

Deep Reinforcement Learning for Visual Object Tracking in Videos

Da Zhang¹, Hamid Maei², Xin Wang¹, and Yuan-Fang Wang¹

¹Department of Computer Science, University of California at Santa Barbara

²Samsung Research America

{dazhang, xwang, yfwang}@cs.ucsb.edu hamid.maei@samsung.com

Abstract

Convolutional neural network (CNN) models have achieved tremendous success in many visual detection and recognition tasks. Unfortunately, visual tracking, a fundamental computer vision problem, is not handled well using the existing CNN models, because most object trackers implemented with CNN **do not effectively leverage temporal and contextual information among consecutive frames**. Recurrent neural network (RNN) models, on the other hand, are often used to process text and voice data due to their ability to learn intrinsic representations of sequential and temporal data. Here, we propose a novel neural network tracking model that is capable of integrating information over time and tracking a selected target in video. It comprises three components: a CNN extracting best tracking features in each video frame, an RNN constructing video memory state, and a reinforcement learning (RL) agent making target location decisions. The tracking problem is formulated as a decision-making process, and our model can be trained with RL algorithms to learn good tracking policies that pay attention to continuous, inter-frame correlation and maximize tracking performance in the long run. We compare our model with an existing neural-network based tracking method and show that the proposed tracking approach works well in various scenarios by performing rigorous validation experiments on artificial video sequences with ground truth. To the best of our knowledge, our tracker is the first neural-network tracker that **combines convolutional and recurrent networks with RL algorithms**.



1. Introduction

Given some object of interest marked in one frame of a video, the goal of *single-object tracking* is to locate this object in subsequent video frames, despite object motion, changes in the camera's viewpoint, and other incidental environmental variations such as lighting and shadows.

Single-object tracking finds applications in many important scenarios. For autonomous driving, a tracker must follow dynamic objects in order to estimate how they are moving and predict where they will end up in the future. For security applications, a surveillance camera must track numerous people as they move through the environments.

CNNs have recently had great success in significantly advancing the state-of-the-art on challenging computer vision tasks such as image classification [22, 15, 33], object detection [10], semantic segmentation [25], and many others [4, 38, 39, 19]. Such great success of CNNs is mostly attributed to their outstanding performance in representing *spatial* information, often identified from collections of unrelated, discrete image frames.

Visual tracking, however, has seen less success using this line of attack since it is often difficult to learn a robust *spatial* and *temporal* representation for continuous video data. Recently, some initial works have been done in using neural networks for visual tracking, and these efforts have produced some impact and improved state-of-the-art accuracy [29, 41, 40]. Although these methods may be sufficient to predict a target's bounding box locations, these classification-based approaches could be limited in learning robust tracking features and maximizing long-term tracking performance without taking temporal coherency and correlation into account. That is, training algorithms specialized for long-term visual tracking is still not fully explored.

To fully exploit the *spatial* and *temporal* information in videos for visual tracking, it is then desirable to add *continuous*, *recurrent* links to augment *discrete*, *convolutional* links and use a reinforcement algorithm to maximize tracking performance in the long run. Here, we propose a novel framework which processes sequential video frames as a whole and directly outputs location predictions of the target object in each frame. Our model integrates convolutional network with recurrent network (Figure 1), and builds up an internal representation of the video. It fuses past recurrent memories with current visual features to make

