

ECO: Efficient Convolution Operators for Tracking

Martin Danelljan, Goutam Bhat, Fahad Shahbaz Khan, Michael Felsberg
Computer Vision Laboratory, Department of Electrical Engineering, Linköping University, Sweden
{martin.danelljan, goutam.bhat, fahad.khan, michael.felsberg}@liu.se

Abstract

In recent years, Discriminative Correlation Filter (DCF) based methods have significantly advanced the state-of-the-art in tracking. However, in the pursuit of ever increasing tracking performance, their characteristic speed and real-time capability have gradually faded. Further, the increasingly complex models, with massive number of trainable parameters, have introduced the risk of severe over-fitting. In this work, we tackle the key causes behind the problems of computational complexity and over-fitting, with the aim of simultaneously improving both speed and performance.

We revisit the core DCF formulation and introduce: (i) a factorized convolution operator, which drastically reduces the number of parameters in the model; (ii) a compact generative model of the training sample distribution, that significantly reduces memory and time complexity, while providing better diversity of samples; (iii) a conservative model update strategy with improved robustness and reduced complexity. We perform comprehensive experiments on four benchmarks: VOT2016, UAV123, OTB-2015, and Temple-Color. When using expensive deep features, our tracker provides a 20-fold speedup and achieves a 13.0% relative gain in Expected Average Overlap compared to the top ranked method [12] in the VOT2016 challenge. Moreover, our fast variant, using hand-crafted features, operates at 60 Hz on a single CPU, while obtaining 65.0% AUC on OTB-2015.

1. Introduction

Generic visual tracking is one of the fundamental problems in computer vision. It is the task of estimating the trajectory of a target in an image sequence, given only its initial state. Online visual tracking plays a crucial role in numerous real-time vision applications, such as smart surveillance systems, autonomous driving, UAV monitoring, intelligent traffic control, and human-computer-interfaces. Due to the online nature of tracking, an ideal tracker should be accurate and robust under the hard computational constraints of real-time vision systems.

In recent years, Discriminative Correlation Filter (DCF)

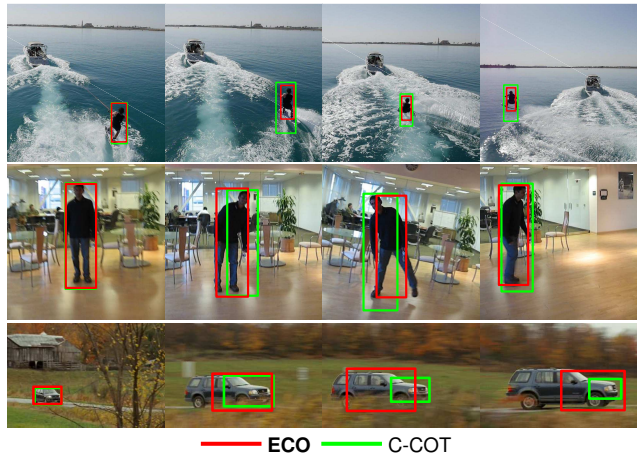


Figure 1. A comparison of our approach ECO with the baseline C-COT [12] on three example sequences. In all three cases, C-COT suffers from severe over-fitting to particular regions of the target. This causes poor target estimation in cases of scale variations (top row), deformations (middle row), and out-of-plane rotations (bottom row). Our ECO tracker successfully tackles the causes of over-fitting, leading to better generalization of the target appearance, while achieving a 20-fold speedup.

based approaches have shown continuous performance improvements in terms of accuracy and robustness on tracking benchmarks [23, 37]. The recent advancement in DCF based tracking performance is driven by the use of multi-dimensional features [13, 15], robust scale estimation [7, 11], non-linear kernels [20], long-term memory components [28], sophisticated learning models [3, 10] and reducing boundary effects [9, 16]. However, these improvements in accuracy come at the price of significant reductions in tracking speed. For instance, the pioneering MOSSE tracker by Bolme *et al.* [4] is about $1000\times$ faster than the recent top-ranked DCF tracker, C-COT [12], in the VOT2016 challenge [23], but obtains only half the accuracy.

As mentioned above, the advancement in DCF tracking performance is predominantly attributed to powerful features and sophisticated learning formulations [8, 12, 27]. This has led to substantially larger models, requiring hundreds of thousands of trainable parameters. On the other hand, such complex and large models have introduced the

risk of severe over-fitting (see figure 1). In this paper, we tackle the issues of over-fitting in recent DCF trackers, while restoring their hallmark real-time capabilities.

1.1. Motivation

We identify three key factors that contribute to *both* increased computational complexity and over-fitting in state-of-the-art DCF trackers.

Model size: The integration of high-dimensional feature maps, such as deep features, has led to a radical increase in the number of appearance model parameters, often beyond the dimensionality of the input image. As an example, C-COT [12] continuously updates about 800,000 parameters during the online learning of the model. Due to the inherent scarcity of training data in tracking, such a high-dimensional parameter space is prone to over-fitting. Further, the high dimensionality causes an increase in the computational complexity, leading to slow tracking speed.

Training set size: State-of-the-art DCF trackers, including C-COT, require a large training sample set to be stored due to their reliance on iterative optimization algorithms. In practice however, the memory size is limited, particularly when using high-dimensional features. A typical strategy for maintaining a feasible memory consumption is to discard the oldest samples. This may however cause over-fitting to recent appearance changes, leading to model drift (see figure 1). Moreover, a large training set increases the computational burden.

Model update: Most DCF-based trackers apply a continuous learning strategy, where the model is updated rigorously in every frame. On the contrary, recent works have shown impressive performance without any model update, using Siamese networks [2]. Motivated by these findings, we argue that the continuous model update in state-of-the-art DCF is excessive and sensitive to sudden changes caused by, *e.g.*, scale variations, deformations, and out-of-plane rotations (see figure 1). This excessive update strategy causes both lower frame-rates and degradation of robustness due to over-fitting to the recent frames.

