

# RATM: Recurrent Attentive Tracking Model

Samira Ebrahimi Kahou, Vincent Michalski, and Roland Memisevic

**Abstract**—We present an attention-based modular neural framework for computer vision. The framework uses a soft attention mechanism allowing models to be trained with gradient descent. It consists of three modules: a recurrent attention module controlling *where* to look in an image or video frame, a feature-extraction module providing a representation of *what* is seen, and an objective module formalizing *why* the model learns its attentive behavior. The attention module allows the model to focus computation on task-related information in the input. We apply the framework to several object tracking tasks and explore various design choices. We experiment with three data sets, bouncing ball, moving digits and the real-world KTH data set. The proposed Recurrent Attentional Tracking Model (RATM) performs well on all three tasks and can generalize to related but previously unseen sequences from a challenging tracking data set.

**Index Terms**—computer vision, deep learning, object tracking, visual attention

## 1 INTRODUCTION

ATTENTION MECHANISMS are one of the biggest trends in deep-learning research and have been successfully applied in a variety of neural-network architectures across different tasks. In computer vision, for instance, attention mechanisms have been used for image generation [1] and image captioning [2]. In natural language processing they have been used for machine translation [3] and sentence summarization [4]. And in computational biology attention was used for subcellular protein localization [5].

In these kinds of applications usually not all information contained in the input data is relevant for the given task. Attention mechanisms allow the neural network to focus on the relevant parts of the input, while ignoring other, potentially distracting, information. Besides enabling models to ignore distracting information, attention mechanisms can be helpful in streaming data scenarios, where the amount of data per frame can be prohibitively large for full processing. In addition, some studies suggest that there is a representational advantage of sequential processing of image parts over a single pass over the whole image (see for example [1], [6], [7], [8], [9], [10]).

Recently, [1] introduced the Deep Recurrent Attentive Writer (DRAW), which involves a Recurrent Neural Network (RNN) that controls a read and a write mechanism based on attention. The read mechanism extracts a parametrized window from the static input image. Similarly, the write mechanism is used to write into a window on an output canvas. This model is trained to sequentially produce a reconstruction of the input image on the canvas. Interestingly, one of the experiments on handwritten digits showed that the read mechanism learns to trace digit contours and the write mechanism generates digits in a continuous motion. This observation hints at the potential of such mechanisms in visual object tracking applications, where the primary goal is to trace the spatio-temporal “contours” of an object as it moves in a video.

Previous work on the application of attention mechanisms for tracking includes [8] and references therein. In contrast to that line of work, we propose a model based on a fully-integrated neural framework, that can be trained end-

to-end using back-propagation. The framework consists of three modules: a recurrent differentiable attention module controlling *where* to look in an image, a feature-extraction module providing a representation of *what* is seen, and an objective module formalizing *why* the model learns its attentive behavior. As we shall show, a suitable surrogate cost in the objective module can provide a supervised learning signal, that allows us to train the network end-to-end, and to learn attentional strategies using simple supervised back-prop without resorting to reinforcement learning or sampling methods.

According to a recent survey of tracking methods [11], many approaches to visual tracking involve a search over multiple window candidates based on a similarity measure in a feature space. Successful methods involving deep learning, such as [12], perform tracking-by-detection, e.g. by using a Convolutional Neural Network (CNN) for foreground-background classification of region proposals. As in most approaches, the method in [12] at each time step samples a number of region proposals (256) from a Gaussian distribution centered on the region of the previous frame. Such methods do not benefit from useful correlations between the target location and the object’s past trajectory. There are deep-learning approaches that consider trajectories by employing particle filters such as [13], which still involves ranking of region proposals (1,000 particles).

In our Recurrent Attentional Tracking Model (RATM), an RNN predicts the position of an object at time  $t$ , given a real-valued hidden state vector. The state vector can summarize the history of observations and predictions of previous time steps. We rely on a *single* prediction per time step instead of using the predicted location as basis for a search over multiple region proposals. This allows for easy integration of our framework’s components and training with simple gradient-based methods.

The main contribution of our work is the introduction of a modular neural framework, that can be trained end-to-end with gradient-based learning methods. Using object tracking as an example application, we explore different settings and provide insights into model design and training. While

the proposed framework is targeted primarily at videos, it can also be applied to sequential processing of still images.

## 2 RECURRENT NEURAL NETWORKS

Recurrent Neural Networks (RNNs) are powerful machine learning models that are used for learning in sequential processing tasks. Advances in understanding the learning dynamics of RNNs enabled their successful application in a wide range of tasks (for example [14], [15], [16], [17], [18], [19]). The standard RNN model consists of an input, a hidden and an output layer as illustrated in Figure 1.

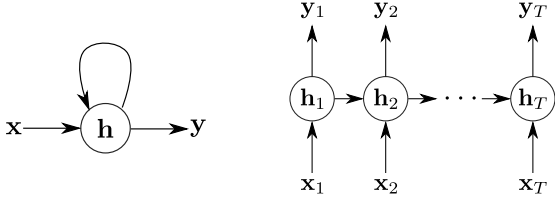


Fig. 1: Left: A simple recurrent network with one input, one output and a hidden layer, which has a feed-back loop connection. Right: The same network unrolled in time for  $T$  time steps.

In each time step  $t$ , the network computes a new hidden state  $\mathbf{h}_t$  based on the previous state  $\mathbf{h}_{t-1}$  and the input  $\mathbf{x}_t$ :

$$\mathbf{h}_t = \sigma(\mathbf{W}_{in}\mathbf{x}_t + \mathbf{W}_{rec}\mathbf{h}_{t-1}), \quad (1)$$

where  $\sigma$  is a non-linear activation function,  $\mathbf{W}_{in}$  is the matrix containing the input-to-hidden weights and  $\mathbf{W}_{rec}$  is the recurrent weight matrix from the hidden layer to itself. At each time step the RNN also generates an output

$$\mathbf{y}_t = \mathbf{W}_{out}\mathbf{h}_t + \mathbf{b}_y, \quad (2)$$

where  $\mathbf{W}_{out}$  is the matrix with weights from the hidden to the output layer.