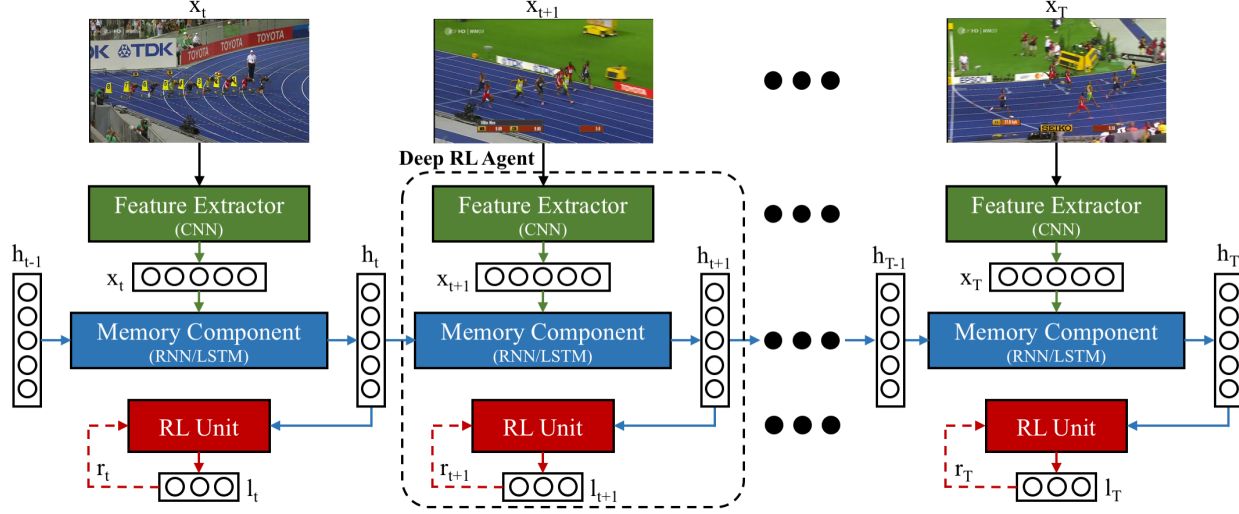


Figure 1: Overview of the proposed network architecture: At each time step t , the **feature extractor** takes an image x_t as input and computes a feature representation i_t . The **memory component** takes the feature representation i_t as input, combining with the internal memory state at previous time step h_{t-1} , and produces the new internal memory state h_t . The **RL unit** uses the internal memory state h_t to produce a target location decision l_t . The agent will receive a scalar reward signal r_t from the environment. The basic RNN iteration is repeated for a variable number of steps and the cumulative rewards $R = \sum_{i=1}^T r_t$ are used to update network parameters such that the long-term tracking performance is maximized.

predictions of the target object’s location over time. We describe an end-to-end optimization procedure that allows the model to be trained to maximize tracking performance in the long run. This procedure uses backpropagation to train the neural-network components and REINFORCE algorithm [42] to train the policy network.

Our algorithm augments traditional CNN with a recurrent convolutional model learning temporal representations and an RL unit making optimal tracking decisions. The main contributions of our work are:

- We propose a convolutional recurrent neural network model, which learns robust spatial representations in each single frame and temporal representations cross multiple frames. The learned features better capture temporal information in video and can be directly applied in a tracking problem.
- Our framework is trained end-to-end with deep RL algorithms, in which the model is optimized to maximize a tracking performance measure which depends on the entire video sequence.
- Our model is trained fully off-line. This off-line training strategy is advantageous in test time, as there is no need to fine-tune the system online to gain efficiency.
- We have validated our algorithms by comparing with an existing neural-network based method and perform-

ing rigorous experiments on artificial video sequences to demonstrate both the validity and applicability.

The rest of the paper is organized as follows. We first review related work in Section 2, and discuss our RL approach for visual tracking in Section 3. Section 4 describes our end-to-end optimization algorithm, and Section 5 demonstrates the experimental results.

2. Related Work

2.1. Visual Tracking Algorithms

Visual Tracking is a fundamental problem in computer vision and has been actively studied for decades. Many methods have been proposed for single-object tracking. For a systematic review and comparison, we refer the readers to a recent benchmark and a tracking challenge report [43, 21]. Generally speaking, most existing trackers belong to either one of two categories: generative trackers and discriminative trackers. Generative methods describe the target appearances using generative models and search for candidate image regions that fit the model best. Some representative methods are based on principal component analysis [32], and sparse representation [26, 44]. On the other hand, discriminative trackers learn to separate the foreground from the background using a classifier. Many advanced machine learning algorithms have been used, including online boosting [11, 12], structured output SVM [14], and multiple-

instance learning [2].

However, almost all existing appearance-based tracking methods rely on low-level, hand-crafted features which are incapable of capturing semantic information of the targets, not robust to significant appearance changes, and only have limited discriminative power.

2.2. Convolutional Neural Networks

CNNs have demonstrated their outstanding representation power in a wide range of computer vision applications [4, 10, 22, 25, 33, 38, 39, 15, 29, 19]. Krizhevsky *et al.* [22] brought significant performance improvement in image classification by training a deep CNN on ImageNet dataset [6]. Region CNN [10] applied a CNN to an object detection task, which can accurately locate and classify objects in natural images.

Despite the impressive success of CNNs applied to still images, the application of deep neural network (DNN) models in visual tracking is not fully explored. Although DNN trackers have improved the state-of-the-art, only a limited number of tracking algorithms using representations from CNNs have been proposed so far [9, 18, 24, 41, 40, 16, 29]. An early tracking algorithm transferred CNNs pretrained on a large-scale dataset constructed for image classification [41]. Although [24] proposed an online learning method based on a pool of CNNs, its accuracy is not particularly better than that of the methods based on hand-crafted features. A recent approach [29] trained a multi-domain neural network and achieved great performance improvement. Unfortunately, existing DNN-based trackers are either very slow at test time due to intensive computational overhead, or they don't effectively leverage cross-frame information in training videos and thus result in bad performance. Our tracker explores temporal information through an RNN and can run very efficiently by performing only a forward pass at run time.

2.3. Recurrent Neural Networks

RNNs are powerful learning models that are often used for learning sequential tasks. Advances in understanding the learning dynamics of RNNs have enabled their successful applications in a wide range of tasks [17, 30, 13, 35, 5, 34, 20, 8, 27]. The standard RNN cell consists of an input, a hidden and an output layer as illustrated in Figure 2.