1.2. Contributions

We propose a novel formulation that addresses the previously listed issues of state-of-the-art DCF trackers. As our first contribution, we introduce a factorized convolution operator that dramatically reduces the number of parameters in the DCF model. Our second contribution is a compact generative model of the training sample space that effectively reduces the number of samples in the learning, while maintaining their diversity. As our final contribution, we introduce an efficient model update strategy, that simultaneously improves tracking speed and robustness.

Comprehensive experiments clearly demonstrate that our approach concurrently improves both tracking performance

and speed, thereby setting a new state-of-the-art on four benchmarks: VOT2016, UAV123, OTB-2015, and Temple-Color. Our approach significantly reduces the number of model parameters by 80%, training samples by 90% and optimization iterations by 80% in the learning, compared to the baseline. On VOT2016, our approach outperforms the top ranked tracker, C-COT [12], in the challenge, while achieving a significantly higher frame-rate. Furthermore, we propose a fast variant of our tracker that maintains competitive performance, with a speed of 60 frames per second (FPS) on a single CPU, thereby being especially suitable for computationally restricted robotics platforms.

2. Baseline Approach: C-COT

In this work, we collectively address the problems of computational complexity and over-fitting in state-of-the-art DCF trackers. We adopt the recently introduced Continuous Convolution Operator Tracker (C-COT) [12] as our baseline. The C-COT obtained the top rank in the recent VOT2016 challenge [23], and has demonstrated outstanding results on other tracking benchmarks [26, 37]. Unlike the standard DCF formulation, Danelljan *et al.* [12] pose the problem of learning the filters in the continuous spatial domain. The generalized formulation in C-COT yields two advantages that are relevant to our work.

The first advantage of C-COT is the natural integration of multi-resolution feature maps, achieved by performing convolutions in the continuous domain. This provides the flexibility of choosing the cell size (*i.e.* resolution) of each visual feature independently, without the need for explicit re-sampling. The second advantage is that the predicted detection scores of the target are directly obtained as a continuous function, enabling accurate sub-grid localization.

Here, we briefly describe the C-COT formulation, adopting the same notation as in [12] for convenience. The C-COT discriminatively learns a convolution filter based on a collection of M training samples $\{x_j\}_1^M \subset \mathcal{X}$. Unlike the standard DCF, each feature layer $x_j^d \in \mathbb{R}^{N_d}$ has an independent resolution N_d .¹ The feature map is transferred to the continuous spatial domain $t \in [0, T)$ by introducing an interpolation model, given by the operator J_d ,

$$J_d\{x^d\}(t) = \sum_{n=0}^{N_d-1} x^d[n]b_d\left(t - \frac{T}{N_d}n\right). \quad (1)$$

Here, b_d is an interpolation kernel with period $T > 0$. The result $J_d\{x^d\}$ is thus an interpolated feature layer, viewed as a continuous T -periodic function. We use $J\{x\}$ to denote the entire interpolated feature map, where $J\{x\}(t) \in \mathbb{R}^D$.

In the C-COT formulation, a continuous T -periodic multi-channel convolution filter $f = (f^1 \dots f^D)$ is trained

¹For clarity, we present the one-dimensional domain formulation. The generalization to higher dimensions, including images, is detailed in [12].

to predict the detection scores $S_f\{x\}(t)$ of the target as,

$$S_f\{x\} = f * J\{x\} = \sum_{d=1}^D f^d * J_d\{x^d\}. \quad (2)$$

The scores are defined in the corresponding image region $t \in [0, T)$ of the feature map $x \in \mathcal{X}$. In (2), the convolution of single-channel T -periodic functions is defined as $f * g(t) = \frac{1}{T} \int_0^T f(t - \tau)g(\tau) d\tau$. The multi-channel convolution $f * J\{x\}$ is obtained by summing the result of all channels, as defined in (2). The filters are learned by minimizing the following objective,

$$E(f) = \sum_{j=1}^M \alpha_j \|S_f\{x_j\} - y_j\|_{L^2}^2 + \sum_{d=1}^D \|w f^d\|_{L^2}^2. \quad (3)$$

The labeled detection scores $y_j(t)$ of sample x_j is set to a periodically repeated Gaussian function. The data term consists of the weighted classification error, given by the L^2 -norm $\|g\|_{L^2}^2 = \frac{1}{T} \int_0^T |g(t)|^2 dt$, where $\alpha_j \geq 0$ is the weight of sample x_j . The regularization integrates a spatial penalty $w(t)$ to mitigate the drawbacks of the periodic assumption, while enabling an extended spatial support [9].

As in previous DCF methods, a more tractable optimization problem is obtained by changing to the Fourier basis. Parseval's formula implies the equivalent loss,

$$E(f) = \sum_{j=1}^M \alpha_j \|\widehat{S_f\{x_j\}} - \hat{y}_j\|_{\ell^2}^2 + \sum_{d=1}^D \|\hat{w} * \hat{f}^d\|_{\ell^2}^2. \quad (4)$$

Here, the hat \hat{g} of a T -periodic function g denotes the Fourier series coefficients $\hat{g}[k] = \frac{1}{T} \int_0^T g(t) e^{-i \frac{2\pi}{T} kt} dt$ and the ℓ^2 -norm is defined by $\|\hat{g}\|_{\ell^2}^2 = \sum_{-\infty}^{\infty} |\hat{g}[k]|^2$. The Fourier coefficients of the detection scores (2) are given by the formula $\widehat{S_f\{x\}} = \sum_{d=1}^D \hat{f}^d X^d \hat{b}_d$, where X^d is the Discrete Fourier Transform (DFT) of x^d .