Although the application of recurrent networks with sophisticated hidden units, such as Long Short-Term Memory (LSTM) [14] or Gated Recurrent Unit (GRU) [18], has become common in recent years (for example [3], [17], [19]), we rely on the simple IRNN proposed by [20], and show that it works well in the context of visual attention. The IRNN corresponds to a standard RNN, where recurrent weights  $\mathbf{W}_{rec}$  are initialized with a scaled version of the identity matrix and the hidden activation function  $\sigma(\cdot)$  is the element-wise Rectified Linear Unit (ReLU) function [21]

$$\sigma(x) = \max(0, x). \quad (3)$$

The initial hidden state  $\mathbf{h}_0$  is initialized as the zero vector. Our experiments are based on the Theano [22], [23] implementation of the IRNN shown to work well for video in [24].

## 3 NEURAL ATTENTION MECHANISMS

Our attention mechanism is a modification of the read mechanism introduced in [1]. It extracts *glimpses* from the input image by applying a grid of two-dimensional Gaussian window filters. Each of the filter responses corresponds to one pixel of the glimpse. An example of the glimpse extraction is shown in Figure 2.

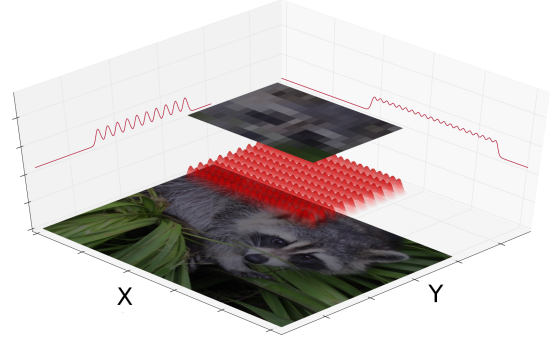


Fig. 2: A  $20 \times 10$  glimpse is extracted from the full image by applying a grid of  $20 \times 10$  two-dimensional Gaussian window filters. The separability of the multi-dimensional Gaussian window allows for efficient computation of the extracted glimpse.

Given an image  $\mathbf{x}$  with  $A$  columns and  $B$  rows, the attention mechanism separately applies a set of  $M$  column filters  $\mathbf{F}_X \in \mathbb{R}^{M \times A}$  and a set of  $N$  row filters  $\mathbf{F}_Y \in \mathbb{R}^{N \times B}$ , extracting an  $M \times N$  glimpse  $\mathbf{p}$ :

$$\mathbf{p} = \mathbf{F}_Y \mathbf{x} \mathbf{F}_X^T. \quad (4)$$

This implicitly computes  $M \times N$  two-dimensional filter responses due to the separability of two-dimensional Gaussian filters. For multi-channel images the same filters are applied to each channel separately.

The sets of one-dimensional row ( $\mathbf{F}_Y$ ) and column ( $\mathbf{F}_X$ ) filters have three parameters each<sup>1</sup>:

- the grid center coordinates  $g_X, g_Y$ ,
- the standard deviation for each axis  $\sigma_X, \sigma_Y$  and
- the stride between grid points on each axis  $\delta_X, \delta_Y$ .

These parameters are dynamically computed as an affine transformation of a vector of activations  $\mathbf{h}$  from a neural network layer:

$$(\tilde{g}_X, \tilde{g}_Y, \tilde{\sigma}_X, \tilde{\sigma}_Y, \tilde{\delta}_X, \tilde{\delta}_Y) = \mathbf{W}\mathbf{h} + \mathbf{b}, \quad (5)$$

where  $\mathbf{W}$  is the transformation matrix and  $\mathbf{b}$  is the bias. This is followed by normalization of the parameters:

$$g_X = \frac{\tilde{g}_X + 1}{2}, \quad g_Y = \frac{\tilde{g}_Y + 1}{2}, \quad (6)$$

$$\delta_X = \frac{A-1}{M-1} \cdot |\tilde{\delta}_X|, \quad \delta_Y = \frac{B-1}{N-1} \cdot |\tilde{\delta}_Y|, \quad (7)$$

$$\sigma_X = |\tilde{\sigma}_X|, \quad \sigma_Y = |\tilde{\sigma}_Y|. \quad (8)$$

The mean coordinates  $\mu_X^i, \mu_Y^j$  of the Gaussian filter at column  $i$ , row  $j$  in the attention grid are computed as follows:

$$\mu_X^i = g_X + (i - \frac{M}{2} - 0.5) \cdot \delta_X, \quad (9)$$

$$\mu_Y^j = g_Y + (j - \frac{N}{2} - 0.5) \cdot \delta_Y \quad (10)$$

1. The original read mechanism in [1] also adds a scalar intensity parameter  $\gamma$ , that is multiplied to filter responses.

Finally, the filter banks  $\mathbf{F}_X$  and  $\mathbf{F}_Y$  are defined by:

$$\mathbf{F}_X[i, a] = \exp\left(-\frac{(a - \mu_X^i)^2}{2\sigma^2}\right), \quad (11)$$

$$\mathbf{F}_Y[j, b] = \exp\left(-\frac{(b - \mu_Y^j)^2}{2\sigma^2}\right) \quad (12)$$

The filters (rows of  $\mathbf{F}_X$  and  $\mathbf{F}_Y$ ) are later normalized to sum to one.

