

End-to-end representation learning for Correlation Filter based tracking

Jack Valmadre Luca Bertinetto João Henriques Andrea Vedaldi Philip H. S. Torr
University of Oxford

Fname.surnameG@eng.ox.ac.uk

Figure 1: Overview of the proposed network architecture, **CFNet**. It is an asymmetric Siamese network: after applying the same convolutional feature transform to both input images, the “training image” is used to learn a linear template, which is then applied to search the “test image” by cross-correlation.

Abstract

The Correlation Filter is an algorithm that trains a linear template to discriminate between images and their translations. It is well suited to object tracking because its formulation in the Fourier domain provides a fast solution, enabling the detector to be re-trained once per frame. Previous works that use the Correlation Filter, however, have adopted features that were either manually designed or trained for a different task. This work is the first to overcome this limitation by interpreting the Correlation Filter learner, which has a closed-form solution, as a differentiable layer in a deep neural network. This enables learning deep features that are tightly coupled to the Correlation Filter. Experiments illustrate that our method has the important practical benefit of allowing lightweight architectures to achieve state-of-the-art performance at high framerates.

1. Introduction

Deep neural networks are a powerful tool for learning image representations in computer vision applications. However, training deep networks online, in order to capture previously unseen object classes from one or few examples,

is challenging. This problem emerges naturally in applications such as visual object tracking, where the goal is to re-detect an object over a video with the sole supervision of a bounding box at the beginning of the sequence. The main challenge is the lack of a-priori knowledge of the target object, which can be of any class.

The simplest approach is to disregard the lack of a-priori knowledge and adapt a pre-trained deep convolutional neural network (CNN) to the target, for example by using stochastic gradient descent (SGD), the workhorse of deep network optimization [31, 25, 35]. The extremely limited training data and large number of parameters make this a difficult learning problem. Furthermore, SGD is quite expensive for online adaptation [31, 25].

A possible answer to these shortcomings is to have no online adaptation of the network. Recent works have focused on learning deep embeddings that can be used as universal object descriptors [3, 12, 28, 17, 5]. These methods use a Siamese CNN, trained offline to discriminate whether two image patches contain the same object or not. The idea is that a powerful embedding will allow the detection (and thus tracking) of objects via similarity, bypassing the online learning problem. However, using a fixed metric to compare appearance prevents the learning algorithm from exploiting any video-specific cues that could be helpful for discrimination.

An alternative strategy is to use instead an online learn-

Equal first authorship.

ing method such as the *Correlation Filter* (CF). The CF is an efficient algorithm that learns to discriminate an image patch from the surrounding patches by solving a large ridge regression problem extremely efficiently [4, 13]. It has proved to be highly successful in object tracking (e.g. [6, 18, 22, 2]), where its efficiency enables a tracker to adapt its internal model of the object on the fly at every frame. It owes its speed to a Fourier domain formulation, which allows the ridge regression problem to be solved with only a few applications of the Fast Fourier Transform (FFT) and cheap element-wise operations. Such a solution is, by design, much more efficient than an iterative solver like SGD, and still allows the discriminator to be tailored to a specific video, contrary to the embedding methods.

The challenge, then, is to combine the online learning efficiency of the CF with the discriminative power of CNN features trained offline. This has been done in several works (e.g. [21, 7, 9, 31]), which have shown that CNNs and CFs are complementary and their combination results in improved performance.

However, in the aforementioned works, the CF is simply applied on top of pre-trained CNN features, without any deep integration of the two methods. End-to-end training of deep architectures is generally preferable to training individual components separately. The reason is that in this manner the free parameters in all components can co-adapt and cooperate to achieve a single objective. Thus it is natural to ask whether a CNN-CF combination can also be trained end-to-end with similar benefits.

The key step in achieving such integration is to interpret the CF as a differentiable CNN layer, so that errors can be propagated through the CF back to the CNN features. This is challenging, because the CF itself is the solution of a learning problem. Hence, this requires to differentiate the solution of a large linear system of equations. This paper provides a closed-form expression for the derivative of the Correlation Filter. Moreover, we demonstrate the practical utility of our approach in training CNN architectures end-to-end.

We present an extensive investigation into the effect of incorporating the CF into the fully-convolutional Siamese framework of Bertinetto *et al.* [3]. We find that the CF does not improve results for networks that are sufficiently deep. However, our method enables ultra-lightweight networks of a few thousand parameters to achieve state-of-the-art performance on multiple benchmarks while running at high framerates.

2. Related work

Since the seminal work of Bolme *et al.* [4], the Correlation Filter has enjoyed great popularity within the tracking community. Notable efforts have been devoted to its improvement, for example by mitigating the effect of periodic

boundaries [10, 15, 8], incorporating multi-resolution feature maps [21, 9] and augmenting the objective with a more robust loss [26]. For the sake of simplicity, in this work we adopt the basic formulation of the Correlation Filter.

Recently, several methods based on Siamese networks have been introduced [28, 12, 3], raising interest in the tracking community for their simplicity and competitive performance. For our method, we prefer to build upon the fully-convolutional Siamese architecture [3], as it enforces the prior that the appearance similarity function should commute with translation.