In practice, the filters f^d are assumed to have finitely many non-zero Fourier coefficients $\{\hat{f}^d[k]\}_{-K_d}^{K_d}$, where $K_d = \lfloor \frac{N_d}{2} \rfloor$. Eq. (4) then becomes a quadratic problem, optimized by solving the normal equations,

$$(A^H \Gamma A + W^H W) \hat{\mathbf{f}} = A^H \Gamma \hat{\mathbf{y}}. \quad (5)$$

Here, $\hat{\mathbf{f}}$ and $\hat{\mathbf{y}}$ are vectorizations of the Fourier coefficients in f^d and y_j , respectively. The matrix A exhibits a sparse structure, with diagonal blocks containing elements of the form $X_j^d[k] \hat{b}_d[k]$. Further, Γ is a diagonal matrix of the weights α_j and W is a convolution matrix with the kernel $\hat{w}[k]$. The C-COT [12] employs the Conjugate Gradient (CG) method [32] to iteratively solve (5), since it was shown to effectively utilize the sparsity structure of the problem.

3. Our Approach

As discussed earlier, over-fitting and computational bottlenecks in the DCF learning stem from common factors. We therefore proceed with a collective treatment of these issues, aiming at both improved performance *and* speed.

Robust learning: As mentioned earlier, the large number of optimized parameters in (3) may cause over-fitting due to limited training data. We alleviate this issue by introducing a factorized convolution formulation in section 3.1. This strategy radically reduces the number of model parameters by 80% in the case of deep features, while increasing tracking performance. Moreover, we propose a compact generative model of the sample distribution in section 3.2, that boosts diversity and avoids the previously discussed problems related to storing a large sample set. Finally, we investigate strategies for updating the model in section 3.3 and conclude that a less frequent update of the filter stabilizes the learning, which results in more robust tracking.

Computational complexity: The learning step is the computational bottleneck in optimization-based DCF trackers, such as C-COT. The computational complexity of the appearance model optimization in C-COT is obtained by analyzing the Conjugate Gradient algorithm applied to (5). The complexity can be expressed as $\mathcal{O}(N_{\text{CG}} D M \bar{K})$,² where N_{CG} is the number of CG iterations and $\bar{K} = \frac{1}{D} \sum_d K_d$ is the average number of Fourier coefficients per filter channel. Motivated by this complexity analysis of the learning, we propose methods for reducing D , M and N_{CG} in sections 3.1, 3.2, and 3.3 respectively.

3.1. Factorized Convolution Operator

We first introduce a factorized convolution approach, with the aim of reducing the number of parameters in the model. We observed that many of the filters f^d learned in C-COT contain negligible energy. This is particularly apparent for high-dimensional deep features, as visualized in figure 2. Such filters hardly contribute to target localization, but still affect the training time. Instead of learning one separate filter for each feature channel d , we use a smaller set of basis filters f^1, \dots, f^C , where $C < D$. The filter for feature layer d is then constructed as a linear combination $\sum_{c=1}^C p_{d,c} f^c$ of the filters f^c using a set of learned coefficients $p_{d,c}$. The coefficients can be compactly represented as a $D \times C$ matrix $P = (p_{d,c})$. The new multi-channel filter can then be written as the matrix-vector product Pf . We obtain the factorized convolution operator,

$$S_{Pf}\{x\} = Pf * J\{x\} = \sum_{c,d} p_{d,c} f^c * J_d\{x^d\} = f * P^T J\{x\}. \quad (6)$$

The last equality follows from the linearity of convolution. The factorized convolution (6) can thus alternatively

²See the supplementary material for a derivation.

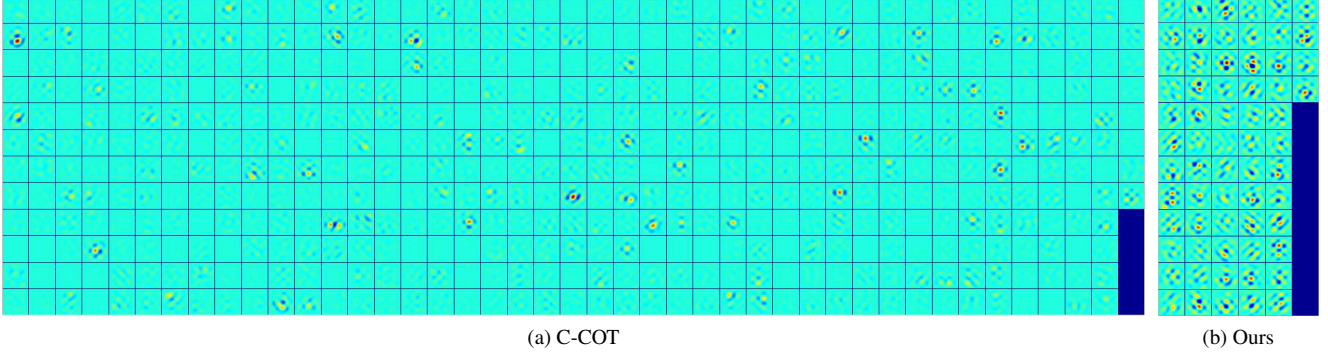


Figure 2. Visualization of the learned filters corresponding to the last convolutional layer in the deep network. We display all the 512 filters f^d learned by the baseline C-COT (a) and the reduced set of 64 filters f^c obtained by our factorized formulation (b). The vast majority of the baseline filters contain negligible energy, indicating irrelevant information in the corresponding feature layers. Our factorized convolution formulation learns a compact set of discriminative basis filters with significant energy, achieving a radical reduction of parameters.

be viewed as a two-step operation where the feature vector $J\{x\}(t)$ at each location t is first multiplied with the matrix P^T . The resulting C -dimensional feature map is then convolved with the filter f . The matrix P^T thus resembles a linear dimensionality reduction operator, as used in *e.g.* [13]. The key difference is that we learn the filter f and matrix P *jointly*, in a discriminative fashion, by minimizing the classification error (3) of the factorized operator (6).

For simplicity, we consider learning the factorized operator (6) from single training sample x . To simplify notation, we use $\hat{z}^d[k] = X^d[k]\hat{b}_d[k]$ to denote the Fourier coefficients of the interpolated feature map $z = J\{x\}$. The corresponding loss in the Fourier domain (4) is derived as,

$$E(f, P) = \left\| \hat{z}^T P \hat{f} - \hat{y} \right\|_{\ell^2}^2 + \sum_{c=1}^C \left\| \hat{w} * \hat{f}^c \right\|_{\ell^2}^2 + \lambda \|P\|_F^2. \quad (7)$$

Here we have added the Frobenius norm of P as a regularization, controlled by the weight parameter λ .