Our read mechanism makes the following modifications to the DRAW read mechanism [1]:

- We allow rectangular (not only square) attention grids and we use separate strides and standard deviations for the  $X$  and  $Y$ -axis. This allows the model to stretch and smooth the glimpse content to correct for distortions introduced by ignoring the original aspect ratio of an input image.
- We use the absolute value function  $abs(x) = |x|$  instead of  $\exp(x)$  to ensure positivity of strides and standard deviations (see Equations 7 and 8). The motivation for this modification is that in our experiments we observed stride and standard deviation parameters to often saturate at low values, causing the attention window to zoom in on a single pixel. This effectively inhibits gradient flow through neighboring pixels of the attention filters. Piecewise linear activation functions have been shown to benefit optimization [21] and the absolute value function is a convenient trade-off between the harsh zeroing of all negative inputs of the ReLU and the extreme saturation for highly negative inputs of the exponential function.
- We drop the additional scalar intensity parameter  $\gamma$ , because we did not observe it to influence the performance in our experiments.

## 4 A MODULAR FRAMEWORK FOR VISION

The proposed modular framework for an attention-based approach to computer vision consists of three components: the attention module (controlling *where* to look), the feature-extraction module (providing a representation of *what* is seen) and the objective module (formalizing *why* the model is learning its attentive behavior). An example architecture for tracking using these modules is described in Section 5.

### 4.1 Feature-extraction module

The feature-extraction module computes a feature representation of a given input glimpse. This representation can be as simple as the identity transformation, i.e. the original pixel representation, or a more sophisticated feature extractor, e.g. an CNN. The extracted features are used by other modules to reason about the visual input. Given a hierarchy of features, such as the activations of layers in an CNN, different features can be passed to the attention and objective modules.

We found that it can be useful to pre-train the feature-extraction module on a large data set, before starting to train the full architecture. After pre-training, the feature extractor’s parameters can either be continued to be updated during end-to-end training, or kept fixed.

Figure 3 shows the symbol used in the following sections to represent a feature-extraction module.

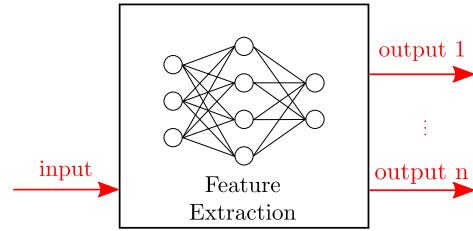


Fig. 3: The symbol for the feature-extraction module. It represents an arbitrary feature extractor, that can have multiple outputs (e.g. for activations from different layers of an CNN).

### 4.2 Attention Module

The attention module is composed of an RNN (see Section 2) and a read mechanism (see Section 3). At each time step, a glimpse is extracted from the current input frame using the attention parameters the RNN predicted in the previous time step (see Section 3). Note that in this context, Equation 5 of the read mechanism corresponds to Equation 2 of the RNN. After the glimpse extraction, the RNN updates its hidden state using the feature representation of the glimpse as input (see Equation 1). Figure 4 shows the symbolic representation used in the following sections to represent the recurrent attention module.

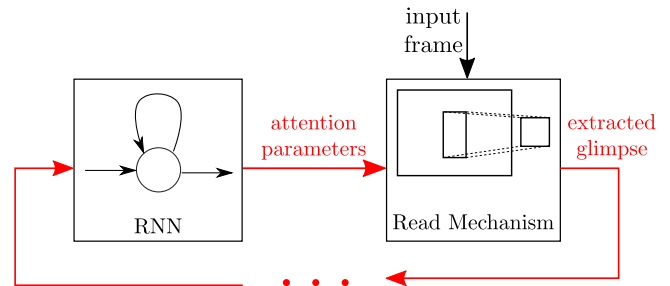


Fig. 4: The symbolic representation of a recurrent attention module, which is composed of an RNN and a read mechanism that extracts a glimpse from the input frame. The extracted glimpse is fed back to the RNN. The dots indicate, that the feed-back connection can involve intermediate processing steps, such as feature extraction.

### 4.3 Objective Module

An objective module guides the model to learn an attentional policy to solve a given task. It outputs a scalar cost, that is computed as function of its target and prediction inputs. There can be multiple objective modules for a single task. A learning algorithm, such as Stochastic Gradient Descent (SGD), uses the sum of cost terms from all objective modules to adapt the parameters of the other modules. Objective modules can receive their input from different parts of the network. For example, if we want to define a penalty between window coordinates, the module would

receive predicted attention parameters from the attention module and target parameters from the trainer.

In all our objective modules we use the Mean Squared Error (MSE) for training:

$$\mathcal{L}_{MSE} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{y}_{target} - \mathbf{y}_{pred}\|_2^2, \quad (13)$$

where  $n$  is the number of training samples,  $\mathbf{y}_{pred}$  is the model’s prediction,  $\mathbf{y}_{target}$  is the target value and  $\|\cdot\|_2^2$  is the squared Euclidean norm. We use the MSE even for classification, as this makes the combination of multiple objectives simpler and worked well. Figure 5 shows the symbol used in the following sections to represent an objective module.

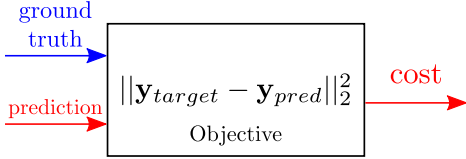


Fig. 5: The symbol for the objective module. It represents the computation of a scalar cost term given prediction and ground truth inputs.

## 5 BUILDING A RECURRENT ATTENTIVE TRACKING MODEL

The task of tracking involves mapping a sequence of input images  $\mathbf{x}_1, \dots, \mathbf{x}_T$  to a sequence of object locations  $\mathbf{y}_1, \dots, \mathbf{y}_T$ . For the prediction  $\hat{\mathbf{y}}_t$  of an object’s location at time  $t$ , the trajectory  $(\mathbf{y}_1, \dots, \mathbf{y}_{t-1})$  usually contains highly relevant contextual information, so it is important to choose a hidden state model which has the capacity to represent this trajectory.