Besides traditional RNN, the application of recurrent networks with sophisticated hidden units, such as Long Short-Term Memory (LSTM) [17] or Gated Recurrent Unit (GRU) [5], has become common in recent years. Mnih *et al.* [27] applied LSTM in visual attention problems by training a decision making network, which sufficiently improves computational efficiency for computer vision tasks. In recent works, Donahue *et al.* [8] proposed a class of recurrent convolutional architectures which are suitable for

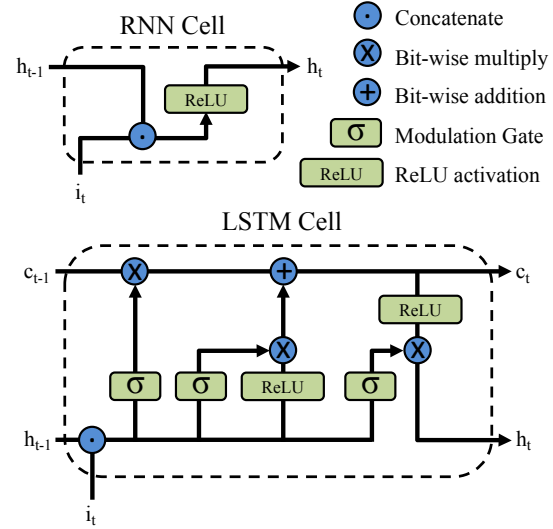


Figure 2: A diagram of standard RNN cell and LSTM cell, and ReLU activation is used here for image processing. The modulation gates in LSTM allow it to better handle long-term and short-term connections. In both cells, the input i_t is integrated with previous hidden state h_{t-1} to generate current hidden state h_t .

large-scale visual understanding tasks, and demonstrated the value of these models for activity recognition, image captioning, and video description. The LSTM cell used in this paper is shown in Figure 2.

2.4. Deep Reinforcement Learning

RL is a learning method based on trial and error, where an agent does not necessarily have *a priori* knowledge about which is the correct action to take. The underlying model that RL learns is a Markov Decision Process (MDP): An agent interacts with the environment and selects an action. Applying the action at a single state, and the environment emits a new state and a reward signal. In order to maximize the expected rewards in long term, the agent learns the best policy to take actions [36].

Some works in computer vision literature and elsewhere e.g., [1, 3, 7, 31, 28, 27] have embraced vision as a sequential decision task. RL is capable of solving sequential decision making tasks especially with the help of good environment representations learned by DNN models. Some recent works included training an RL agent to play Atari games from raw input pixels [28], and a recurrent visual attention model [27] to reduce computation for computer vision tasks. Our work is similar to the attention model described in [27], but we designed our own network architecture especially for solving the visual tracking problem by combining CNN, RNN and RL algorithms.

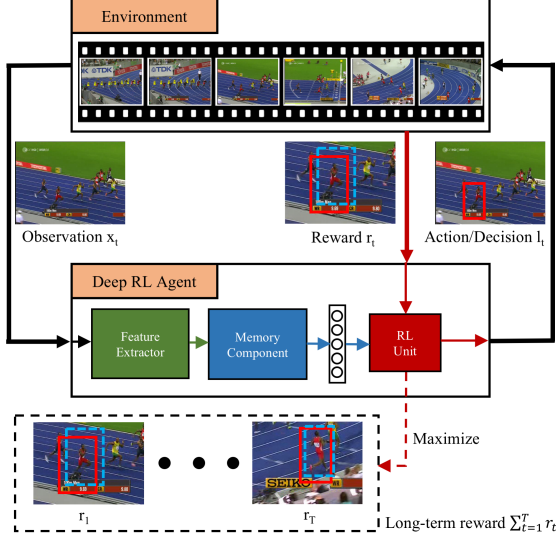


Figure 3: A general work flow of the deep RL agent: At each time step t , the agent observes the environment as x_t , takes an action l_t on where to put the bounding box and receives a scalar reward r_t from the environment which measures the goodness of the tracking performance. During training, parameters of the agent will be updated to maximize the long-term tracking reward $\sum_{t=1}^T r_t$.

Our work formulates visual tracking as a sequential decision making process and leverages RL to solve this problem. We describe an end-to-end optimization procedure that allows the model to be trained directly with respect to a visual tracking task and to maximize a long-term tracking performance measure which depends on the entire sequence of predictions made by the model. We will describe our framework in detail in next section.

3. Deep RL Agent

In this paper we formulate the tracking problem as the sequential decision process of a goal-directed agent interacting with the visual environment: As shown in Figure 3, the visual environment is a sequence of video frames, the agent is modeled as our neural network, and the decision or action is to decide the location of target object. At each point in time, the agent observes the environment and extracts representative features from observations. Since the agent only observes one image frame at a time, it needs to integrate information over time in order to determine how to take actions most effectively. At each step, the agent receives a scalar reward which depends on its actions, and the goal of the agent is to maximize the total long-term rewards. In our system, the agent is built around an RNN as shown in Figure 1. At each time step, it processes the input frame,

integrates information over time, and chooses how to act correspondingly.

