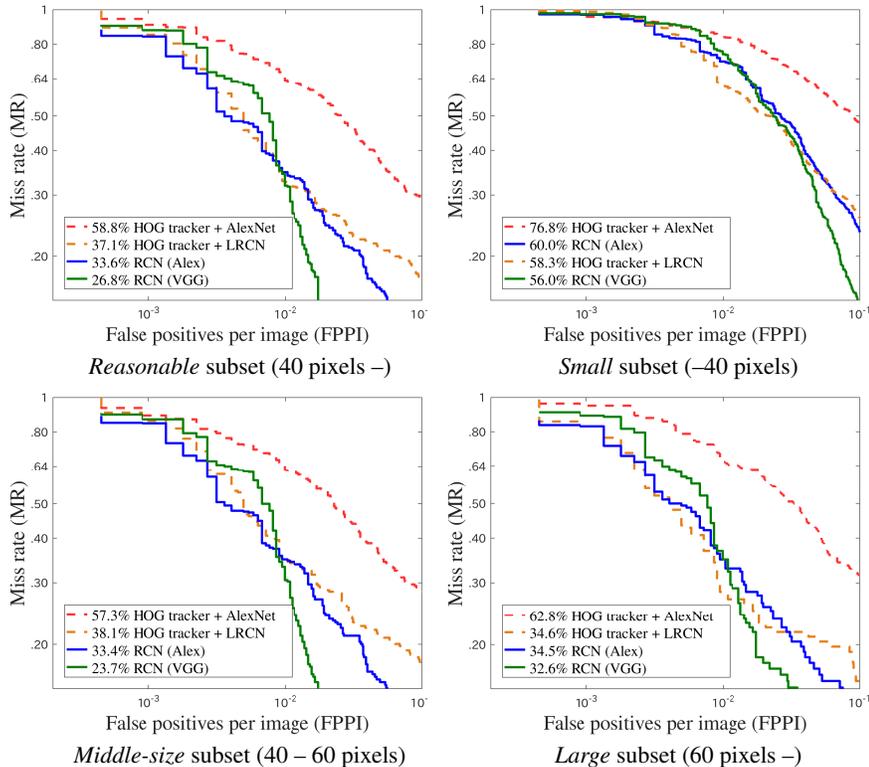
Figure 4: Detection results. The lower left is better. Our RCN (VGG) outperformed all of the other methods with deeper convolutional layers, and our RCN (Alex) outperformed the previous method with the same convolutional layer depth on three subsets. The subsets are distinguish by the sizes of birds in the images.
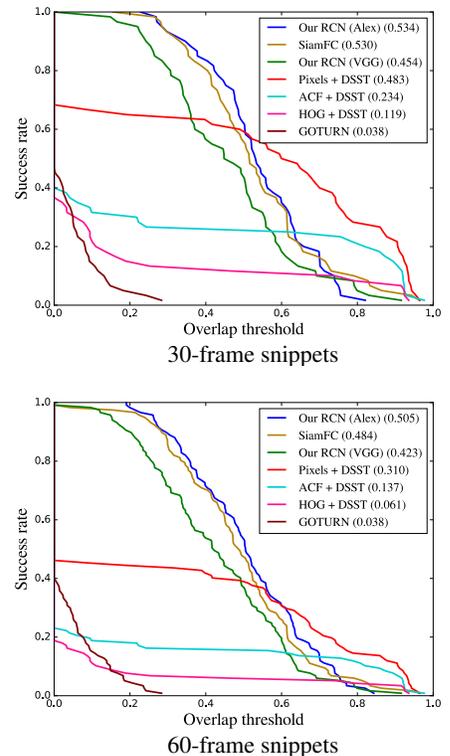


Figure 5: Tracking results. The upper right is better. The proposed methods outperformed DSST trackers with various handcrafted features and the ImageNet-pretrained deep trackers.

are summarized in Table 1.

**Evaluation metric** To evaluate detection performance, we used the number of false positives per image (FPPI) and the log average miss rate (MR). These metrics were based on single-image detection; *i.e.*, they were calculated only on given test frames that were sampled discretely. Detection was performed on the given test frames and, for our method, tracking of all candidates was conducted in some of the subsequent frames. We used the toolkit provided for the Caltech Pedestrian Detection Benchmark [13] to calculate the scores and plot the curves in Fig. 4.