Unlike the original formulation (4), our new loss (7) is a non-linear least squares problem. Due to the bi-linearity of $\hat{z}^T P \hat{f}$, the loss (7) is similar to a matrix factorization problem [21]. Popular optimization strategies for these applications, including Alternating Least Squares, are however not feasible due to the parameter size and online nature of our problem. Instead, we employ Gauss-Newton [32] and use the Conjugate Gradient method to optimize the quadratic subproblems. The Gauss-Newton method is derived by linearizing the residuals in (7) using a first order Taylor series expansion. Here, this corresponds to approximating the bilinear term $\hat{z}^T P \hat{f}$ around the current estimate (\hat{f}_i, P_i) as,

$$\begin{aligned} \hat{z}^T (P_i + \Delta P) (\hat{f}_i + \Delta \hat{f}) &\approx \hat{z}^T P_i \hat{f}_{i,\Delta} + \hat{z}^T \Delta P \hat{f}_i \\ &= \hat{z}^T P_i \hat{f}_{i,\Delta} + (\hat{f}_i \otimes \hat{z})^T \text{vec}(\Delta P). \end{aligned} \quad (8)$$

Here, we set $\hat{f}_{i,\Delta} = \hat{f}_i + \Delta \hat{f}$. In the last equality, the Kronecker product \otimes is used to obtain a vectorization of the matrix step ΔP .

The Gauss-Newton subproblem at iteration i is derived by substituting the first-order approximation (8) into (7),

$$\begin{aligned} \tilde{E}(\hat{f}_{i,\Delta}, \Delta P) &= \left\| \hat{z}^T P_i \hat{f}_{i,\Delta} + (\hat{f}_i \otimes \hat{z})^T \text{vec}(\Delta P) - \hat{y} \right\|_{\ell^2}^2 \\ &+ \sum_{c=1}^C \left\| \hat{w} * \hat{f}_{i,\Delta}^c \right\|_{\ell^2}^2 + \mu \|P_i + \Delta P\|_F^2. \end{aligned} \quad (9)$$

Since the filter f is constrained to have finitely many non-zero Fourier coefficients, eq. (9) is a linear least squares problem. The corresponding normal equations have a partly similar structure to (5), with additional components corresponding to the matrix increment ΔP variable.³ We employ the Conjugate Gradient method to optimize each Gauss-Newton subproblem to obtain the new filter $\hat{f}_{i,\Delta}^*$ and matrix increment ΔP^* . The filter and matrix estimates are then updated as $\hat{f}_{i+1} = \hat{f}_{i,\Delta}^*$ and $P_{i+1} = P_i + \Delta P^*$.

The main objective of our factorized convolution operation is to reduce the computational and memory complexity of the tracker. Due to the adaptability of the filter, the matrix P can be learned just from the first frame. This has two important implications. Firstly, only the projected feature map $P^T J\{x_j\}$ requires storage, leading to significant memory savings. Secondly, the filter can be updated in subsequent frames using the projected feature maps $P^T J\{x_j\}$ as input to the method described in section 2. This reduces the linear complexity in the feature dimensionality D to the filter dimensionality C , *i.e.* $\mathcal{O}(N_{CG} C M \bar{K})$.

3.2. Generative Sample Space Model

Here, we propose a compact generative model of the sample set that averts the earlier discussed issues of storing a large set of recent training samples. Most DCF trackers, such as SRDCF [9] and C-COT [12], add one training sample x_j in each frame j . The weights are typically set to decay exponentially $\alpha_j \sim (1 - \gamma)^{M-j}$, controlled by

³See supplementary material for the derivation of the normal equations.



Figure 3. Visualization of the training set representation in the baseline C-COT (bottom row) and our method (top row). In C-COT, the training set consists of a sequence of consecutive samples. This introduces large redundancies due to slow change in appearance, while previous aspects of the appearance are forgotten. This can cause over-fitting to recent samples. Instead, we model the training data as a mixture of Gaussian components, where each component represent a different aspect of the appearance. Our approach yields a compact yet diverse representation of the data, thereby reducing the risk of over-fitting.

the learning rate γ . If the number of samples has reached a maximum limit M_{\max} , the sample with the smallest weight α_j is replaced. This strategy however requires a large sample limit M_{\max} to obtain a representative sample set.

We observe that collecting a new sample in each frame leads to large redundancies in the sample set, as visualized in figure 3. The standard sampling strategy (bottom row) populates the whole training set with similar samples x_j , despite containing almost the same information. Instead, we propose to use a probabilistic generative model of the sample set that achieves a compact description of the samples by eliminating redundancy and enhancing variety (top).

Our approach is based on the joint probability distribution $p(x, y)$ of the sample feature maps x and corresponding desired outputs scores y . Given $p(x, y)$, the intuitive objective is to find the filter that minimizes the expected correlation error. This is obtained by replacing (3) with

$$E(f) = \mathbb{E} \left\{ \|S_f\{x\} - y\|_{L^2}^2 \right\} + \sum_{d=1}^D \|wf^d\|_{L^2}^2. \quad (10)$$

Here, the expectation \mathbb{E} is evaluated over the joint sample distribution $p(x, y)$. Note that the original loss (3) is obtained as a special case by estimating the sample distribution as $p(x, y) = \sum_{j=1}^M \alpha_j \delta_{x_j, y_j}(x, y)$, where δ_{x_j, y_j} denotes the Dirac impulse at the training sample (x_j, y_j) .⁴ Instead, we propose to estimate a compact model of the sam-

⁴We can without loss of generality assume the weights α_j sum to one.

ple distribution $p(x, y)$ that leads to a more efficient approximation of the expected loss (10).