In more detail, the proposed deep RL agent consists of three components: the feature extractor (processes input frame and generates feature representations), the memory component (integrates information overtime and formulates and updates the internal state) and the RL unit (chooses how to act and maximizes the overall reward). These are described below.

3.1. Feature Extractor

At each time step t the agent observes the environment in the form of an image x_t , on which the feature extractor computes a representation. This representation can be as simple as a linear transformation or a more sophisticated feature computed by a CNN, and is used by other components to integrate sequential information and take appropriate actions. The extracted features are denoted as:

$$i_t = f_i(x_t; W_i) \quad (1)$$

where i_t is a function of x_t given parameters W_i , and the functional relationship is defined by f_i .

The same feature extractor will be applied to every single frame in an video, thus the feature extractor has to be generic enough to capture representative features for different objects in different contexts. Given varying tracking tasks, the feature extractor can either be pre-trained, transferred from other models, or initialized and trained from scratch. To make the generic extractor adaptive over time, its parameters W_i will continue to be updated during end-to-end training which provides better descriptors to maximize the overall rewards.

3.2. Memory Component

The memory component maintains an internal memory state which summarizes information extracted from the history of past observations. This hidden state encodes the knowledge of the environment and is fundamental for the agent to take actions.

At one single time step t , the feature vector of the input frame i_t is fed into an RNN. The recurrent neural network updates its internal memory vector h_t based on the previous memory vector h_{t-1} and the current input feature i_t provided by the feature extractor:

$$h_t = f_h(h_{t-1}, i_t; W_h) \quad (2)$$

where f_h is a recurrent transformation function such as GRU, LSTM or a simple logistic function, parameterized by the parameter W_h .

We use the LSTM cell in Figure 2 which has proven to be better at handling long-term and short-term connections, and robust for processing different sequential data.

3.3. RL Unit

The instrumental component is an RL unit taking actions based on the internal memory state and receiving rewards from the environment. At each time step t , the agent performs an action to locate the target object l_t . In this work, the location actions are chosen stochastically from a distribution parameterized by RL unit $f_l(h_t; W_l)$ at time t :

$$l_t \sim p(\cdot | f_l(h_t; W_l)) \quad (3)$$

where p is the probability distribution function determined by $f_l(h_t; W_l)$, the external input is the internal memory state h_t and f_l is the functional relationship defined by the RL unit.

More specifically, the distribution $p(\cdot | f_l(h_t; W_l))$ is defined as a multi-variate Gaussian. The variance of the Gaussian is fixed and decided by our experiments, while the mean value $\mu = f_l(h_t; W_l)$ is the output of the RL unit. Given different tracking tasks, the output of the RL unit can either be the center of a target object (x_t, y_t) , a circle (x_t, y_t, r_t) centered at (x_t, y_t) with radius r_t or a rectangular bounding box (x_t, y_t, w_t, h_t) whose top left corner is (x_t, y_t) and width and height are w_t and h_t respectively.

During training, the agent will receive a reward signal r_t from the environment after executing an action at time t . In this work, the scalar reward signal is defined as follows:

$$r_t = -avg(|l_t - g_t|) - max(|l_t - g_t|) \quad (4)$$

where l_t is the location outputted by RL unit, g_t is the target ground truth at time t , $avg(\cdot)$ computes the mean value of a given matrix and $max(\cdot)$ computes the maximum value.

The training objective is to maximize the sum of the reward signal: $R = \sum_{t=1}^T r_t$. By definition, the reward in Equation 4 measures the closeness between predicted location l_t and ground-truth location g_t . Our training algorithm is designed to maximize R towards 0 such that both average and maximum errors made by the deep RL agent is minimized. We will discuss training details in the next section.

4. Training the Deep RL Agent

Training this network to maximize the overall tracking performance is a non-trivial task, and we leverage the REINFORCE algorithm [42] from the RL community to solve this problem.

4.1. Gradient Approximation

Our deep RL agent is parameterized by $W = \{W_i, W_h, W_l\}$ and we aim to learn these parameters to maximize the total tracking reward the agent can expect when taking different actions. More specifically, the objective of the agent is to learn a policy function $\pi(l_t | z_{1:t}; W)$ with parameters W that, at each step t , maps the history of past interactions with the environment $z_{1:t} =$

$x_1, l_1, \dots, l_{t-1}, x_t$ (a sequence of past observations and actions taken by the agent) to a distribution over actions for the current time step. Here, the policy π is defined by our neural network architecture, and the history of interactions $z_{1:t}$ is summarized in the internal memory state h_t . For simplicity, we will use $Z_t = z_{1:t}$ to indicate all histories up to time t , thus, the policy function can be written as $\pi(l_t | Z_t; W)$.