We also tested tracking accuracy separately from detection on the bird detection dataset. We fed the ground-truth bounding boxes in the first frames to our network and other trackers, aiming to evaluate our joint detection and tracking network as a tracker. We conducted one-path evaluation (OPE), tracking by using ground truth bounding boxes given only in the first frame of the snippets without re-initialization, re-detection, or trajectory fusion. To remove very short trajectories to evaluate trackers, we selected ground truth trajectories longer than 90 frames (three seconds at 30 fps) from the annotation of the bird dataset. We plotted success rates versus overlap thresholds. The curves in Fig. 5 show the proportion of the estimated bounding boxes whose overlaps with the ground truths were higher than the thresholds.

**Object proposals** We used a different strategy for each dataset to generate object proposals for pre-processing. In the bird dataset, we extracted the moving object by background subtraction [63]. The extracted regions were provided by the authors with the dataset; therefore, we could compare the networks fairly, regardless of the hyperparameters or the detailed tuning of the background subtraction. In the UAV dataset, we used the HOG3D-based sliding window detector provided by the authors of [46].

**Compared methods** *RCN (Alex)* and *RCN (VGG)* are two implementations of the proposed method using the convolutional layers from AlexNet [33] and VGG16Net [50]. *HOG tracker+AlexNet* and *HOG tracker+LRCN* are baselines for the bird dataset provided by [52]. The former is a combination of the HOG-based [9] discriminative scale-space tracker (DSST [10, 11]) and convnets that classify the tracked candidates into positives and negatives. The latter is a combination of DSST and the CNN-LSTM tandem model [14]. They used five frames following the test frames, for fair comparison, and our method used the same number of frames in the detection evaluation.

For evaluating the tracking performance, we included

other combinations of the DSST and hand-crafted features for further analysis. *HOG+DSST* is the original version in [10]. *ACF+DSST* replaces the classical HOG with more discriminative aggregated channel features [12]. The ACF is similar to HOG, but is more powerful because of the additional gradient magnitude and LUV channels for orientation histograms. *Pixel+DSST* is a simplified version that uses RGB values of raw pixels instead of gradient-based features. We also included ImageNet-pretrained convolutional trackers, namely, correlation-based SiamFC [2] and regression-based GOTURN [22]. They are based on the convolutional architecture of AlexNet.

**Results**  The results of detection on the bird dataset are shown in Fig. 4. The curves are for four subsets of the test set, which consists of birds of different sizes, namely *reasonable* (over 40 pixels square), *small* (smaller than 40 pixels square), *mid-sized* (40–60 pixels square), and *large* (over 60 pixels square).

On all subsets, the proposed method, *RCN (VGG)* showed the smallest average miss rate (MR) of the tested detectors. The improvements were -10.3 percentage points on *Reasonable*, -2.3 percentage points on *Small*, -14.4 on *Mid-sized*, and -2.0 percentage points on *Large* subset, in comparison with the previous best published method *HOG tracker+LRCN*.

A comparison of *HOG tracker+LRCN* and proposed *RCN (Alex)* is also important, because these share the same convolutional architecture. Our *RCN (Alex)* performed better on all of the subset except *Small*, without deepening the network. The margins are -3.5 percentage points on *Reasonable*, -4.7 percentage points on *Mid-sized* subset, and -0.1 percentage points on *Large* subset. Examples of the test frames and results are shown in Fig. 6 (more examples are in the supplementary material).

A comparison of *RCN (Alex)* and *RCN (VGG)* provides an interesting insight. *RCN (Alex)* is more robust against smaller FPPI values in spite of the lower average performance than that of *RCN (VGG)*. *RCN (Alex)* showed a smaller MR than *RCN (VGG)* when the FPPI was lower than $10^{-2}$. A possible reason is that a deeper network is less generalizable because of many parameters; thus, it may miss-classify new negatives more often in the test set than in the shallower one.

The results of tracking on the bird dataset are shown in Fig. 5. We found that gradient-based features were inefficient on this dataset. HOG-based DSST missed the target even in 30-frame short tracking (but this was already longer than in [52] for detection). We assume that this was because of the way the HOG normalizes the gradients, which might render it over-sensitive to low-contrast but complex background patterns, like clouds. We found that replacing HOG with ACF and utilizing gradient magnitudes and LUV values benefited the DSST on the bird dataset. However,

Table 2: Performance differences as a result of varying models and parameters. MR represents the log-average miss rate in the *reasonable* subset of the bird dataset, and diff. represents its difference from the baseline. $k$ denotes the kernel size of the ConvLSTM.