At its core, the Correlation Filter layer that we introduce amounts to computing the solution to a regularized deconvolution problem, not to be confused with upsampling convolution layers that are sometimes referred to as “deconvolution layers” [20]. Before it became apparent that algorithms such as SGD are sufficient for training deep networks, Zeiler *et al.* [34] introduced a deep architecture in which each layer solves a convolutional sparse coding problem. In contrast, our problem has a closed-form solution since the Correlation Filter employs quadratic regularization rather than 1-norm regularization.

The idea of back-propagating gradients through the solution to an optimization problem during training has been previously investigated. Ionescu *et al.* [14] and Murray [24] have presented back-propagation forms for the SVD and Cholesky decomposition respectively, enabling gradient descent to be applied to a network that computes the solution to either a system of linear equations or an eigenvalue problem. Our work can be understood as an efficient back-propagation procedure through the solution to a system of linear equations, where the matrix has circulant structure.

When the solution to the optimization problem is obtained iteratively, an alternative is to treat the iterations as a Recurrent Neural Network, and to explicitly unroll a fixed number of iterations [36]. Maclaurin *et al.* [23] go further and back-propagate gradients through an entire SGD learning procedure, although this is computationally demanding and requires judicious bookkeeping. Gould *et al.* [11] have recently considered differentiating the solution to general arg min problems without restricting themselves to iterative procedures. However, these methods are unnecessary in the case of the Correlation Filter, as it has a closed-form solution.

Back-propagating through a learning algorithm invites a comparison to meta-learning. Recent works [30, 1] have proposed feed-forward architectures that can be interpreted as learning algorithms, enabling optimization by gradient descent. Rather than adopt an abstract definition of learning, this paper propagates gradients through a conventional learning problem that is already widely used.

3. Method

We briefly introduce a framework for learning embeddings with Siamese networks (Section 3.1) and the use of such an embedding for object tracking (Section 3.2) before presenting the CFNet architecture (Section 3.3). We subsequently derive the expressions for evaluation and back-propagation of the main new ingredient in our networks, the Correlation Filter layer, which performs online learning in the forward pass (Section 3.4).

3.1. Fully-convolutional Siamese networks

Our starting point is a network similar to that of [3], which we later modify in order to allow the model to be interpreted as a Correlation Filter tracker. The fully-convolutional Siamese framework considers pairs (x, z) comprising a training image x and a test image z ¹. The image x represents the object of interest (*e.g.* an image patch centered on the target object in the first video frame), while z is typically larger and represents the search area (*e.g.* the next video frame).

Both inputs are processed by a CNN f with learnable parameters θ . This yields two feature maps, which are then cross-correlated:

$$g(x, z) = f(x) \star f(z). \quad (1)$$

Eq. 1 amounts to performing an exhaustive search of the pattern x over the test image z . The goal is for the maximum value of the response map (left-hand side of eq. 1) to correspond to the target location.

To achieve this goal, the network is trained offline with millions of random pairs (x_i, z_i) taken from a collection of videos. Each example has a spatial map of labels c_i with values in $\{-1, 1\}$, with the true object location belonging to the positive class and all others to the negative class. Training proceeds by minimizing an element-wise logistic loss over the training set:

$$\arg \min_{\theta} \sum_i (g(x_i, z_i) - c_i). \quad (2)$$

3.2. Tracking algorithm

The network itself only provides a function to measure the similarity of two image patches. To apply this network to object tracking, it is necessary to combine this with a procedure that describes the logic of the tracker. Similar to [3], we employ a simplistic tracking algorithm to assess the utility of the similarity function.

Online tracking is performed by simply evaluating the network in forward-mode. The feature representation of the target object is compared to that of the search region, which

¹Note that this differs from [3], in which the target object and search area were instead denoted z and x respectively.

is obtained in each new frame by extracting a window centred at the previously estimated position, with an area that is four times the size of the object. The new position of the object is taken to be the location with the highest score.

The original fully-convolutional Siamese network simply compared every frame to the initial appearance of the object. In contrast, we compute a new template in each frame and then combine this with the previous template in a moving average.

3.3. Correlation Filter networks

We propose to modify the baseline Siamese network of eq. 1 with a Correlation Filter block between x and the cross-correlation operator. The resulting architecture is illustrated in Figure 1. This change can be formalized as:

$$h_{s,b}(x, z) = s \cdot (f(x) \star f(z)) + b \quad (3)$$

The CF block $w = f(x)$ computes a standard CF template w from the training feature map $x = f(x)$ by solving a ridge regression problem in the Fourier domain [13]. Its effect can be understood as crafting a discriminative template that is robust against translations. It is necessary to introduce scalar parameters s and b (scale and bias) to make the score range suitable for logistic regression. Offline training is then performed in the same way as for a Siamese network (Section 3.1), replacing g with h in eq. 2.

We found that it was important to provide the Correlation Filter with a large region of context in the training image, which is consistent with the findings of Danelljan et al. [8] and Kiani et al. [15]. To reduce the effect of circular boundaries, the feature map x is pre-multiplied by a cosine window [4] and the final template is cropped [29].

Notice that the forward pass of the architecture in Figure 1 corresponds exactly to the operation of a standard CF tracker [13, 6, 22, 3] with CNN features, as proposed in previous work [21, 7]. However, these earlier networks were not trained end-to-end. The novelty is to compute the derivative of the CF template with respect to its input so that a network incorporating a CF can be trained end-to-end.