To put it in a formal way, the policy of deep RL agent $\pi(l_t | Z_t; W)$ induces a distribution over possible interactions Z_t and we aim to maximize the total reward R under this distribution, thus, the objective is defined as:

$$G(W) = E_{p(z_{1:T}; W)} \left[\sum_{t=1}^T r_t \right] = E_{p(Z_T; W)} [R] \quad (5)$$

where $p(z_{1:T}; W)$ is the distribution over possible interactions parameterized by W .

This formulation involves an expectation over high-dimensional interactions which is hard to solve in traditional supervised manner. Here, we bring techniques from the RL community to solve this problem, as shown in [42], the gradient can be first simplified by taking the derivative over log-probability of the policy function π :

$$\nabla_W G = \sum_{t=1}^T E_{p(Z_T; W)} [\nabla_W \ln \pi(l_t | Z_t; W) R] \quad (6)$$

and the expectation can be further approximated by an episodic algorithm: since the action is drawn from probabilistic distributions, one can execute the same policy for many episodes and approximate expectation by taking the average, thus

$$\nabla_W G \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_W \ln \pi(l_t^i | Z_t^i; W) R^i \quad (7)$$

where R^i 's are cumulative rewards obtained by running the current policy π for N episodes, $i = 1 \dots N$.

The above training rule is known as the episodic REINFORCE [42] algorithm, and it involves running the deep RL agent with its current policy to obtain samples of interactions and then updating parameters W of the agent such that the log-probability of chosen actions that have led to high overall rewards is increased.

In practice, although Equation 7 computes a good estimation of the gradient $\nabla_W G$, when applied to train the deep RL agent, the training process is hard to converge due to the high variance of this gradient estimation. Thus, in order to obtain an unbiased low-variance gradient estimation, a common method is to subtract a reinforcement baseline from the cumulative rewards R :

$$\nabla_W G \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_W \ln \pi(l_t^i | Z_t^i; W) (R_t^i - b_t) \quad (8)$$

where b_t is called reinforcement baseline in the RL literature, it is natural to select $b_t = E_\pi[R_t]$, and this form of baseline is known as the value function [37]. This estimation maintains the same expectation with Equation 7 while sufficiently reduces the variance.

4.2. Training with Backpropagation

The only remaining part to compute the gradient in Equation 8 is to compute the gradient over log-probability of the policy function $\nabla_W \ln \pi(l_t|Z_t; W)$. To simplify notation, we focus on one single time step and omit usual unit index subscript throughout. In our network design, the policy function π outputs the target location l which is drawn from a Gaussian distribution centered at μ with fixed variance σ , and μ is the output of the deep RL agent parameterized by W . The density function g determining the output l on any single trial is given by:

$$g(l, \mu, \delta) = \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} e^{-\frac{(l-\mu)^2}{2\sigma^2}} \quad (9)$$

Based on REINFORCE algorithm [42], the gradient of the policy function with respect to μ is given by the gradient of the density function:

$$\nabla_\mu \ln \pi = \frac{\partial \ln g}{\partial \mu} = \frac{l - \mu}{\sigma^2} \quad (10)$$

since μ is the output of deep RL agent parameterized by W , the gradients with respect to network weights W can be easily computed by standard backpropagation.

4.3. Overall Procedure

The overall procedure of our training algorithm is presented in Algorithm 1. the network parameters W are first randomly initialized to define our initial policy. Then, we take first T frames from one training video to be the input of our network. We execute current policy N times, compute gradients and update network parameters. Next, we take consecutive T frames from the same video and apply the same training procedure. We repeat this for all training videos in our dataset, and we stop when we reach the maximum number of epochs or the cumulative reward ceases to increase.

During testing, the network parameters W are fixed and no online fine-tuning is needed. The procedure at test time is as simple as computing one forward pass of our deep RL agent, *i.e.*, given a test video, the deep RL agent predicts the location of target object in every single frame by sequentially processing the video data.

5. Experiments

We evaluated the proposed approach of visual object tracking on artificially generated datasets. We varied the

Algorithm 1 Deep RL agent training algorithm

Input: Training videos $\{v_1, \dots, v_M\}$ with ground-truth

Output: Network weights W

- 1: Randomly initialize weights W_i, W_h and W_l
 - 2: Start from the first frame in training dataset
 - 3: **repeat**
 - 4: Sequentially select T frames $\{x_1, \dots, x_T\}$
 - 5: Extract features $\{i_1, \dots, i_T\}$
 - 6: Generate internal memory states $\{h_1, \dots, h_T\}$
 - 7: Compute network output $\{\mu_1, \dots, \mu_T\}$
 - 8: Randomly sample predictions for N episodes $\{l_1^{1:N}, \dots, l_T^{1:N}\}$ according to Equation 9
 - 9: Calculates rewards $\{r_1^{1:N}, \dots, r_T^{1:N}\}$ based on Equation 4
 - 10: $b_t \leftarrow \frac{1}{N} \sum_{i=1}^N r_t^i, t = 1, \dots, T$
 - 11: Computes gradient according to Equation 8
 - 12: Update W using backpropagation
 - 13: Move on to next T frames
 - 14: **until** reward doesn't increase
-

configurations of generating these datasets to validate the performance of our proposed tracker in different scenarios.