|  | Network config. | MR | diff. |
|---|---|---|---|
| RCN (Alex) | | | |
| k = 3 | A + B + C + D | **0.336** | 0 |
| k = 1 | A + B + C + D | 0.346 | + 0.010 |
| k = 5 | A + B + C + D | 0.347 | + 0.011 |
| RCN (VGG) | | | |
| k = 3 | A + B + C + D | **0.268** | 0 |
| ConvGRU k = 3 | A + B + C + D | 0.271 | + 0.003 |
| w/o tracking | A + B + D | 0.321 | + 0.053 |
| w/o ConvLSTM | A + C + D | 0.344 | + 0.076 |
| Single frame | A + D | 0.332 | + 0.064 |

the simpler pixel-DSST outperformed the ACF-DSST by a large margin.

The trajectories provided by our network were more robust than all of DSST variations tested. This shows that representations learned through detection tasks also work better in tracking than hand-crafted gradient features do. It also worth noting that our trajectories were less accurate than those obtained through the feature-based DSSTs when they did not miss the target. When bounding-box overlaps larger than 0.6 were needed, the success rates were smaller than those of the DSSTs on both 30- and 60-frame tracking. This is because our network used the correlation involving pooled representation, the resolution of which was 32 times smaller than that of the original images. In addition, Our RCN (Alex) outperformed two existing convnet-based trackers (GOTURN and SiamFC). Examples of the tracking results are presented in the supplementary material.

The resulting ROC curves of drone detection are shown in Fig. 7. We report the results of a shallower AlexNet-based version of our RCN, because of the size of the training data. We also show the curve yielded by AlexNet after single-frame pre-training without LSTM or tracking, which we refer to as *Our AlexNet only*. This simple implementation slightly outperformed the baseline in [46] without auxiliary multi-frame information by tracking or motion compensation. Our network was different in that it was deeper and larger, and had been pre-trained in ImageNet. It is interesting that pre-training in the ImageNet classification is useful even in this domain of small, grayscale UAV detection. The ConvLSTM and joint tracking consistently improved in detection performance (-4.3 percentage points). However, the performance gain was smaller than that on the bird dataset. The reason seemed to be that the amount motion information in the UAV dataset was limited because the objects were rigid, in contrast to the articulated deformation in birds. Examples of the results are shown in Fig. 8.

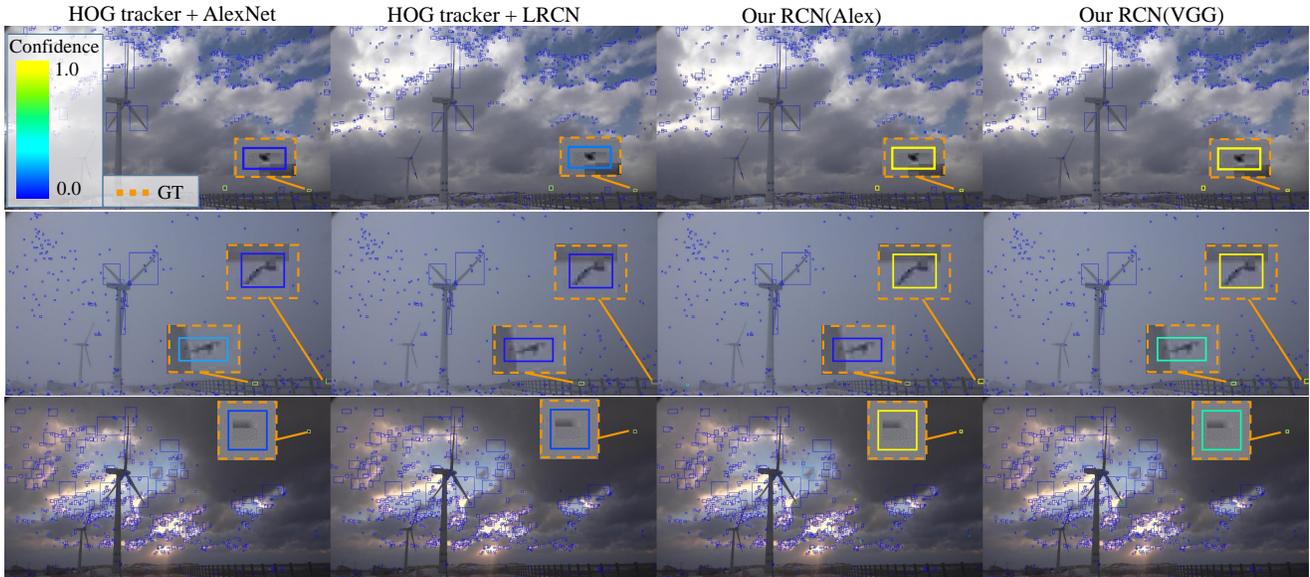**Hyperparameters and ablation**  Here, we report the fluc-