We observe that the shape of the desired correlation output y for a sample x is predetermined, here as a Gaussian function. The label functions y_j in (3) only differ by a translation that aligns the peak with the target center. This alignment is equivalently performed by shifting the feature map x . We can thus assume that the target is centered in the image region and that all $y = y_0$ are identical. Hence, the sample distribution can be factorized as $p(x, y) = p(x)\delta_{y_0}(y)$ and we only need to estimate $p(x)$. For this purpose we employ a Gaussian Mixture Model (GMM) such that $p(x) = \sum_{l=1}^L \pi_l \mathcal{N}(x; \mu_l; I)$. Here, L is the number of Gaussian components $\mathcal{N}(x; \mu_l; I)$, π_l is the prior weight of component l , and $\mu_l \in \mathcal{X}$ is its mean. The covariance matrix is set to the identity matrix I to avoid costly inference in the high-dimensional sample space.

To update the GMM, we use a simplified version of the online algorithm by Declercq and Piater [14]. Given a new sample x_j , we first initialize a new component m with $\pi_m = \gamma$ and $\mu_m = x_j$ (concatenate in [14]). If the number of components exceeds the limit L , we simplify the GMM. We discard a component if its weight π_l is below a threshold. Otherwise, we merge the two closest components k and l into a common component n [14],

$$\pi_n = \pi_k + \pi_l, \quad \mu_n = \frac{\pi_k \mu_k + \pi_l \mu_l}{\pi_k + \pi_l}. \quad (11)$$

The required distance comparisons $\|\mu_k - \mu_l\|$ are efficiently computed in the Fourier domain using Parseval's formula.

Finally, the expected loss (10) is approximated as,

$$E(f) = \sum_{l=1}^L \pi_l \|S_f\{\mu_l\} - y_0\|_{L^2}^2 + \sum_{d=1}^D \|wf^d\|_{L^2}^2. \quad (12)$$

Note that the Gaussian means μ_l and the prior weights π_l directly replace x_j and α_j , respectively, in (3). So, the same training strategy as described in section 2 can be applied.

The key difference in complexity compared to (3) is that the number of samples has decreased from M to L . In our experiments, we show that the number of components L can be set to $M/8$, while obtaining an improved tracking performance. Our sample distribution model $p(x, y)$ is combined with the factorized convolution from section 3.1 by replacing the sample x with the projected sample $P^T Jx$. The projection does not affect our formulation since the matrix P is constant after the first frame.

3.3. Model Update Strategy

The standard approach in DCF based tracking is to update the model in each frame [4, 9, 20]. In C-COT, this implies optimizing (3) after each new sample is added, by iteratively solving the normal equations (5). Iterative optimization based DCF methods exploit that the loss function changes gradually between frames. The current estimate of the filter therefore provides a good initialization of the iterative search. Still, updating the filter in each frame have a severe impact on the computational load.

Instead of updating the model in a continuous fashion every frame, we use a sparser updating scheme, which is a common practice in non-DCF trackers [31, 38]. Intuitively, an optimization process should only be started once sufficient change in the objective has occurred. However, finding such conditions is non-trivial and may lead to unnecessarily complex heuristics. Moreover, optimality conditions based on the gradient of the loss (3), given by the residual of (5), are expensive to evaluate in practice. We therefore avoid explicitly detecting changes in the objective and simply update the filter by starting the optimization process in every N_S th frame. The parameter N_S determines how often the filter is updated, where $N_S = 1$ corresponds to optimizing the filter in every frame, as in standard DCF methods. In every N_S th frame, we perform a fixed number of N_{CG} Conjugate Gradient iterations to refine the model. As a result, the average number of CG iterations per frame is reduced to N_{CG}/N_S , which has a substantial effect on the overall computational complexity of the learning. Note that N_S does not affect the updating of the sample space model, introduced in section 3.2, which is updated every frame.

To our initial surprise, we observed that a moderately infrequent update of the model ($N_S \approx 5$) generally led to improved tracking results. We mainly attribute this effect to reduced over-fitting to the recent training samples. By postponing the model update a few frames, the loss is updated

	Conv-1	Conv-5	HOG	CN
Feature dimension, D	96	512	31	11
Filter dimension, C	16	64	10	3

Table 1. The settings of the proposed factorized convolution approach, as employed in our experiments. For each feature, we show the dimensionality D and the number of filters C .

by adding a new mini-batch to the training samples, instead of only a single one. This might contribute to stabilizing the learning, especially in scenarios where a new sample is affected by sudden changes, such as out-of-plane rotations, deformations, clutter, and occlusions (see figure 1).

While increasing N_S leads to reduced computations, it may also reduce the convergence speed of the optimization, resulting in a less discriminative model. A naive compensation by increasing the number of CG iterations N_{CG} would counteract the achieved computational gains. Instead, we aim to achieve a faster convergence by better adapting the CG algorithm to online tracking, where the loss changes dynamically. This is obtained by substituting the standard Fletcher-Reeves formula to the Polak-Ribière formula [34] for finding the momentum factor, since it has shown improved convergence rates for inexact and flexible preconditioning [18], which have similarities to our scenario.

4. Experiments

We validate our proposed formulation by performing comprehensive experiments on four benchmarks: VOT2016 [23], UAV123 [29], OTB-2015 [37], and TempleColor [26].

4.1. Implementation Details

Our tracker is implemented in Matlab. We apply the same feature representation as C-COT, namely a combination of the first (Conv-1) and last (Conv-5) convolutional layer in the VGG-m network [5], along with HOG [6] and Color Names (CN) [35]. For the factorized convolution presented in section 3.1, we learn one coefficient matrix P for each feature type. The settings for each feature is summarized in table 1. The regularization parameter λ in (7) is set to $2 \cdot 10^{-7}$. The loss (7) is optimized in the first frame using 10 Gauss-Newton iterations and 20 CG iterations for the subproblems (9). In the first iteration $i = 0$, the filter \hat{f}_0 is initialized to zero. To preserve the deterministic property of the tracker, we initialize the coefficient matrix P_0 by PCA, though we found random initialization to be equally robust.