5.1. Evaluation Metrics

We used two performance measures to quantitatively evaluate our tracking model on test data. The first one was the average central pixel error

$$e = \sqrt{(x_{pred} - x_{gt})^2 + (y_{pred} - y_{gt})^2} \quad (11)$$

where (x_{gt}, y_{gt}) and (x_{pred}, y_{pred}) were the ground truth and predicted central pixel coordinate in a 2D image, and the error was computed as the pixel distance between these two points. The error was then averaged among all predictions in test data.

Average central pixel error measured how far away the center of prediction deviated from the center of ground truth, but it didn't measure the scale difference between the prediction and the ground truth. Here, we used the average Intersection-over-Union (IoU)

$$IoU = \frac{|b_{pred} \cap b_{gt}|}{|b_{pred} \cup b_{gt}|} \quad (12)$$

where b_{gt} and b_{pred} were the ground truth and predicted bounding boxes, and IoU was computed as the intersection area divided by union area between b_{gt} and b_{pred} . Commonly, b_{gt} and b_{pred} were defined as rectangles.

5.2. Implementation Details

We generated videos from MNIST images of handwritten digits [23] by placing randomly-drawn digits in a larger

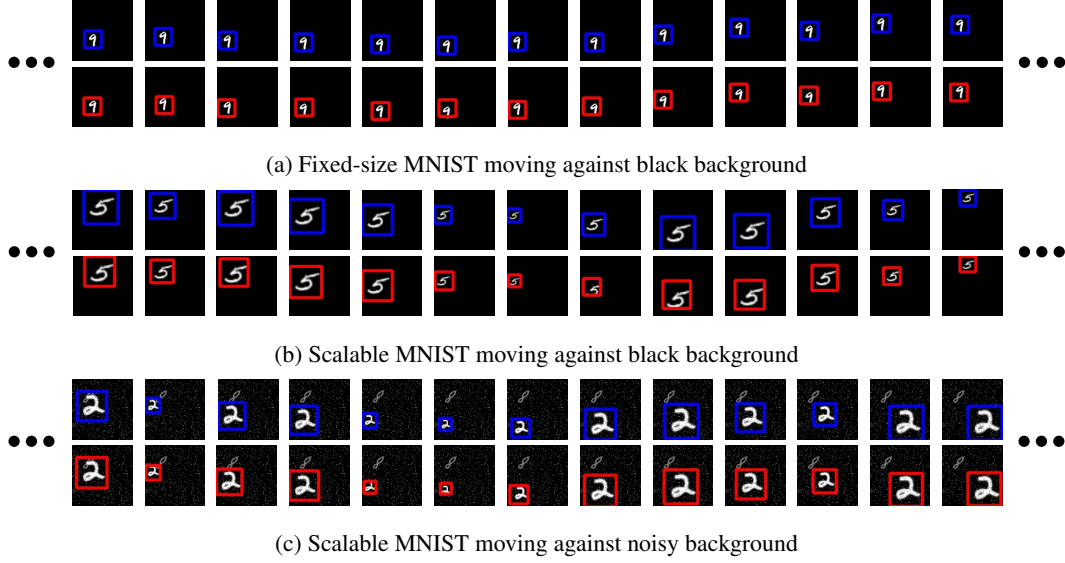


Figure 4: Qualitative results on moving MNIST datasets. The figure shows one test sequence for three different experimental settings. In each experiment, the blue rectangle (first row) indicates the location of ground truth and the red rectangle (second row) indicates the location of our predictions.

Experiment	Average central error (pixels)	Average IoU (%)
Single-digit in RATM [20]	-	63.53
Fixed-size moving MNIST in black background	4.13	72.14
Scalable moving MNIST in black background	3.57	75.85
Scalable moving MNIST in noisy background	6.54	64.29

Table 1: Evaluation measurements on test data.

100×100 canvas and moving the digits from one frame to the next. We respected the same data distribution for training and testing as in the original MNIST dataset, i.e., digits were drawn from the training split to generate the training sequences and from the test split to generate the test sequences. Our tracker demonstrated good performance in challenging scenarios and achieved better results comparing with [20] in the same experimental settings.

The CNN, RNN and RL parts were all implemented using the TensorFlow toolbox¹. We first described the design choices that were common to all our experiments:

Feature extractor: In order to capture a robust representation against small variations in the MNIST dataset, a shallow hierarchical CNN was used as the feature extractor. The CNN structure had two convolutional layers with filter sizes of $8 \times 5 \times 5$ and $16 \times 5 \times 5$, each followed by a 2×2 maxpooling layer and the ReLU activation function.

Memory component: An LSTM was used here as the memory component. The LSTM structure was similar to the one proposed in [17] with all gated controls except that

¹<https://www.tensorflow.org/>

the activation function was replaced with ReLU. The LSTM used in our network is shown in Figure 2. The LSTM state vector h had dimensionality 256.