3.4. Correlation Filter

We now show how to back-propagate gradients through the Correlation Filter solution efficiently and in closed form via the Fourier domain.

Formulation. Given a scalar-valued image $x \in \mathbb{R}^{m \times m}$, the Correlation Filter is the template $w \in \mathbb{R}^{m \times m}$ whose inner product with each circular shift of the image x_{-u} is as close as possible to a desired response $y[u]$ [13], minimizing

$$\sum_u (x_{-u} \star w - y[u])^2 = \|w \star x - y\|^2. \quad (4)$$

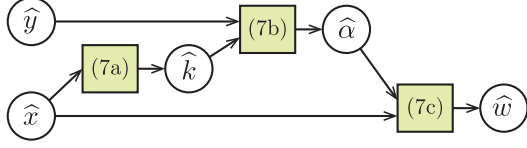


Figure 2: Internal computational graph for the Correlation Filter. The boxes denote functions, which are defined in eq. 7, and the circles denote variables.

Here, $U = \{0, \dots, m-1\}^2$ is the domain of the image, $y \in \mathbb{R}^{m \times m}$ is a signal whose u -th element is $y[u]$, and $\delta[t] = \delta[t - \cdot]$ is the translated Dirac delta function. In this section, we use \circ to denote circular convolution and \star to denote circular cross-correlation. Recall that convolution with the translated function is equivalent to translation $(x \circ \delta)[t] = x[t - \cdot \bmod m]$. Incorporating quadratic regularization to prevent overfitting, the problem is to find

$$\arg \min_w \frac{1}{2n} \|w \circ x - y\|^2 + \frac{\lambda}{2} \|w\|^2 \quad (5)$$

where $n = |U|$ is the effective number of examples.

The optimal template w must satisfy the system of equations (obtained via the Lagrangian dual, see Appendix C, supplementary material)

$$\begin{aligned} k &= \frac{1}{n} (x \circ x) + \\ k &= \frac{1}{n} y \\ w &= x \end{aligned} \quad (6)$$

where k can be interpreted as the signal that defines a circulant linear kernel matrix, and $\mathbf{1}$ is a signal comprised of the Lagrange multipliers of a constrained optimization problem that is equivalent to eq. 5. The solution to eq. 6 can be computed efficiently in the Fourier domain [13],

$$k = \frac{1}{n} (x \circ x) + \quad (7a)$$

$$= \frac{1}{n} k^{-1} y \quad (7b)$$

$$w = x \quad (7c)$$

where we use $\tilde{x} = F x$ to denote the Discrete Fourier Transform of a variable, x^* to denote the complex conjugate, \odot to denote element-wise multiplication and \oslash to denote a signal of ones. The inverse of element-wise multiplication is element-wise scalar inversion. Notice that the operations in eq. 7 are more efficiently computed in the Fourier domain, since they involve element-wise operations instead of more expensive convolutions or matrix operators (eq. 6). Moreover, the inverse convolution problem (to find \tilde{k} such that $k \oslash \tilde{k} = \frac{1}{n} y$) is the solution to a diagonal system of equations in the Fourier domain (eq. 7b).

Back-propagation. We adopt the notation that if $x \in \mathbb{R}^n$ is a variable in a computational graph that computes a final scalar loss R , then $\frac{\partial R}{\partial x}$ denotes the vector of partial derivatives $(\frac{\partial R}{\partial x})_i = \frac{\partial R}{\partial x_i}$. If $y \in \mathbb{R}^m$ is another variable in the graph, which is computed directly from x according to $y = f(x)$, then the so-called *back-propagation map* for the function f is a linear map from $\frac{\partial R}{\partial y}$ to $\frac{\partial R}{\partial x}$.

Appendix D gives a tutorial review of the mathematical background. In short, the back-propagation map is the linear map which is the adjoint of the differential. This property was used by Ionescu *et al.* [14] to compute back-propagation maps using matrix differential calculus. While they used the matrix inner product $\langle X, Y \rangle = \text{tr}(X^T Y)$ to find the adjoint, we use Parseval's theorem, which states that the Fourier transform is unitary (except for a scale factor) and therefore preserves inner products $\langle x, y \rangle = \langle \tilde{x}, \tilde{y} \rangle$.

To find the linear map for back-propagation through the Correlation Filter, we first take the differentials of the system of equations in eq. 6 that defines the template w

$$\begin{aligned} dk &= \frac{1}{n} (dx \circ x + x \circ dx) \\ dk &= k \odot d = \frac{1}{n} dy \\ dw &= d \circ x + x \circ dx \end{aligned} \quad (8)$$

and then take the Fourier transform of each equation and rearrange to give the differential of each dependent variable in Figure 2 as a linear function (in the Fourier domain) of the differentials of its input variables

$$dk = \frac{1}{n} (dx \circ x + x \circ dx) \quad (9a)$$

$$d = k^{-1} \odot \frac{1}{n} dy - dk \quad (9b)$$

$$dw = d \circ x + x \circ dx \quad (9c)$$

Note that while these are complex equations, that is simply because they are the Fourier transforms of real equations. The derivatives themselves are all computed with respect to real variables.