For the sample space model, presented in section 3.2, we set the learning rate to $\gamma = 0.012$. The number of components are set to $L = 50$, which represents an 8-fold reduction compared to the number of samples ($M = 400$) used in C-COT. We update the filter in every $N_S = 6$ frame (section 3.3). We use the same number of $N_{CG} = 5$ Conjugate Gradient iterations as in C-COT. Note that *all* parameters settings are kept fixed for all videos in a dataset.

	Baseline C-COT \implies (Sec. 2)	Factorized Convolution \implies (Sec. 3.1)	Sample Space Model \implies (Sec. 3.2)	Model Update \implies (Sec. 3.3)
EAO	0.331	0.342	0.352	0.374
FPS	0.3	1.1	2.6	6.0
Compl. change	-	$D \rightarrow C$	$M \rightarrow L$	$N_{CG} \rightarrow \frac{N_{CG}}{N_S}$
Compl. red.	-	$6\times$	$8\times$	$6\times$

Table 2. Analysis of our approach on the VOT2016. The impact of progressively integrating one contribution at the time, from left to right, is displayed. We show the performance in Expected Average Overlap (EAO) and speed in FPS (benchmarked on a single CPU). We also summarize the reduction in learning complexity $\mathcal{O}(N_{CG}DM\bar{K})$ obtained in each step, both symbolically and in absolute numbers (bottom row) using our settings. Our contributions systematically improve both performance and speed.

4.2. Baseline Comparison

Here, we analyze our approach on the VOT2016 benchmark by demonstrating the impact of progressively integrating our contributions. The VOT2016 dataset consists of 60 videos compiled from a set of more than 300 videos. The performance is evaluated both in terms of accuracy (average overlap during successful tracking) and robustness (failure rate). The overall performance is evaluated using Expected Average Overlap (EAO) which accounts for both accuracy and robustness. We refer to [24] for details.

Table 2 shows an analysis of our contributions. The integration of our factorized convolution into the baseline leads to a performance improvement and a significant reduction in complexity ($6\times$). The sample space model further improves the performance by a relative gain of 2.9% in EAO, while reducing the learning complexity by a factor of 8. Additionally incorporating our proposed model update elevates us to an EAO score of 0.374, leading to a final relative gain of 13.0% compared to the baseline. In table 2 we also show the impact on the tracker speed achieved by our contributions. For a fair comparison, we report the FPS measured on a single CPU for all entries in the table, without accounting for feature extraction time. Each of our contributions systematically improves the speed of the tracker, combining to a 20-fold final gain compared to the baseline. When including all steps (also feature extraction), the GPU version of our tracker operates at 8 FPS.

We found the settings in table 1 to be insensitive to minor changes. Substantial gain in speed can be obtained by reducing the number of filters C , at the cost of a slight reduction in performance. To further analyze the impact of our jointly learned factorized convolution approach, we compare with applying PCA in the first frame to obtain the matrix P . PCA degrades the EAO from 0.331 to 0.319, while our discriminative learning based method achieves 0.342.

We observed that our sample model provides consistently better results compared to the training sample set management employed in C-COT when using the same number of components and samples ($L = M$). This

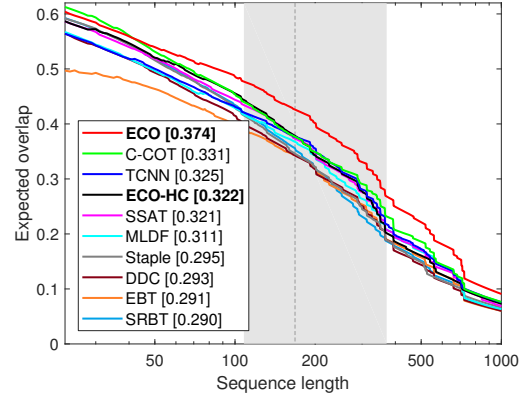


Figure 4. Expected Average Overlap (EAO) curve on VOT2016. Only the top 10 trackers are shown for clarity. The EAO measure, computed as the average EAO over typical sequence lengths (grey region), is displayed in the legend (see [24] for details).

	SRBT [23]	EBT [39]	DDC [23]	Staple [1]	MLDF [23]	SSAT [23]	TCNN [30]	C-COT [12]	ECO-HC Ours	ECO Ours
EAO	0.290	0.291	0.293	0.295	0.311	0.321	0.325	0.331	0.322	0.374
Fail. rt.	1.25	0.90	1.23	1.35	0.83	1.04	0.96	0.85	1.08	0.72
Acc.	0.50	0.44	0.53	0.54	0.48	0.57	0.54	0.52	0.53	0.54
EFO	3.69	3.01	0.20	11.14	1.48	0.48	1.05	0.51	15.13	4.53

Table 3. State-of-the-art in terms of expected average overlap (EAO), robustness (failure rate), accuracy, and speed (in EFO units) on the VOT2016 dataset. Only the top-10 trackers are shown. Our deep feature based ECO achieve superior EAO, while our hand-crafted feature version (ECO-HC) has the best speed.

is particularly evident for a smaller number of components/samples: When reducing the number of samples from $M = 400$ to $M = 50$ in the standard approach, the EAO decreases from 0.342 to 0.338 (-1.2%). Instead, when using our approach with $L = 50$ components, the EAO increases by $+2.9\%$ to 0.351. In case of the model update, we observed an upward trend in performance when increasing N_S from 1 to 6. When increasing N_S further, a gradual downward trend was observed. We therefore use $N_S = 6$ throughout our experiments.

4.3. State-of-the-art Comparison

Here, we compare our approach with state-of-the-art trackers on four challenging tracking benchmarks. Detailed results are provided in the supplementary material.