RL unit: This component was implemented as a 256-neuron hidden layer with ReLU activation function which took state vector h as input and generated the mean value μ of the location policy l . In different experimental settings, the network output μ could either be the central target location (x, y) or the target square (x, y, w) . The location policy was sampled from a Gaussian distribution with mean μ and variance σ during training phase, and we found that $\sigma = 1$ was good for both randomness and certainty in our experiment. During test phase, we directly used the output mean value μ as prediction which was the same as setting $\sigma = 0$.

5.3. Evaluation on Moving MNIST

Fixed-size moving MNIST against a black background: We first tested the ability of our deep RL agent to find good tracking policies for fixed-size MNIST digits moving in black background. In this experiment, we generated videos, each with a single 28×28 MNIST digit

moving in a random walk manner with momentum. The dataset consisted of 500 training sequences and 100 test sequences. Each testing video sequence contained 100 frames while each training sequence only contained 20 frames. The training algorithm was the same as Algorithm 1, we used $T = 10$ and $N = 5$ as these hyper-parameters provided the best tracking performance. The initial learning rate was 0.00001 and we annealed the learning rate linearly from its initial value to 0 over the course of training. We trained the model up to 2000 epochs, or until the cumulative tracking reward R stopped increasing. The trained model was directly applied to the test sequences with no online fine-tuning. Since the scale of MNIST digits didn't change in this experiment, the output of deep RL agent had two neurons directly encoding the center location (x, y) of target object. The qualitative result is shown in Figure 4a, and the second row of Table 1 shows the evaluation result.

Scalable moving MNIST against a black background: The second problem we considered was tracking a scalable MNIST digit moving in a black background, and we wanted to see whether our model was capable of predicting scale changes. In this experiment, we generated new sequences for training and testing. We used the same script as in the previous experiment to generate this dataset except that the 28×28 MNIST digits were replaced by scalable MNIST digits: the side length of a MNIST digit was defined as $s_t = ks_{t-1}$ where k was a random number uniformly sampled from $(-0.5, 2)$, and we limited the side length s in range $(14, 56)$ to make sure that the digit did not become either too large or too small. The network architectures were slightly changed such that there were 3 output neurons which encoded not only the center location (x, y) but also the side length w , and we didn't modify the CNN and RNN architecture. The network was randomly initialized and trained from scratch following the same training and testing procedures as discussed in the previous experiment. Figure 4b shows the tracking result, and the third row of Table 1 presents the quantitative evaluation metrics.

Scalable moving MNIST against a noisy background: One of the most challenging aspects of object tracking is the presence of cluttered and noisy background. It is interesting to investigate how robust our tracker is in this scenario. In this experiment, we leveraged the same dataset in the second experiment, but added random background noise for each training and testing video. Technically, we added white noise dots by randomly selecting 500 points in the 100×100 canvas, and we also added one MNIST digit as distraction in the background. We used the same network design as in the second experiment and followed the same training procedure to train our network from scratch. Our testing results showed that our tracker was robust enough to not only predict center locations but also kept track of scale changes. The qualitative and quantitative results are shown

in Figure 4c and fourth row of Table 1 respectively.

5.4. Discussion

We show that the proposed tracking approach works well in three different scenarios by performing rigorous validation experiments on artificial video sequences with ground truth. We have found in our experiments that the proposed framework possesses:

- **Generalizability:** In all experiments, we evaluate a trained model on test sequences that are longer than the training sequences, and our model has proven to be able to generalize to longer test sequences. This is not only because of we apply RNN in our design, but also due to the RL algorithm that allows us to take long-term rewards into consideration.
- **Adaptability:** By formulating tracking as a decision making process, the number of neurons in the output layer entirely depends on tracking complexity. We have shown that our model can achieve reasonable performance by changing only the last-layer architecture and keeping other components untouched. This is a desirable feature for machine learning models since a general model can be applied to different tasks with minimal alteration.
- **Robustness:** The features are captured by the CNN and RNN in our tracking model, and our experiments have shown that by leveraging the same CNN and RNN architectures, our model is not only able to generate more complex predictions but is also robust to background noise. This shows that our network design is robust and informative in terms of learning good tracking features.

We will extend our framework for real-world tracking datasets in our future work.

6. Conclusion

In this paper, we proposed a novel neural network tracking model based on a recurrent convolutional network trained with deep RL algorithm. To the best of our knowledge, we are the first to bring RL into CNN and RNN to solve visual tracking problems. The entire network is end-to-end trainable offline, and the full pipeline is jointly tuned to maximize the tracking quality. The deep RL algorithm directly optimizes a long-term tracking performance measure which depends on the whole tracking video sequence. The proposed deep RL agent maps raw image pixels to target location decisions, and the spatial temporal representations learned by our model are general and robust to accommodate different tracking scenarios. Extensive experiments have validated the applicability of our proposed tracker.