### 5.1 Architecture

At each time step, the recurrent attention module outputs a glimpse from the current input frame using the attention parameters predicted at the previous time step. Optionally, a feature-extraction module extracts a representation of the glimpse and feeds it back to the attention module, which updates its hidden state. The tracking behavior can be learned in various ways:

- One can penalize the difference between the glimpse content and a ground truth image. This can be done in the raw pixel space for simple data sets, which do not show much variation in the objects appearance. This loss is defined as

$$\mathcal{L}_{pixel} = \|\hat{\mathbf{p}} - \mathbf{p}\|_2^2, \quad (14)$$

where  $\hat{\mathbf{p}}$  is the glimpse extracted by the attention mechanism and  $\mathbf{p}$  is the ground truth image. For objects with more variance in appearance, a distance measure between features extracted from the glimpse and from the ground truth image, is more appropriate:

$$\mathcal{L}_{feat} = \|f(\hat{\mathbf{p}}) - f(\mathbf{p})\|_2^2, \quad (15)$$

where  $f(\cdot)$  is the function computed by a feature-extraction module.

- Alternatively, a penalty term can also be defined directly on the attention parameters. For instance, the distance between the center of the ground truth bounding box and the attention mechanism’s  $\hat{\mathbf{g}} = (g_x, g_y)$  parameters can be used as a localization loss

$$\mathcal{L}_{loc} = \|\hat{\mathbf{g}} - \mathbf{g}\|_2^2, \quad (16)$$

We explore several variations of this architecture in Section 6.

## 5.2 Evaluation of Tracking Performance

Tracking models can be evaluated quantitatively on test data using the average Intersection-over-Union (IoU) [25]

$$IoU = \frac{|B_{gt} \cap B_{pred}|}{|B_{gt} \cup B_{pred}|}, \quad (17)$$

where  $B_{gt}$  and  $B_{pred}$  are the ground truth and predicted bounding boxes. A predicted bounding box for RATM is defined as the rectangle between the corner points of the attention grid. This definition of predicted bounding boxes ignores the fact that each point in the glimpse is a weighted sum of pixels around the grid points and the boxes are smaller than the region seen by the attention module. While this might affect the performance under the average IoU metric, the average IoU still serves as a reasonable metric for the soft attention mechanism’s performance in tracking.

## 6 EXPERIMENTAL RESULTS

For an initial study, we use generated data, as described in Sections 6.1 and 6.2, to explore some design choices without being limited by the number of available training sequences. In Section 6.3, we show how one can apply the RATM in a real-world context.

### 6.1 Bouncing Balls

For our initial experiment, we generated videos of a bouncing ball using the script released with [26]. The videos have 32 frames of resolution  $20 \times 20$ . We used 100,000 videos for training and 10,000 for testing. The **attention module**

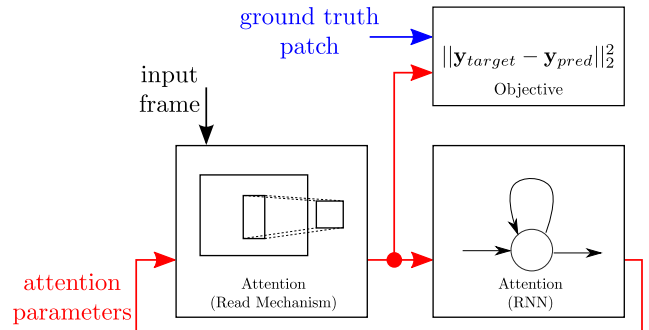


Fig. 6: The architecture used for bouncing balls experiments.

has 64 hidden units in its RNN and its read mechanism extracts glimpses of size  $5 \times 5$ . The attention parameters are initialized to a random glimpse in the first frame. The

input to the RNN are the raw pixels of the glimpse, i.e. the **feature-extraction module** here corresponds to the identity transformation. The **objective module** computes the MSE between the glimpse at the last time step and a target patch, which is simply a cropped white ball image, since shape and color of the object are constant across the whole data set. Figure 6 shows a schematic of this architecture.

For **learning**, we use SGD with a mini-batch size of 16, a learning rate of 0.01 and gradient clipping [15] with a threshold of 1 for 200 epochs. Figure 7 shows results of tracking a ball in a test sequence. RATM is able to learn the correct tracking behaviour only using the penalty on the last frame. We also trained a version with the objective module computing the average MSE between glimpses of all time steps and the target patch. An example tracking sequence of this experiment is shown in Figure 8. The first two rows of Table 1 show the test performance of the model trained with only penalizing the last frame during training. The first row shows the average IoU of the last frame and the second shows the average IoU over all 32 frames of test sequences. The third row shows the average IoU over all frames of the model trained with the penalty on all frames.

The model trained with the penalty at every time step is able to track a bouncing ball for sequences that are much longer than the training sequences. We generated videos that are almost ten times longer (300 frames) and RATM reliably tracks the ball until the last frame. An example is uploaded as part of the supplementary material.

The dynamics in this data-set are rather limited, but as a proof-of-concept they show that the model is able to learn tracking behavior end-to-end. We describe more challenging tasks in the following sections.

## 6.2 MNIST

To increase the difficulty of the tracking task, we move to more challenging data sets, containing more than a single type of object (ten digits), each with variation. We generate videos from  $28 \times 28$  MNIST images of handwritten digits [27] by placing randomly-drawn digits in a larger  $100 \times 100$  canvas with black background and moving the digits from one frame to the next. We respected the same data split for training and testing as in the original MNIST dataset, i.e. digits were drawn from the training split to generate training sequences and from the test split for generation of the test sequences.