The adjoints of these linear maps define the overall back-propagation map from $\frac{\partial R}{\partial w}$ to $\frac{\partial R}{\partial x}$ and $\frac{\partial R}{\partial y}$. We defer the derivation to Appendix B and present here the final result,

$$\begin{aligned} \frac{\partial R}{\partial x} &= \tilde{x}^* \odot \left(\frac{\partial R}{\partial w} \right) \\ \frac{\partial R}{\partial y} &= \frac{1}{n} k^{-1} \odot \frac{\partial R}{\partial w} \\ \frac{\partial R}{\partial x} &= \tilde{x}^* \odot \left(\frac{\partial R}{\partial w} + \frac{2}{n} x \odot \text{Re}\{ \tilde{k} \} \right) \end{aligned} \quad (10)$$

It is necessary to compute forward Fourier transforms at the start and inverse transforms at the end. The extension to multi-channel images is trivial and given in Appendix E (supplementary material).

As an interesting aside, we remark that, since we have the gradient of the loss with respect to the “desired” response y , it is actually possible to optimize for this parameter rather than specify it manually. However, in practice we did not find learning this parameter to improve the tracking accuracy compared to the conventional choice of a fixed Gaussian response [4, 13].

4. Experiments

The principal aim of our experiments is to investigate the effect of incorporating the Correlation Filter during training. We first compare against the symmetric Siamese architecture of Bertinetto *et al.* [3]. We then compare the end-to-end trained CFNet to a variant where the features are replaced with features that were trained for a different task. Finally, we demonstrate that our method achieves state-of-the-art results.

4.1. Evaluation criteria

Popular tracking benchmarks like VOT [16] and OTB [32, 33] have made all ground truth annotations available and do not enforce a validation/test split. However, in order to avoid overfitting to the test set in design choices and hyperparameter selection, we consider OTB-2013, OTB-50 and OTB-100 as our *test set* and 129 videos from VOT-2014, VOT-2016 and Temple-Color [19] as our *validation set*, excluding any videos which were already assigned to the test set. We perform all of our tracking experiments in Sections 4.2, 4.3 and 4.4 on the validation set with the same set of “natural” hyperparameters, which are reasonable for all methods and not tuned for any particular method.

As in the OTB benchmark [32, 33], we quantify the performance of the tracker on a sequence in terms of the average overlap (intersection over union) of the predicted and ground truth rectangles in all frames. The success rate of a tracker at a given threshold corresponds to the fraction of frames in which the overlap with the ground truth is at least . This is computed for a uniform range of 100 thresholds between 0 and 1, effectively constructing the cumulative distribution function. Trackers are compared using the area under this curve.

Mimicking the TRE (Temporal Robustness Evaluation) mode of OTB, we choose three equispaced points per sequence and run the tracker from each until the end. Differently from the OTB evaluation, when the target is *lost* (*i.e.* the overlap with the ground truth becomes zero) the tracker is terminated and an overlap of zero is reported for all remaining frames.

Despite the large number of videos, we still find that the performance of similarity networks varies considerably as training progresses. To mitigate this effect, we average the final tracking results that are obtained using the parameters of the network at epochs 55, 60, . . . , 95, 100 (the final

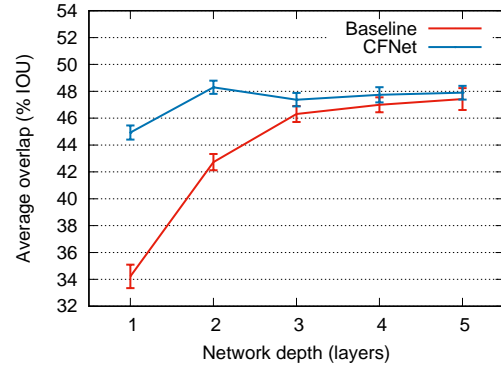


Figure 3: Tracker accuracy for different network depths, on the 129 videos of the validation set. Error bars indicate two standard deviations. Refer to section 4.2 for more details. All figures best viewed in colour.

epoch) to reduce the variance. These ten results are used to estimate the standard deviation of the distribution of results, providing error bars for most figures in this section. While it would be preferable to train all networks to convergence multiple times with different random seeds, this would require significantly more resources.

4.2. Comparison to Siamese baseline

Figures 3 and 4 compare the accuracy of both methods on the validation set for networks of varying depth. The feature extraction network of depth n is terminated after the n -th linear layer, including the following ReLU but not the following pooling layer (if any).

Our baseline diverges slightly from [3] in two ways. Firstly, we reduce the total stride of the network from 8 to 4 (2 at conv1, 2 at pool1) to avoid training Correlation Filters with small feature maps. Secondly, we always restrict the final layer to 32 output channels in order to preserve the high speed of the method with larger feature maps. These changes did not have a negative effect on the tracking performance of SiamFC.

The results show that CFNet is significantly better than the baseline when shallow networks are used to compute features. Specifically, it brings a relative improvement of 31% and 13% for networks of depth one and two respectively. At depths three, four and five, the difference is much less meaningful. CFNet is relatively unaffected by the depth of the network, whereas the performance of the baseline increases steadily and significantly with depth. It seems that the ability of the Correlation Filter to adapt the distance metric to the content of the training image is less important given a sufficiently expressive embedding function.