References

- [1] B. Alexe, N. Heess, Y. W. Teh, and V. Ferrari. Searching for objects driven by context. In *Advances in Neural Information Processing Systems*, pages 881–889, 2012. 3
- [2] B. Babenko, M.-H. Yang, and S. Belongie. Robust object tracking with online multiple instance learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1619–1632, 2011. 3
- [3] N. J. Butko and J. R. Movellan. Optimal scanning for faster object detection. In *Computer vision and pattern recognition, 2009. cvpr 2009. ieee conference on*, pages 2751–2758. IEEE, 2009. 3
- [4] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. *arXiv preprint arXiv:1405.3531*, 2014. 1, 3
- [5] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014. 3
- [6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009. 3
- [7] M. Denil, L. Bazzani, H. Larochelle, and N. de Freitas. Learning where to attend with deep architectures for image tracking. *Neural computation*, 24(8):2151–2184, 2012. 3
- [8] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2625–2634, 2015. 3
- [9] J. Fan, W. Xu, Y. Wu, and Y. Gong. Human tracking using convolutional neural networks. *IEEE Transactions on Neural Networks*, 21(10):1610–1623, 2010. 3
- [10] 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. 1, 3
- [11] H. Grabner, M. Grabner, and H. Bischof. Real-time tracking via on-line boosting. In *BMVC*, volume 1, page 6, 2006. 2
- [12] H. Grabner, C. Leistner, and H. Bischof. Semi-supervised on-line boosting for robust tracking. In *European conference on computer vision*, pages 234–247. Springer, 2008. 2
- [13] A. Graves, A.-r. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6645–6649. IEEE, 2013. 3
- [14] S. Hare, A. Saffari, and P. H. Torr. Struck: Structured output tracking with kernels. In *2011 International Conference on Computer Vision*, pages 263–270. IEEE, 2011. 2
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 1, 3
- [16] D. Held, S. Thrun, and S. Savarese. Learning to track at 100 fps with deep regression networks. *arXiv preprint arXiv:1604.01802*, 2016. 3
- [17] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. 3, 7
- [18] S. Hong, T. You, S. Kwak, and B. Han. Online tracking by learning discriminative saliency map with convolutional neural network. *arXiv preprint arXiv:1502.06796*, 2015. 3
- [19] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. *arXiv preprint arXiv:1603.08155*, 2016. 1, 3
- [20] S. E. Kahou, V. Michalski, and R. Memisevic. Ratm: Recurrent attentive tracking model. *arXiv preprint arXiv:1510.08660*, 2015. 3, 7
- [21] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, L. Cehovin, G. Fernandez, T. Vojir, G. Hager, G. Nebehay, and R. Pflugfelder. The visual object tracking vot2015 challenge results. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1–23, 2015. 2
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 3
- [23] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 6
- [24] H. Li, Y. Li, and F. Porikli. Deeptack: Learning discriminative feature representations online for robust visual tracking. *IEEE Transactions on Image Processing*, 25(4):1834–1848. 3
- [25] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 1, 3
- [26] X. Mei and H. Ling. Robust visual tracking using ℓ_1 minimization. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1436–1443. IEEE, 2009. 2
- [27] V. Mnih, N. Heess, A. Graves, et al. Recurrent models of visual attention. In *Advances in Neural Information Processing Systems*, pages 2204–2212, 2014. 3
- [28] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013. 3
- [29] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking. *arXiv preprint arXiv:1510.07945*, 2015. 1, 3
- [30] R. Pascanu, T. Mikolov, and Y. Bengio. On the difficulty of training recurrent neural networks. *ICML (3)*, 28:1310–1318, 2013. 3
- [31] M. Ranzato. On learning where to look. *arXiv preprint arXiv:1405.5488*, 2014. 3
- [32] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang. Incremental learning for robust visual tracking. *International Journal of Computer Vision*, 77(1-3):125–141, 2008. 2
- [33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1, 3

- [34] N. Srivastava, E. Mansimov, and R. Salakhutdinov. Unsupervised learning of video representations using lstms. *CoRR, abs/1502.04681*, 2, 2015. 3
- [35] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014. 3
- [36] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*, volume 1. 3
- [37] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour, et al. Policy gradient methods for reinforcement learning with function approximation. 6
- [38] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014. 1, 3
- [39] A. Toshev and C. Szegedy. Deeppose: Human pose estimation via deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1653–1660, 2014. 1, 3
- [40] L. Wang, W. Ouyang, X. Wang, and H. Lu. Visual tracking with fully convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3119–3127, 2015. 1, 3
- [41] N. Wang, S. Li, A. Gupta, and D.-Y. Yeung. Transferring rich feature hierarchies for robust visual tracking. *arXiv preprint arXiv:1501.04587*, 2015. 1, 3
- [42] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992. 2, 5, 6
- [43] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013. 2
- [44] T. Zhang, B. Ghanem, S. Liu, and N. Ahuja. Robust visual tracking via multi-task sparse learning. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2042–2049. IEEE, 2012. 2