Figure 9 shows the schematic of RATM for the MNIST experiments. The **attention module** is similar to the one used in Section 6.1, except that its RNN has 100 hidden units and the size of the glimpse is  $28 \times 28$  (the size of the MNIST images and the CNN input layer).

In the bouncing balls experiment we were able to generate a reliable training signal using pixel-based similarity. However, the variation in the MNIST data set requires a representation that is robust against small variations to guide the training. For this reason, our **feature-extraction module** consists of a (relatively shallow) CNN, that is pre-trained on classification of MNIST digits. Note, that this CNN is only used during training. The CNN structure has two convolutional layers with filter bank sizes of  $32 \times 5 \times 5$ , each followed by a  $2 \times 2$  maxpooling layer, 0.25 dropout [28],

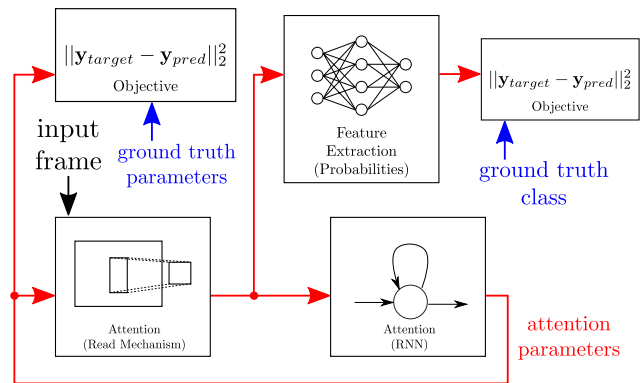


Fig. 9: The architecture used for MNIST experiments.

and the ReLU activation function. These layers are followed by a 10-unit softmax layer for classification. The CNN was trained using SGD with a mini-batch size of 128, a learning rate 0.01, momentum of 0.9 and gradient clipping with a threshold of 5.0 to reach a validation accuracy of 99%.

This CNN classifier is used to extract class probabilities for each glimpse and its parameters remain fixed after pre-training. One of the **objective modules** computes the loss using these probabilities and the target class. Since training did not converge to a useful solution using only this loss, we first introduced an additional objective module penalizing the distances between the upper-left and lower-right bounding-box corners and the corresponding target coordinates. While this also led to unsatisfactory results, we found that replacing the bounding box objective module with one that penalized only grid center coordinates worked well. One possible explanation is, that the grid center penalty does not constrain the stride. Therefore, the glimpse is free to explore without being forced to zoom in. The two penalties on misclassification and on grid center distance, helped the model to reliably find and track the digit. The localization term helped in the early stages of training to guide RATM to track the digits, whereas the classification term encourages the model to properly zoom into the image to maximize classification accuracy.

For **learning** we use SGD with mini-batch size of 32, a learning rate of 0.001, momentum of 0.9 and gradient clipping with a threshold of 1 for 32,000 gradient descent steps.

### 6.2.1 Single-Digit

In the first MNIST experiment, we generate videos, each with a single digit moving in a random walk with momentum. The data set consists of 100,000 training sequences and 10,000 test sequences. The initial glimpse roughly covers the whole frame.

Training is done on sequences with only 10 frames. The classification and localization penalties were applied at every time-step. At test time, the CNN is switched off and we let the model track test sequences of 30 frames. The fourth row of Table 1 shows the average IoU over all frames of the test sequences. Figure 10 shows one test sample.



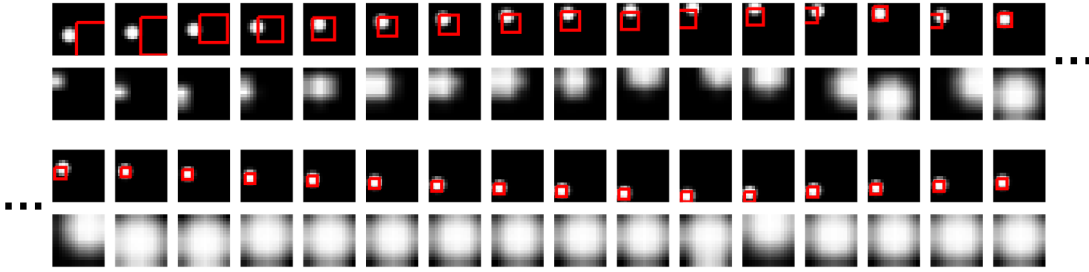


Fig. 7: An example of tracking on the bouncing ball data set. The first row shows the first 16 frames of the sequence with a red rectangle indicating the location of the glimpse. The second row contains the extracted glimpses. The third and fourth row show the continuation of the sequence.

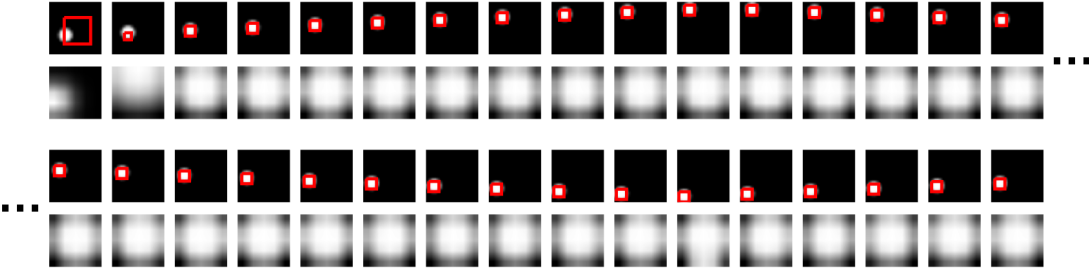


Fig. 8: Tracking result on a test sequence from a model trained with the MSE penalty at every time step. The first row shows the first 16 frames of the sequence with a red rectangle indicating the location of the glimpse. The second row contains the extracted glimpses. The third and fourth row show the continuation of the sequence.