VOT2016 Dataset: In table 3 we compare our approach, in terms of expected average overlap (EAO), robustness, accuracy and speed (in EFO units), with the top-ranked trackers in the VOT2016 challenge. The first-ranked performer in VOT2016 challenge, C-COT, provides an EAO score of 0.331. Our approach achieves a relative gain of 13.0% in EAO compared to C-COT. Further, our ECO tracker achieves the best failure rate of 0.72 while maintaining a competitive accuracy. We also report the total speed in terms of EFO, which normalizes the speed with respect to hardware performance. Note that EFO also takes feature ex-

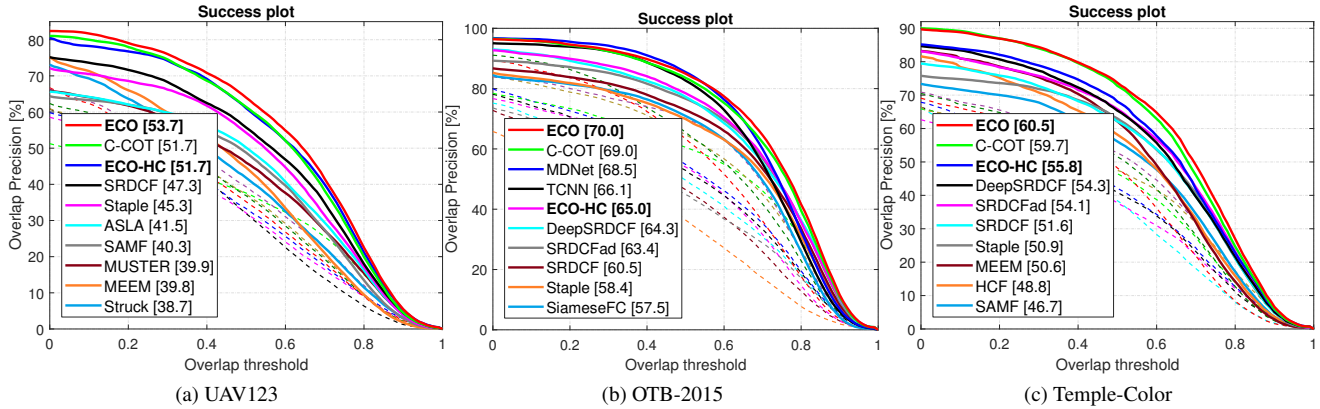


Figure 5. Success plots on the UAV-123 (a), OTB-2015 (b) and TempleColor (c) datasets. Only the top 10 trackers are shown in the legend for clarity. The AUC score of each tracker is shown in the legend. Our approach significantly improves the state-of-the-art on all datasets.

traction time into account, a major additive complexity that is independent of our DCF improvements. In the comparison, our tracker ECO-HC using only hand-crafted features (HOG and Color Names) achieves the best speed. Among the top three trackers in the challenge, which are all based on deep features, TCNN [30] obtains the best speed with an EFO of 1.05. Our deep feature version (ECO) achieves an almost 5-fold speedup in EFO and a relative performance improvement of 15.1% in EAO compared to TCNN. Figure 4 displays the EAO curves of the top-10 trackers.

UAV123 Dataset: Aerial tracking using unmanned aerial vehicles (UAVs) has received much attention recently, with many vision applications, including wild-life monitoring, search and rescue, navigation, and crowd surveillance. In these applications, persistent UAV navigation is required, for which real-time tracking output is crucial. In such cases, the desired tracker should be accurate and robust, while operating in real-time under limited hardware capabilities, *e.g.*, CPUs or mobile GPU platforms. We therefore introduce a real-time variant of our method (ECO-HC), based on hand-crafted features (HOG and Color Names), operating at 60 FPS on a single i7 CPU (including feature extraction).

We evaluate our trackers on the recently introduced aerial video benchmark, UAV123 [29], for low altitude UAV target tracking. The dataset consists of 123 aerial videos with more than 110K frames. The trackers are evaluated using success plot [36], calculated as percentage of frames with an intersection-over-union (IOU) overlap exceeding a threshold. Trackers are ranked using the area-under-the-curve (AUC) score. Figure 5a shows the success plot over all the 123 videos in the dataset. We compare with all tracking results reported in [29] and further add Staple [1], due to its high frame-rate, and C-COT [12]. Among the top 5 compared trackers, only Staple runs at real-time, with an AUC score of 45.3%. Our ECO-HC tracker also operates in real-time (60 FPS), with an AUC score of 51.7%, significantly outperforming Staple by 6.4%. C-COT obtains an AUC score of 51.7%. Our ECO outperforms C-COT,

achieving an AUC score of 53.7%, using same features.

OTB2015 Dataset: We compare our tracker with 20 state-of-the-art methods: TLD [22], Struck [19], CFLB [16], ACT [13], TGPR [17], KCF [20], DSST [7], SAMF [25], MEEM [38], DAT [33], LCT [28], HCF [27], SRDCF [9], SRDCFad [10], DeepSRDCF [8], Staple [1], MDNet [31], SiameseFC [2], TCNN [30] and C-COT [12].

Figure 5b shows the success plot over all the 100 videos in the OTB-2015 dataset [37]. Among the compared trackers using hand-crafted features, SRDCFad provides the best results with an AUC score of 63.4%. Our proposed method, ECO-HC, also employing hand-crafted features outperforms SRDCFad with an AUC score of 65.0%, while running on a CPU with a speed of 60 FPS. Among the compared deep feature trackers, C-COT, MDNet and TCNN provide the best results with AUC scores of 69.0%, 68.5% and 66.1% respectively. Our approach ECO, provides the best performance with an AUC score of 70.0%.

TempleColor Dataset: In figure 5c we present results on the TempleColor dataset [26] containing 128 videos. Our method again achieves a substantial improvement over C-COT, with a gain of 0.8% in AUC.

5. Conclusions

We revisit the core DCF formulation to counter the issues of over-fitting and computational complexity. We introduce a factorized convolution operator to reduce the number of parameters in the model. We also propose a compact generative model of the training sample distribution to drastically reduce memory and time complexity of the learning, while enhancing sample diversity. Lastly, we suggest a simple yet effective model update strategy that reduces over-fitting to recent samples. Experiments on four datasets demonstrate state-of-the-art performance with improved frame rate.