Figure 6: Example frames of results of detection on the bird dataset [52]. The dotted yellow boxes show ground truths, enlarged to avoid overlapping and keep them visible. The confidence scores of vague birds are increased and that of non-bird regions are decreased by our RCN detector. The contrast was modified for visibility in the zoomed-up samples.
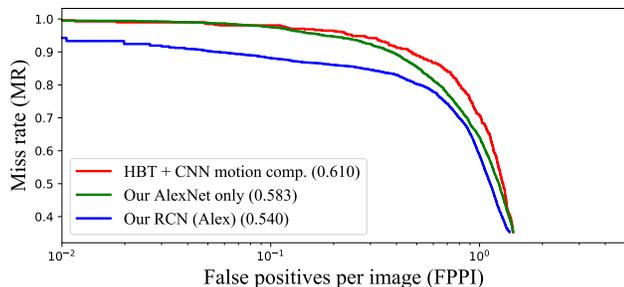


Figure 7: Detection results on the UAV dataset [46]. RCN performed the best.

tuation on performance for different settings of the networks and hyperparameters. We investigated the following factors: 1) kernel size in ConvLSTM, 2) ConvGRU vs. ConvLSTM, 3) w/o tracking, 4) w/o ConvLSTM, and 5) single frame detection. The kernel size controlled the receptive field of a memory cell. Second, we see the effect of simplifying ConvLSTM to ConvGRU. Third, we removed the joint tracker to see how useful multi-frame information was without stabilization. Fourth, we removed the recurrent part and averaged the confidence scores through time, to see the importance of the recurrent part. Finally, we used the network as a single-frame detector. The results are summarized in Table 2. Here *Network config. means which modules in Fig. 2 are active.* All of the results were in the *reasonable* subset of the bird dataset. The best kernel size was $k = 3$ in RCN (Alex). Larger and smaller kernels adversely affected performance slightly but not critically (+0.011 and + 0.010 MR). The performance of the ConvGRU was slightly worse than that of ConvLSTM (+0.003 MR), possibly because the input was pre-processed by convolutional layers



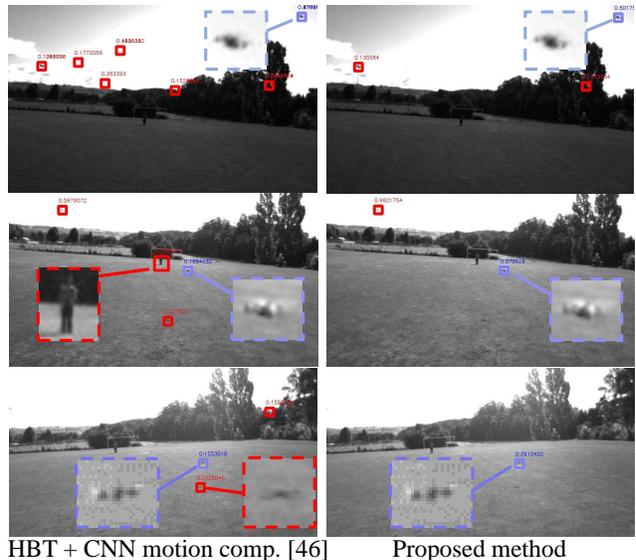HBT + CNN motion comp. [46]          Proposed method

Figure 8: Sample frames of detection results on the UAV dataset [46]. The blue boxes show correct detections and the red ones show misdetections. Our method caused fewer misdetections when the detectors thresholds were set to give roughly the same MR.

and the burden on the recurrent part was smaller. Lack of stabilization, recurrent parts, or multi-frame cues led to critical degradations in performance (+0.053, +0.076 MR and +0.064 MR), which in turn demonstrates effectiveness of the proposed network design.

## 5. Conclusion

We introduced the *Recurrent Correlation Network*, a novel joint detection and tracking framework that exploit motion information of small flying objects. In experiments, we tackled two recently developed datasets consisting of images of small flying objects, where the use of multi-frame information is inevitable due to poor per-frame visual information. The results showed that in such situations, multi-frame information exploited by the ConvLSTM and tracking-based motion compensation yields better detection performance. In future work, we will try to extend the framework to multi-class small object detection in videos.