The CF layer can be understood to encode prior knowledge of the test-time procedure. This prior may become redundant or even overly restrictive when enough model ca-

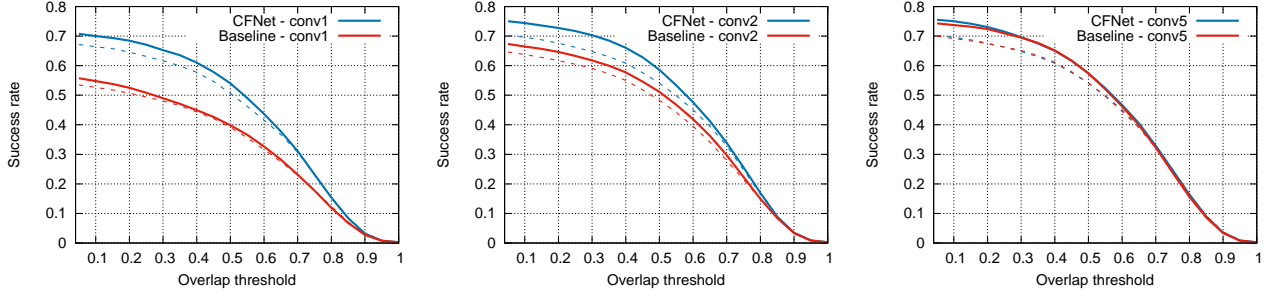


Figure 4: Success rates of rectangle overlap for individual trackers on the validation set. Solid and dotted lines represent methods that update the template with a running average learning rate of 0.01 and 0, respectively.

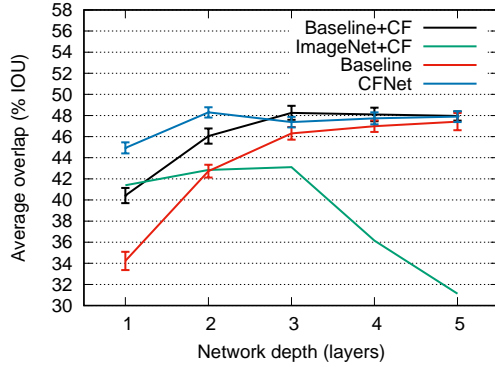


Figure 5: Accuracy of a Correlation Filter tracker when using features obtained via different methods. Error bars indicate two standard deviations. Refer to Section 4.3 for details.

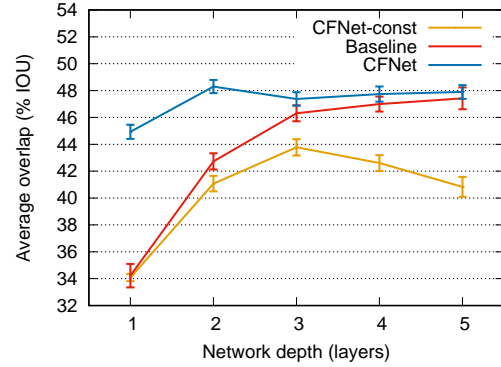


Figure 6: Comparison of CFNet to a “constant” variant of the architecture, in which the Lagrange multipliers do not depend on the image (section 4.4). Error bars indicate two standard deviations.

capacity and data are available. We believe this explains the saturation of CFNet performance when more than two convolutional layers are used.

Figure 4 additionally shows that updating the template is always helpful, for both Baseline and CFNet architectures, at any depth.

4.3. Feature transfer experiment

The motivation for this work was the hypothesis that incorporating the CF during training will result in features that are better suited to tracking with a CF. We now compare our end-to-end trained CFNet to variants that use features from alternative sources: *Baseline+CF* and *ImageNet+CF*. The results are presented in Figure 5.

To obtain the curve *Baseline+CF* we trained a baseline Siamese network of the desired depth and then combined those features with a CF during tracking. Results show that taking the CF into account during offline training is critical at depth one and two. However, it seems redundant when more convolutional layers are added, since using features from the *Baseline* in conjunction with the CF achieves sim-

ilar performance.

The *ImageNet+CF* variant employs features taken from a network trained to solve the ImageNet classification challenge [27]. The results show that these features, which are often the first choice for combining CFs with CNNs [7, 9, 21, 25, 31, 35], are significantly worse than those learned by *CFNet* and the *Baseline* experiment. The particularly poor performance of these features at deeper layers is somewhat unsurprising, since these layers are expected to have greater invariance to position when trained for classification.

4.4. Importance of adaptation

For a multi-channel CF, each channel p of the template w can be obtained as $w_p = x_p$, where x is itself a function of the exemplar x (Appendix C, supplementary material). To verify the importance of the online adaptation that solving a ridge regression problem at test time should provide, we propose a “constant” version of the Correlation Filter (*CFNet-const*) where the vector of Lagrange multipliers is instead a parameter of the network that is learned offline and remains fixed at test time.