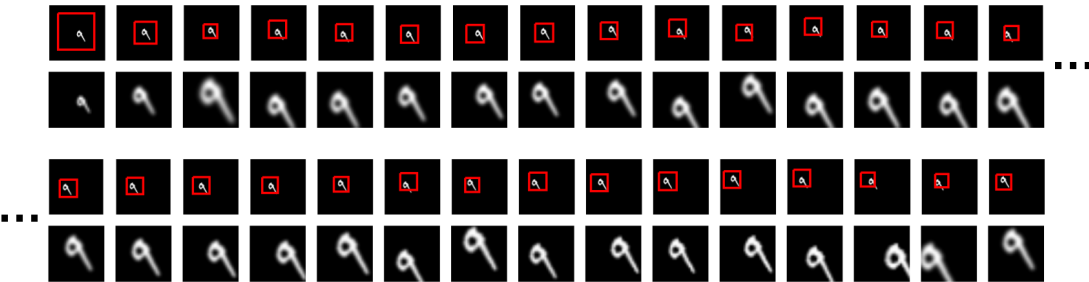


Fig. 10: Tracking one digit. The first and second row show the sequence and corresponding extracted glimpses, respectively. The red rectangle indicates the location of the glimpse in the frame. The third and fourth row are the continuation. Prediction works well far beyond the training horizon of 10 frames.

### 6.2.2 Multi-Digit

It is interesting to investigate how robust RATM is in presence of another moving digit in the background. To this end, we generated new sequences by modifying the bouncing balls script released with [26]. The balls were replaced by randomly drawn MNIST digits. We also added a random walk with momentum to the motion vectors. We generated 100,000 sequences for training and 5,000 for testing.

Here, the bias for attention parameters is not a learn-able parameter. For each video, the bias is set such that the initial glimpse is centered on the digit to be tracked. Width and height are set to about 80% of the frame size. The model was also trained on 10 frame sequences and was able to focus on digits for at least 15 frames on test data. Figure 11 shows tracking results on a test sequence. The fifth row of Table 1 shows the average IoU of all test sequences over 30 frames.

### 6.3 Tracking humans in video

To evaluate the performance on a real-world data set, we train RATM to track humans in the KTH action recognition data set [29], which has a reasonably large number of sequences. We selected the three activity categories, which show considerable motion: walking, running and jogging. We used the bounding boxes provided by [30], which were not hand-labeled and contain noise, such as bounding boxes around the shadow instead of the subject itself.

For the **feature-extraction module** in this experiment, we trained a CNN on binary – human vs. background – classification of  $28 \times 28$  grayscale patches. To generate training data for this CNN, we cropped positive patches from annotated subjects in the ETH pedestrian [31] and INRIA person [32] data sets. Negative patches were sampled from the KITTI detection benchmark [33]. This yielded 21,134 positive and 29,923 negative patches, of which we used 20,000 per class for training. The architecture of the CNN is as follows: two



Fig. 11: Tracking one of two digits. The first and second row show the sequence and corresponding extracted glimpses, respectively. The red rectangle indicates the location of the glimpse in the frame. The third and fourth row are the continuation. Prediction works well for sequences twice as long as the training sequences with 10 frames.

convolutional layers with filter bank sizes  $128 \times 5 \times 5$  and  $64 \times 3 \times 3$ , each followed by  $2 \times 2$  max-pooling and a ReLU activation. After the convolutional layers, we added one fully-connected ReLU-layer with 256 hiddens and the output softmax-layer of size 2. For **pre-training**, we used SGD with a mini-batch size of 64, a learning rate of 0.01, momentum of 0.9 and gradient clipping with a threshold of 1. We performed early stopping with a held-out validation set sampled randomly from the combined data set.

As this real-world data set has more variation than the previous data sets, the **attention module**'s RNN can also benefit from a richer feature representation. Therefore, the ReLU activations of the second convolutional layer of the feature-extraction module are used as input to the attention module. The RNN has 32 hidden units. This low number of hidden units was selected to avoid overfitting, as the number of sequences (1, 200 short sequences) in this data set is much lower than in the synthetic data sets. We initialize the attention parameters for the first time step with the first frame's target window. The initial and target bounding boxes are scaled up by a factor of 1.5 and the predicted bounding boxes are scaled back down with factor  $\frac{1}{1.5}$  for testing. This was necessary, because the training data for the feature-extraction module had significantly larger bounding box annotations.

The inputs to the **objective module** are the ReLU activations of the fully-connected layer, extracted from the predicted window and from the target window. The computed cost is the MSE between the two feature vectors. We also tried using the cosine distance between two feature vectors, but did not observe any improvement in performance. The target window is extracted using the same read mechanism as in the attention module. Simply cropping the target bounding boxes would have yielded local image statistics that are too different from windows extracted using the read mechanism. Figure 12 shows the schematic of the architecture used in this experiment.

For **learning**, we used SGD with a mini-batch size of 16, a learning rate of 0.001 and gradient clipping with a threshold of 1.0. In this experiment we also added a weight-decay regularization term to the cost function that penalizes the sum of the squared Frobenius norms of the RNN weight matrices from the input to the hidden layer and from the hidden layer to the attention parameters. The

squared Frobenius norm is defined as

$$\|\mathbf{A}\|_F^2 = \sum_i^m \sum_j^n |a_{ij}|^2, \quad (18)$$

where  $a_{ij}$  is the element at row  $i$ , column  $j$  in matrix  $\mathbf{A}$ . This regularization term improved the stability during learning. As another stabilization measure, we started training with short five-frame sequences and increased the length of sequences by one frame every 160 gradient descent steps.