## Acknowledgement

## References

[1] M. Andriluka, S. Roth, and B. Schiele. People-tracking-by-detection and people-detection-by-tracking. In *CVPR*, pages 1–8, 2008. 3

[2] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr. Fully-convolutional siamese networks for object tracking. In *ECCVW*, pages 850–865, 2016. 2, 3, 7

[3] R. Bunel, F. Davoine, and P. Xu. Detection of pedestrians at far distance. In *International Conference on Robotics and Automation (ICRA)*, pages 2326–2331. IEEE, 2016. 2

[4] C. Chen, M.-Y. Liu, O. Tuzel, and J. Xiao. R-cnn for small object detection. In *ACCV*, 2016. 2

[5] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Gated feedback recurrent neural networks. In *ICML*, pages 2067–2075, 2015. 4

[6] R. T. Collins, A. J. Lipton, T. Kanade, H. Fujiyoshi, D. Duggins, Y. Tsin, D. Tolliver, N. Enomoto, O. Hasegawa, P. Burt, et al. A system for video surveillance and monitoring. *VSAM final report*, pages 1–68, 2000. 2

[7] A. Coluccia, M. Ghenescu, T. Piatrik, G. De Cubber, A. Schumann, L. Sommer, J. Klatte, T. Schuchert, J. Beyerer, M. Farhadi, et al. Drone-vs-bird detection challenge at ieee avss2017. In *Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6. IEEE, 2017. 2

[8] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object detection via region-based fully convolutional networks. In *NIPS*, pages 379–387, 2016. 2

[9] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, volume 1, pages 886–893. IEEE, 2005. 6

[10] M. Danelljan, G. Häger, F. Khan, and M. Felsberg. Accurate scale estimation for robust visual tracking. In *BMVC*, 2014. 6, 7

[11] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg. Discriminative scale space tracking. *PAMI*, 39(8):1561–1575, 2017. 1, 6

[12] P. Dollár, R. Appel, S. Belongie, and P. Perona. Fast feature pyramids for object detection. *PAMI*, 36(8):1532–1545, 2014. 7

[13] P. Dollár, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *PAMI*, 34(4):743–761, 2012. 2, 6

[14] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, pages 2625–2634, 2015. 3, 4, 6

[15] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, pages 2758–2766, 2015. 3

[16] C. Feichtenhofer, A. Pinz, and A. Zisserman. Detect to track and track to detect. In *ICCV*, 2017. 2

[17] R. Girshick. Fast R-CNN. In *ICCV*, pages 1440–1448, 2015. 2

[18] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 2, 3

[19] S. Gladh, M. Danelljan, F. S. Khan, and M. Felsberg. Deep motion features for visual tracking. In *ICPR*, 2016. 3

[20] D. Gordon, A. Farhadi, and D. Fox. Re3: Real-time recurrent regression networks for object tracking. *arXiv preprint arXiv:1705.06368*, 2017. 3

[21] Y. He, S. Shirakabe, Y. Satoh, and H. Kataoka. Human action recognition without human. In *ECCVW*, pages 11–17. Springer, 2016. 3

[22] D. Held, S. Thrun, and S. Savarese. Learning to track at 100 fps with deep regression networks. In *ECCV*, 2016. 2, 7

[23] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 3, 4

[24] B. K. Horn and B. G. Schunck. Determining optical flow. *Artificial intelligence*, 17(1-3):185–203, 1981. 3

[25] J. Hosang, M. Omran, R. Benenson, and B. Schiele. Taking a deeper look at pedestrians. In *CVPR*, pages 4073–4082, 2015. 2

[26] P. Hu and D. Ramanan. Finding tiny faces. In *CVPR*, 2017. 2

[27] C. Huang, B. Wu, and R. Nevatia. Robust object tracking by hierarchical association of detection responses. In *ECCV*, pages 788–801. Springer, 2008. 3

[28] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACMMM*, pages 675–678, 2014. 5

[29] Z. Kalal, K. Mikolajczyk, and J. Matas. Tracking-learning-detection. *PAMI*, 34(7):1409–1422, 2010. 3

[30] K. Kang, H. Li, J. Yan, X. Zeng, B. Yang, T. Xiao, C. Zhang, Z. Wang, R. Wang, X. Wang, et al. T-CNN: Tubelets with convolutional neural networks for object detection from videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 2017. 2