Method	speed (fps.)	OTB-2013				OTB-50				OTB-100			
		OPE		TRE		OPE		TRE		OPE		TRE	
		IoU	prec.	IoU	prec.	IoU	prec.	IoU	prec.	IoU	prec.	IoU	prec.
CFNet-conv1	83	57.8	71.4	58.6	71.7	48.8	61.3	51.0	63.6	53.6	65.8	55.9	67.6
CFNet-conv2	75	61.1	74.6	64.0	77.9	53.0	66.0	56.5	70.2	56.8	69.3	60.6	73.2
Baseline+CF-conv3	67	61.0	74.8	<u>63.1</u>	76.8	<u>53.8</u>	<u>66.5</u>	57.4	70.8	58.9	71.1	<u>61.1</u>	<u>73.4</u>
CFNet-conv5	43	61.1	73.6	62.6	75.7	53.9	67.0	56.6	70.1	58.6	71.1	60.8	72.7
Baseline-conv5	52	61.8	<u>75.3</u>	64.0	<u>77.3</u>	51.7	64.1	56.1	69.1	<u>58.8</u>	<u>71.4</u>	61.6	73.7
SiamFC-3s [3]		60.7	73.5	61.8	75.0	51.6	63.9	55.5	69.2	58.2	70.2	60.5	72.8
Staple [2]		60.0	72.5	61.7	74.2	50.9	63.4	54.1	67.5	58.1	71.6	60.4	72.8
LCT [22]		<u>61.2</u>	78.0	59.4	74.2	49.2	62.5	49.5	61.7	56.2	69.2	56.9	68.2
SAMF [18]		–	–	–	–	46.2	60.7	51.4	65.6	53.9	69.0	57.7	71.4
DSST [6]		55.4	67.5	56.6	68.4	45.2	56.6	48.4	60.1	51.3	63.1	–	–

Table 1: Performance as overlap (IoU) and precision produced by the OTB toolkit for the OTB-2013, OTB-50 and OTB-100 datasets. The **first** and second best results are highlighted in each column. For details refer to Section 4.5.

Figure 6 compares CFNet to its constant variant. CFNet is consistently better, demonstrating that in order to improve over the baseline Siamese network it is paramount to back-propagate through the solution to the inverse convolution problem that defines the Lagrange multipliers.

4.5. Comparison with the state-of-the-art

We use the OTB-2013/50/100 benchmarks to confirm that our results are on par with the state-of-the-art. All numbers in this section are obtained using the OTB toolkit [32]. We report the results for the three best instantiations of CFNet from Figure 5 (*CFNet-conv2*, *CFNet-conv5*, *Baseline+CF-conv3*), the best variant of the baseline (*Baseline-conv5*) and the most promising single-layer network (*CFNet-conv1*). We compare our methods against state-of-the-art trackers that can operate in real-time: SiamFC-3s [3], Staple [2] and LCT [22]. We also include the recent SAMF [18] and DSST [6] for reference.

For the evaluation of this section, we use a different set of tracking hyperparameters per architecture, chosen to maximize the performance on the validation set after a random search of 300 iterations. More details are provided in the supplementary material. For the few greyscale sequences present in OTB, we re-train each architecture using exclusively greyscale images.

Both overlap (IoU) and precision scores [33] are reported for OPE (one pass) and TRE (temporal robustness) evaluations. For OPE, the tracker is simply run once on each sequence, from the start to the end. For TRE, the tracker is instead started from twenty different starting points, and run until the end from each. We observed that this ensures more robust and reliable results compared to OPE.

Similarly to the analysis on the validation set, *CFNet-conv2* is among the top performers and its accuracy rivals that of *Baseline-conv5*, which possesses approximately 30× as many parameters. In general, our best proposed CFNet variants are superior (albeit modestly) to the state-of-the-

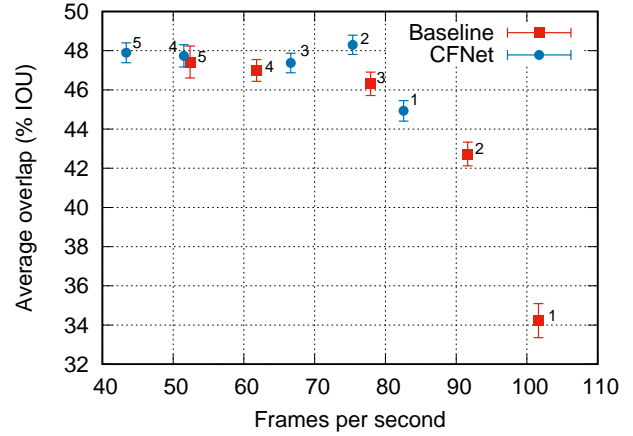


Figure 7: Tracker accuracy versus speed for CFNet and Siamese baseline. Labels indicate network depth. CFNet enables better accuracy to be obtained at higher speeds using shallower networks. Error bars indicate two standard deviations. Refer to section 4.6 for details.

art. In order to focus on the impact of our contribution, we decided to avoid including orthogonal improvements which can often be found in the tracking literature (e.g. bounding box regression [25], ensembling of multiple cues [22, 2], optical flow [28]).

4.6. Speed and practical benefits

The previous sections have demonstrated that there is a clear benefit to integrating Correlation Filters into Siamese networks when the feature extraction network is relatively shallow. Shallow networks are practically advantageous in that they require fewer operations and less memory to evaluate and store. To understand the trade-off, Figure 7 reports the speed and accuracy of both CFNet and the baseline for

varying network depth².