For evaluation, we performed a leave-one-subject-out experiment. For each of the 25 subjects in KTH, we used the remaining 24 for training and validation. A validation subject was selected randomly and used for early stopping. The reported number in the sixth row of Table 1 is the IoU on full-length videos of the test subject averaged over all frames of each left-out subject and then averaged over all subjects.

Figure 13 shows examples of test sequences for the classes *jogging* and *walking*. Note, that the region captured by the glimpses is larger than the bounding boxes, because the model internally scales the width and height by factor 1.5 and the Gaussian sampling kernels of the attention mechanism extend beyond the bounding box. An interesting observation is that RATM scales up the noisy initial bounding box in Figure 13 (bottom example), which covers only a small part of the subject. This likely results from pre-training the feature-extraction module on full images of persons. We observed a similar behavior for multiple other samples. Although the evaluation assumes that the target bounding boxes provided by [30] are accurate, RATM is able to recover from such noise.

To show how the model generalizes to unseen videos containing humans, we let it predict some sequences of the TB-100 tracking benchmark [34]. For this experiment, we picked one of the 25 KTH model, that had a reasonably stable learning curve (IoU over epochs). As an example, Figure 14 shows every seventh predicted frame of the *Dancer* sequence and every tenth predicted frame of the sequences *Skater2*, *BlurBody* and *Human2*. For the first two examples, *Dancer* and *Skater2*, RATM tracks the subjects reliably through the whole length of the sequence. This is interesting, as the tracking model was only trained on sequences of up to 30 frames length and the variation in this data is quite different from KTH. The *BlurBody* and *Human2* sequences are more challenging, including extreme

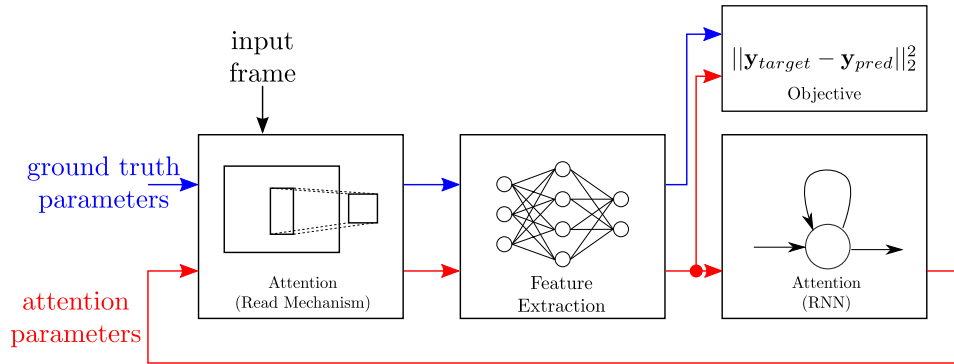


Fig. 12: The architecture used for KTH experiments.



Fig. 13: Two examples of tracking on the KTH data set. The layout for each example is as follows: the first row shows 15 frames of one test sequence with a red rectangle indicating the location of the glimpse. The second row contains the extracted glimpses. The third and fourth row show the continuation of the sequence. We only show every second frame.

camera motion and/or occlusions, causing the model to fail on parts of the sequence. Interestingly in some cases it seems to recover.

In general, the model shows the tendency to grow the window, when it loses a subject. This might be explained by instability of the RNN dynamics and blurry glimpses due to flat Gaussians in the attention mechanism. These challenges will be discussed further in Section 7.

## 7 DISCUSSION

We propose a novel neural framework including a soft attention mechanism for vision, and demonstrate its application to several tracking tasks. Contrary to most existing similar approaches, RATM only processes a small window of each frame. The selection of this window is controlled by a learned attentional behavior. Our experiments explore several design decisions that help overcome challenges associated with adapting the model to new data sets. Several observations in the real-world scenario in Section 6.3, are important for applications of attention mechanisms in computer vision in general:

- The model can be trained on noisy bounding box annotation of videos and at test time recover from

noisy initialization. This might be related to the pre-training of the feature-extraction module on static images. The information about the appearance of humans is transferred to the attention module, which learns to adapt the horizontal and vertical strides among other parameters of the glimpse to match this appearance.

- The trained human tracker seems to generalize to related but more challenging data.

## 8 DIRECTIONS FOR FUTURE RESEARCH

The modular neural architecture is fully differentiable, allowing end-to-end training. End-to-end training allows the discovery of spatio-temporal patterns, which would be hard to learn with separate training of feature extraction and attention modules. In future work we plan to selectively combine multiple data sets from different tasks, e.g. activity recognition, tracking and detection. This makes it possible to benefit from synergies between tasks [35], and can help overcome data set limitations.

One could also try to find alternatives for the chosen modules, e.g. replacing the read mechanism with *spatial transformers* [36]. Spatial transformers offer a more general read mechanism, that can learn to align glimpses using





Fig. 14: Predictions of a KTH model on sequences from the TB-100 benchmark. From top to bottom we show the sequences *Dancer*, *Skater2*, *BlurBody* and *Human2*. To save space, we only show every seventh frame of the *Dancer* predictions and every tenth frame of the other sequences. The layout for each sequence is as follows: The first row shows 15 frames of one test sequence with a red rectangle indicating the location of the predicted glimpse. The second row contains the extracted glimpses. The third and fourth row show the continuation of the sequence.

TABLE 1: Average Intersection-over-Union scores on test data.

Experiment	Average IoU (over # frames)
Bouncing Balls (training penalty only on last frame)	69.15 (1, only last frame)
Bouncing Balls (training penalty only on last frame)	54.65 (32)
Bouncing Balls (training penalty on all frames)	66.86 (32)
MNIST (single-digit)	63.53 (30)
MNIST (multi-digit)	51.62 (30)
KTH (average leave-one-subject-out)	55.03 (full length of test sequences)

various types of transformations. The application of Spatial Transformers in RNNs for digit recognition has been explored in [37].