[31] K. Kang, W. Ouyang, H. Li, and X. Wang. Object detection from video tubelets with convolutional neural networks. In *CVPR*, 2016. 2

[32] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014. 3

[33] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012. 6

[34] J. Li, X. Liang, Y. Wei, T. Xu, J. Feng, and S. Yan. Perceptual generative adversarial networks for small object detection. In *CVPR*, 2017. 2

[35] K. Li, Y. Kong, and Y. Fu. Multi-stream deep similarity learning networks for visual tracking. *IJCAI*, 2017. 2

[36] Z. Li, E. Gavves, M. Jain, and C. G. Snoek. VideoLSTM colves, attends and flows for action recognition. *arXiv preprint arXiv:1607.01794*, 2016. 3

[37] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. 2

[38] B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *IJCAI*, pages 674–679, 1981. 3

[39] A. Milan, S. H. Rezatofighi, A. R. Dick, I. D. Reid, and K. Schindler. Online multi-target tracking using recurrent neural networks. In *AAAI*, pages 4225–4232, 2017. 3

[40] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking. In *CVPR*, pages 4293–4302, 2016. 2

[41] G. Ning, Z. Zhang, C. Huang, Z. He, X. Ren, and H. Wang. Spatially supervised recurrent convolutional neural networks for visual object tracking. In *IEEE International Symposium on Circuits and Systems*, 2017. 3

[42] P. Ondruska and I. Posner. Deep tracking: Seeing beyond seeing using recurrent neural networks. In *AAAI*, 2016. 3

[43] D. Park, C. L. Zitnick, D. Ramanan, and P. Dollár. Exploring weak stabilization for motion feature extraction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2882–2889, 2013. 2, 3

[44] J. Redmon and A. Farhadi. YOLO9000: Better, faster, stronger. In *CVPR*, 2017. 2

[45] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015. 2

[46] A. Rozantsev, V. Lepetit, and P. Fua. Detecting flying objects using a single moving camera. *PAMI*, 39(5):879–892, 2017. 2, 3, 5, 6, 7, 8

[47] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 115(3):211–252, 2015. 2

[48] A. Schumann, L. Sommer, J. Klatte, T. Schuchert, and J. Beyerer. Deep cross-domain flying object classification for robust uav detection. In *Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6. IEEE, 2017. 2

[49] M. Siam, S. Valipour, M. Jagersand, and N. Ray. Convolutional gated recurrent networks for video segmentation. *arXiv preprint arXiv:1611.05435*, 2016. 4

[50] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 6

[51] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 3

[52] T. T. Trinh, R. Yoshihashi, R. Kawakami, M. Iida, and T. Naemura. Bird detection near wind turbines from high-resolution video using lstm networks. In *World Wind Energy Conference (WWEC)*, 2016. 2, 5, 6, 7, 8

[53] J. Valmadre, L. Bertinetto, J. F. Henriques, A. Vedaldi, and P. H. Torr. End-to-end representation learning for correlation filter based tracking. *CVPR*, 2017. 2, 3

[54] L. Wang, Y. Qiao, and X. Tang. Action recognition with trajectory-pooled deep- descriptors. In *CVPR*, pages 4305–4314, 2015. 3

[55] L. Wang, L. Zhang, and Z. Yi. Trajectory predictor by using recurrent neural networks in visual tracking. *IEEE Transactions on Cybernetics*, 2017. 3

[56] P. Weinzaepfel, Z. Harchaoui, and C. Schmid. Learning to track for spatio-temporal action localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3164–3172, 2015. 3

[57] P. J. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural networks*, 1(4):339–356, 1988. 3

[58] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280, 1989. 5

[59] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *NIPS*, pages 802–810, 2015. 2, 4

[60] L. Zhang, L. Lin, X. Liang, and K. He. Is faster r-cnn doing well for pedestrian detection? In *European Conference on Computer Vision*, pages 443–457. Springer, 2016. 2

[61] X. Zhu, Y. Wang, J. Dai, L. Yuan, and Y. Wei. Flow-guided feature aggregation for video object detection. In *ICCV*, 2017. 2, 3

[62] X. Zhu, Y. Xiong, J. Dai, L. Yuan, and Y. Wei. Deep feature flow for video recognition. In *CVPR*, 2017. 2, 3

[63] Z. Zivkovic. Improved adaptive gaussian mixture model for background subtraction. In *ICPR*, volume 2, pages 28–31. IEEE, 2004. 6