Acknowledgments: This work has been supported by SSF (SymbiCloud), VR (EMC², starting grant 2016-05543), SNIC, WASP, Visual Sweden, and Nvidia.

References

- [1] L. Bertinetto, J. Valmadre, S. Golodetz, O. Miksik, and P. H. S. Torr. Staple: Complementary learners for real-time tracking. In *CVPR*, 2016.
- [2] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr. Fully-convolutional siamese networks for object tracking. In *ECCV workshop*, 2016.
- [3] A. Bibi, M. Mueller, and B. Ghanem. Target response adaptation for correlation filter tracking. In *ECCV*, 2016.
- [4] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui. Visual object tracking using adaptive correlation filters. In *CVPR*, 2010.
- [5] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. In *BMVC*, 2014.
- [6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005.
- [7] M. Danelljan, G. Häger, F. Shahbaz Khan, and M. Felsberg. Accurate scale estimation for robust visual tracking. In *BMVC*, 2014.
- [8] M. Danelljan, G. Häger, F. Shahbaz Khan, and M. Felsberg. Convolutional features for correlation filter based visual tracking. In *ICCV Workshop*, 2015.
- [9] M. Danelljan, G. Häger, F. Shahbaz Khan, and M. Felsberg. Learning spatially regularized correlation filters for visual tracking. In *ICCV*, 2015.
- [10] M. Danelljan, G. Häger, F. Shahbaz Khan, and M. Felsberg. Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking. In *CVPR*, 2016.
- [11] M. Danelljan, G. Häger, F. Shahbaz Khan, and M. Felsberg. Discriminative scale space tracking. *TPAMI*, PP(99), 2016.
- [12] M. Danelljan, A. Robinson, F. Shahbaz Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. In *ECCV*, 2016.
- [13] M. Danelljan, F. Shahbaz Khan, M. Felsberg, and J. van de Weijer. Adaptive color attributes for real-time visual tracking. In *CVPR*, 2014.
- [14] A. Declercq and J. H. Piater. Online learning of gaussian mixture models - a two-level approach. In *VISAPP*, 2008.
- [15] H. K. Galoogahi, T. Sim, and S. Lucey. Multi-channel correlation filters. In *ICCV*, 2013.
- [16] H. K. Galoogahi, T. Sim, and S. Lucey. Correlation filters with limited boundaries. In *CVPR*, 2015.
- [17] J. Gao, H. Ling, W. Hu, and J. Xing. Transfer learning based visual tracking with gaussian process regression. In *ECCV*, 2014.
- [18] G. H. Golub and Q. Ye. Inexact preconditioned conjugate gradient method with inner-outer iteration. *SIAM J. Scientific Computing*, 21(4):1305–1320, 1999.
- [19] S. Hare, A. Saffari, and P. Torr. Struck: Structured output tracking with kernels. In *ICCV*, 2011.
- [20] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista. High-speed tracking with kernelized correlation filters. *TPAMI*, 37(3):583–596, 2015.
- [21] J. Hyeong Hong and A. Fitzgibbon. Secrets of matrix factorization: Approximations, numerics, manifold optimization and random restarts. In *ICCV*, 2015.
- [22] Z. Kalal, J. Matas, and K. Mikolajczyk. P-n learning: Bootstrapping binary classifiers by structural constraints. In *CVPR*, 2010.
- [23] M. Kristan, A. Leonardis, J. Matas, R. Felsberg, Pflugfelder, M., L. Čehovin, G. Vojír, T. Häger, and et al. The visual object tracking vot2016 challenge results. In *ECCV workshop*, 2016.
- [24] M. Kristan, J. Matas, A. Leonardis, M. Felsberg, L. Čehovin, G. Fernández, T. Vojír, G. Nebehay, R. Pflugfelder, and G. Häger. The visual object tracking vot2015 challenge results. In *ICCV workshop*, 2015.
- [25] Y. Li and J. Zhu. A scale adaptive kernel correlation filter tracker with feature integration. In *ECCV Workshop*, 2014.
- [26] P. Liang, E. Blasch, and H. Ling. Encoding color information for visual tracking: Algorithms and benchmark. *TIP*, 24(12):5630–5644, 2015.
- [27] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang. Hierarchical convolutional features for visual tracking. In *ICCV*, 2015.
- [28] C. Ma, X. Yang, C. Zhang, and M.-H. Yang. Long-term correlation tracking. In *CVPR*, 2015.
- [29] M. Mueller, N. Smith, and B. Ghanem. A benchmark and simulator for uav tracking. In *ECCV*, 2016.
- [30] H. Nam, M. Baek, and B. Han. Modeling and propagating cnns in a tree structure for visual tracking. *CoRR*, abs/1608.07242, 2016.
- [31] H. Nam and B. Han. Learning multi-domain convolutional neural networks for visual tracking. In *CVPR*, 2016.
- [32] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- [33] H. Possegger, T. Mauthner, and H. Bischof. In defense of color-based model-free tracking. In *CVPR*, 2015.
- [34] J. R. Shewchuk. An introduction to the conjugate gradient method without the agonizing pain. Technical report, Pittsburgh, PA, USA, 1994.
- [35] J. van de Weijer, C. Schmid, J. J. Verbeek, and D. Larlus. Learning color names for real-world applications. *TIP*, 18(7):1512–1524, 2009.
- [36] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark. In *CVPR*, 2013.
- [37] Y. Wu, J. Lim, and M.-H. Yang. Object tracking benchmark. *TPAMI*, 37(9):1834–1848, 2015.
- [38] J. Zhang, S. Ma, and S. Sclaroff. MEEM: robust tracking via multiple experts using entropy minimization. In *ECCV*, 2014.
- [39] G. Zhu, F. Porikli, and H. Li. Beyond local search: Tracking objects everywhere with instance-specific proposals. In *CVPR*, 2016.