## ACKNOWLEDGMENTS

The authors would like to thank the developers of Theano [22], [23]. We thank Kishore Konda, Jörg Bornschein and Pierre-Luc St-Charles for helpful discussions. This work was supported by an NSERC Discovery Award, CIFAR, FQRNT and the German BMBF, project 01GQ0841.

## REFERENCES

- [1] K. Gregor, I. Danihelka, A. Graves, and D. Wierstra, "DRAW: A recurrent neural network for image generation," *CoRR*, vol. abs/1502.04623, 2015. [Online]. Available: <http://arxiv.org/abs/1502.04623>
- [2] K. Xu, J. Ba, R. Kiros, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," *arXiv preprint arXiv:1502.03044*, 2015.
- [3] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [4] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," *arXiv preprint arXiv:1509.00685*, 2015.
- [5] S. K. Sønderby, C. K. Sønderby, H. Nielsen, and O. Winther, "Convolutional lstm networks for subcellular localization of proteins," in *Algorithms for computational biology*. Springer, 2015, pp. 68–80.
- [6] V. Mnih, N. Heess, A. Graves *et al.*, "Recurrent models of visual attention," in *Advances in Neural Information Processing Systems*, 2014, pp. 2204–2212.
- [7] H. Larochelle and G. E. Hinton, "Learning to combine foveal glimpses with a third-order boltzmann machine," in *Advances in neural information processing systems*, 2010, pp. 1243–1251.
- [8] M. Denil, L. Bazzani, H. Larochelle, and N. de Freitas, "Learning where to attend with deep architectures for image tracking," *CoRR*, vol. abs/1109.3737, 2011. [Online]. Available: <http://arxiv.org/abs/1109.3737>
- [9] M. Ranzato, "On learning where to look," *arXiv preprint arXiv:1405.5488*, 2014.
- [10] P. Sermanet, A. Frome, and E. Real, "Attention for fine-grained categorization," *arXiv preprint arXiv:1412.7054*, 2014.
- [11] A. W. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah, "Visual tracking: An experimental survey," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 36, no. 7, pp. 1442–1468, 2014.
- [12] H. Nam and B. Han, "Learning multi-domain convolutional neural networks for visual tracking," *CoRR*, vol. abs/1510.07945, 2015.
- [13] N. Wang and D.-Y. Yeung, "Learning a deep compact image representation for visual tracking," in *Advances in neural information processing systems*, 2013, pp. 809–817.
- [14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [15] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," *arXiv preprint arXiv:1211.5063*, 2012.
- [16] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 6645–6649.
- [17] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in neural information processing systems*, 2014, pp. 3104–3112.
- [18] 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.
- [19] N. Srivastava, E. Mansimov, and R. Salakhutdinov, "Unsupervised learning of video representations using lstms," *arXiv preprint arXiv:1502.04681*, 2015.
- [20] Q. V. Le, N. Jaitly, and G. E. Hinton, "A simple way to initialize recurrent neural networks of rectified linear units," *arXiv preprint arXiv:1504.00941*, 2015.
- [21] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 2010, pp. 807–814.
- [22] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. Goodfellow, A. Bergeron, N. Bouchard, D. Warde-Farley, and Y. Bengio, "Theano: new features and speed improvements," *arXiv preprint arXiv:1211.5590*, 2012.
- [23] J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley, and Y. Bengio, "Theano: a cpu and gpu math expression compiler," in *Proceedings of the Python for scientific computing conference (SciPy)*, vol. 4. Austin, TX, 2010, p. 3.
- [24] S. Ebrahimi Kahou, V. Michalski, K. Konda, R. Memisevic, and C. Pal, "Recurrent neural networks for emotion recognition in video," in *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, ser. ICMI '15. New York, NY, USA: ACM, 2015, pp. 467–474. [Online]. Available: <http://doi.acm.org/10.1145/2818346.2830596>
- [25] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [26] I. Sutskever, G. E. Hinton, and G. W. Taylor, "The recurrent temporal restricted boltzmann machine," in *Advances in Neural Information Processing Systems*, 2009, pp. 1601–1608.
- [27] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [28] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.
- [29] C. Schüldt, I. Laptev, and B. Caputo, "Recognizing human actions: a local svm approach," in *Pattern Recognition, 2004. ICPR 2004. Proceedings of the 17th International Conference on*, vol. 3. IEEE, 2004, pp. 32–36.
- [30] Z. Jiang, Z. Lin, and L. S. Davis, "Recognizing human actions by learning and matching shape-motion prototype trees," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 34, no. 3, pp. 533–547, 2012.
- [31] A. Ess, B. Leibe, K. Schindler, and L. van Gool, "A mobile vision system for robust multi-person tracking," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR'08)*. IEEE Press, June 2008.
- [32] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *International Conference on Computer Vision & Pattern Recognition*, C. Schmid, S. Soatto, and C. Tomasi, Eds., vol. 2, INRIA Rhône-Alpes, ZIRST-655, av. de l'Europe, Montbonnot-38334, June 2005, pp. 886–893. [Online]. Available: <http://lear.inrialpes.fr/pubs/2005/DT05>
- [33] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.
- [34] Y. Wu, J. Lim, and M.-H. Yang, "Object tracking benchmark," *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 37, no. 9, pp. 1834–1848, 2015.
- [35] R. Caruana, "Multitask learning," *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [36] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 2008–2016.
- [37] S. K. Sønderby, C. K. Sønderby, L. Maaløe, and O. Winther, "Recurrent spatial transformer networks," *arXiv preprint arXiv:1509.05329*, 2015.