This plot suggests that the two-layer CFNet could be the most interesting variant for practitioners requiring an accurate tracking algorithm that operates at high framerates. It runs at 75 frames per second and has less than 4% of the parameters of the five-layer baseline, requiring only 600kB to store. This may be of particular interest for embedded devices with limited memory. In contrast, methods like DeepSRDCF [7] and C-COT [9], which use out-of-the-box deep features for the Correlation Filter, run orders of magnitude slower. Even the one-layer CFNet remains competitive despite having less than 1% of the parameters of the five-layer baseline and requiring under 100kB to store.

5. Conclusion

This work proposes the Correlation Filter network, an asymmetric architecture that back-propagates gradients through an online learning algorithm to optimize the underlying feature representation. This is made feasible by establishing an efficient back-propagation map for the solution to a system of circulant equations.

Our empirical investigation reveals that, for a sufficiently deep Siamese network, adding a Correlation Filter layer does not significantly improve the tracking accuracy. We believe this is testament to the power of deep learning given sufficient training data. However, incorporating the Correlation Filter into a similarity network during training does enable shallow networks to rival their slower, deeper counterparts.

Future research may include extensions to account for adaptation over time, and back-propagating gradients through learning problems for related tasks such as one-shot learning and domain adaptation.

A. Implementation details

We follow the procedure of [3] to minimize the loss (equation 2) through SGD, with the Xavier-improved parameters initialization and using mini-batches of size 8. We use all the 3862 training videos of ImageNet Video [27], containing more than 1 million annotated frames, with multiple objects per frame. Training is conducted for 100 epochs, each sampling approximately 12 pairs (x_i, z_i) from each video, randomly extracted so that they are at most 100 frames apart.

During tracking, a spatial cosine window is multiplied with the score map to penalize large displacements. Tracking in scale space is achieved by evaluating the network at the scale of the previous object and at one adjacent scale on either side, with a geometric step of 1.04. Updating the scale is discouraged by multiplying the responses of the

scaled object by 0.97. To avoid abrupt transitions of object size, scale is updated using a rolling average with learning rate 0.6.

Code and results are available online ³.

B. Back-propagation for the Correlation Filter

As described in Appendix D (supplementary material), the back-propagation map is the adjoint of the linear maps that is the differential. These linear maps for the Correlation Filter are presented in eq. 9. We are free to obtain these adjoint maps in the Fourier domain since Parseval's theorem provides the preservation of inner products. Let J_1 denote the map $dx \rightarrow dk$ in eq. 9a. Hence manipulation of the inner product

$$\begin{aligned} F dk, F J_1(dx) &= dk, \frac{1}{n}(dx \cdot x + x \cdot dx) \\ &= \frac{1}{n} dx, dk \cdot x + dk \cdot x, dx \\ &= dx, \frac{2}{n} \text{Re}\{dk\} \cdot x \end{aligned} \quad (11)$$

gives the back-propagation map

$$x = \frac{2}{n} x \cdot \text{Re}\{k\} \cdot \quad (12)$$

Similarly, for the linear map $dk, dy \rightarrow d$ in eq. 9b,

$$\begin{aligned} F d, F J_2(dk, dy) &= d, k^{-1}[\frac{1}{n} dy - dk \cdot] \\ &= \frac{1}{n} k^{-1} d, dy + -k^{-1} d, dk \end{aligned} \quad (13)$$

the back-propagation maps are

$$y = \frac{1}{n} k^{-1} \quad (14)$$

$$k = -k^{-1} \quad (15)$$

and for the linear map $dx, d \rightarrow dw$ in eq. 9c,

$$\begin{aligned} F dw, F J_3(dx, d) &= dw, d \cdot x + dx \cdot \\ &= d, dw \cdot x + dw \cdot, dx \end{aligned} \quad (16)$$

the back-propagation maps are

$$w = x \cdot (w) \cdot, \quad (17)$$

$$x = w \cdot \quad (18)$$

The two expressions for x above are combined to give the back-propagation map for the entire Correlation Filter block in eq. 10.

Acknowledgements. This research was supported by Apical Ltd., EPSRC grant Seebibyte EP/M013774/1 and ERC grants ERC-2012-AdG 321162-HELIOS, HELIOS-DFR00200, “Integrated and Detailed Image Understanding” (EP/L024683/1) and ERC 677195-IDIU.

²The speed was measured using a 4.0GHz Intel i7 CPU and an NVIDIA Titan X GPU.

³www.robots.ox.ac.uk/~luca/cfnet.html

References

- [1] L. Bertinetto, J. F. Henriques, J. Valmadre, P. H. S. Torr, and A. Vedaldi. Learning feed-forward one-shot learners. In *NIPS 2016*, pages 523–531, 2016. [2](#)
- [2] L. Bertinetto, J. Valmadre, S. Golodetz, O. Miksik, and P. H. S. Torr. Staple: Complementary learners for real-time tracking. In *CVPR 2016*, pages 1401–1409, 2016. [2](#), [7](#)
- [3] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. S. Torr. Fully-convolutional Siamese networks for object tracking. In *ECCV 2016 Workshops*, pages 850–865, 2016. [1](#), [2](#), [3](#), [5](#), [7](#), [8](#)
- [4] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui. Visual object tracking using adaptive correlation filters. In *CVPR 2010*, 2010. [2](#), [3](#), [5](#)
- [5] K. Chen and W. Tao. Once for all: A two-flow convolutional neural network for visual tracking. *arXiv preprint arXiv:1604.07507*, 2016. [1](#)
- [6] M. Danelljan, G. Häger, F. Khan, and M. Felsberg. Accurate scale estimation for robust visual tracking. In *BMVC 2014*, 2014. [2](#), [3](#), [7](#)
- [7] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg. Convolutional features for correlation filter based visual tracking. In *ICCV 2015 Workshops*, pages 58–66, 2015. [2](#), [3](#), [6](#), [8](#)
- [8] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg. Learning spatially regularized correlation filters for visual tracking. In *ICCV 2015*, pages 4310–4318, 2015. [2](#), [3](#)
- [9] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *ECCV 2016*, pages 472–488, 2016. [2](#), [6](#), [8](#)
- [10] J. A. Fernandez and B. Vijayakumar. Zero-aliasing correlation filters. In *International Symposium on Image and Signal Processing and Analysis 2013*, pages 101–106, 2013. [2](#)
- [11] S. Gould, B. Fernando, A. Cherian, P. Anderson, R. S. Cruz, and E. Guo. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*, 2016. [2](#)
- [12] D. Held, S. Thrun, and S. Savarese. Learning to track at 100 fps with deep regression networks. In *ECCV 2016*, pages 749–765. Springer, 2016. [1](#), [2](#)
- [13] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. High-speed tracking with kernelized correlation filters. *IEEE TPAMI*, 37(3):583–596, 2015. [2](#), [3](#), [4](#), [5](#)
- [14] C. Ionescu, O. Vantzou, and C. Sminchisescu. Matrix back-propagation for deep networks with structured layers. In *ICCV 2015*, pages 2965–2973, 2015. [2](#), [4](#)
- [15] H. Kiani Galoogahi, T. Sim, and S. Lucey. Correlation filters with limited boundaries. In *CVPR 2015*, pages 4630–4638, 2015. [2](#), [3](#)
- [16] M. Kristan, A. Leonardis, J. Matas, M. Felsberg, R. Pflugfelder, L. Čehovin, T. Vojir, G. Häger, A. Lukežič, G. Fernández, et al. The Visual Object Tracking VOT2016 challenge results. 2016. [5](#)
- [17] L. Leal-Taixé, C. Canton-Ferrer, and K. Schindler. Learning by tracking: Siamese CNN for robust target association. In *CVPR 2016 Workshops*, pages 33–40, 2016. [1](#)
- [18] Y. Li and J. Zhu. A scale adaptive kernel correlation filter tracker with feature integration. In *ECCV 2014*, pages 254–265, 2014. [2](#), [7](#)
- [19] P. Liang, E. Blasch, and H. Ling. Encoding color information for visual tracking: Algorithms and benchmark. *IEEE Transactions on Image Processing*, 24(12):5630–5644, 2015. [5](#)
- [20] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR 2015*, pages 3431–3440, 2015. [2](#)
- [21] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang. Hierarchical convolutional features for visual tracking. In *ICCV 2015*, pages 3074–3082, 2015. [2](#), [3](#), [6](#)
- [22] C. Ma, X. Yang, C. Zhang, and M.-H. Yang. Long-term correlation tracking. In *CVPR 2015*, pages 5388–5396, 2015. [2](#), [3](#), [7](#)
- [23] D. Maclaurin, D. Duvenaud, and R. P. Adams. Gradient-based hyperparameter optimization through reversible learning. In *ICML 2015*, 2015. [2](#)
- [24] I. Murray. Differentiation of the Cholesky decomposition. *arXiv preprint arXiv:1602.07527*, 2016. [2](#)
- [25] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking. In *CVPR 2016*, pages 4293–4302, 2016. [1](#), [6](#), [7](#)
- [26] A. Rodriguez, V. N. Boddeti, B. V. K. V. Kumar, and A. Mahalanobis. Maximum margin correlation filter: A new approach for localization and classification. *IEEE Transactions on Image Processing*, 22(2):631–643, 2013. [2](#)
- [27] 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. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [6](#), [8](#)
- [28] R. Tao, E. Gavves, and A. W. M. Smeulders. Siamese instance search for tracking. In *CVPR 2016*, pages 1420–1429, 2016. [1](#), [2](#), [7](#)
- [29] J. Valmadre, S. Sridharan, and S. Lucey. Learning detectors quickly with stationary statistics. In *ACCV 2014*, pages 99–114. Springer, 2014. [3](#)
- [30] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *NIPS 2016*, pages 3630–3638, 2016. [2](#)
- [31] 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](#), [2](#), [6](#)
- [32] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark. In *CVPR 2015*, pages 2411–2418, 2013. [5](#), [7](#)
- [33] Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark. *TPAMI*, 37(9):1834–1848, 2015. [5](#), [7](#)
- [34] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus. Deconvolutional networks. In *CVPR 2010*, pages 2528–2535, 2010. [2](#)
- [35] M. Zhai, M. J. Roshtkhari, and G. Mori. Deep learning of appearance models for online object tracking. *arXiv preprint arXiv:1607.02568*, 2016. [1](#), [6](#)
- [36] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. S. Torr. Conditional random fields as recurrent neural networks. In *ICCV 2015*, pages 1529–1537, 2015. [2](#)