

Multi-kernel Correlation Filter for Visual Tracking

Ming TANG and Jiayi FENG

National Lab of Pattern Recognition, Institute of Automation,
Chinese Academy of Sciences, Beijing 100190, China

Ftangm, jyfeng@nlpr.ia.ac.cn

Abstract

Correlation filter based trackers are ranked top in terms of performances. Nevertheless, they only employ a single kernel at a time. In this paper, we will derive a multi-kernel correlation filter (MKCF) based tracker which fully takes advantage of the invariance-discriminative power spectrums of various features to further improve the performance. Moreover, it may easily introduce location and representation errors to search several discrete scales for the proper one of the object bounding box, because normally the discrete candidate scales are determined and the corresponding feature pyramid are generated ahead of searching. In this paper, we will propose a novel and efficient scale estimation method based on optimal bisection search and fast evaluation of features. Our scale estimation method is the first one that uses the truly minimal number of layers of feature pyramid and avoids constructing the pyramid before searching for proper scales.

1. Introduction

Visual object tracking is one of the most challenging problems in computer vision [32, 18, 20, 26, 23]. To adapt to unpredictable variations of object appearance and background during tracking, one strategy is to select a single strong feature that is robust to any variation. However, this has been known to be difficult [33, 11], especially for the model-free tracking task in which no prior knowledge about the target object is known except for the initial frame. Therefore, designing an efficient scheme to combine several complementary features is a natural alternative.

In such a natural scheme, different features would capture different channels of target information and result in a better performance [35, 37, 21, 7]. Since any feature possesses its own invariance and discriminative power, what really distinguishes one feature from another is its location at

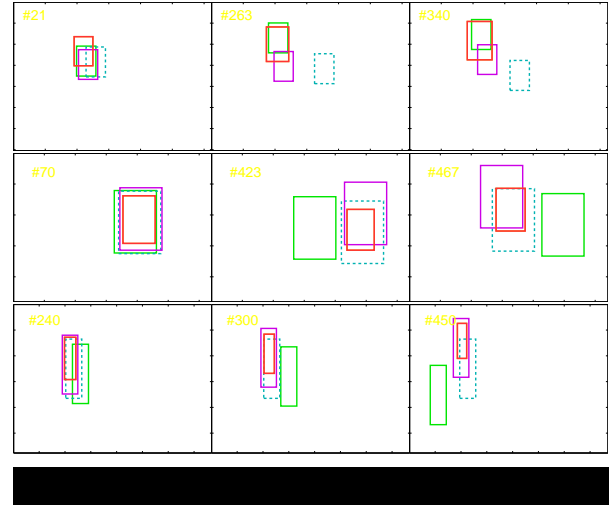


Figure 1. Qualitative comparisons of our Multi-kernel Correlation Filter (MKCF) with state-of-the-art trackers, Struck [13], KCF [15], and CN2 [7] on 3 sequences [36]. The figure is best viewed in color.

the invariance-discriminative power spectrum [33]. And the location may vary for different types of variations. Complementary features should be selected in such a way that they are not at the same location of the spectrum for the same appearance variation. In other words, complementary features must be at different spectrum locations under the same variation, especially in long term tracking because the object appearance and its local background may vary substantially. Concatenation of several features into a single kernel space, however, may confuse their characteristics on the invariance-discriminative power spectrums, and is not a good choice. Consequently, a multiple kernel feature, i.e., a combination of several kernels, one for each feature, is a natural choice.

In recent years, correlation filter based trackers have been proposed [4, 15, 7, 6, 16] and achieved top performance. Bolme *et al.*[4] proposed a correlation filter based tracker, called Minimum Output Sum of Squared Error (MOSSE), with classical signal processing techniques.

This work was supported by Natural Science Foundation of China. Grant No. 61375035.

The corresponding author of the paper is Ming TANG.

They used a base image patch and several circulant virtual ones to train the filter directly in the Fourier domain. Henriques *et al.* [15, 16] utilized the circulant structure produced by a base sample to propose a kernelized correlation filter (KCF). The KCF used a single kernel and enabled fast learning with fast Fourier transforms instead of costly matrix operation, providing the highest tracking speed [36]. Danelljan *et al.* [7] extended the KCF with low-dimensional adaptive color channels and achieved state-of-the-art performance in a comprehensive evaluation. Nevertheless, the correlation filter based tracker can not yet utilize multiple kernels simultaneously. In this paper, we will extend the correlation filter based tracker further to its multiple kernel version and demonstrate its superior performance. Fig. 1 shows a qualitative comparison to indicate that our approach, MKCF, outperforms other trackers in challenging sequences Shaking, Trellis, and Walking2 [36].

Another challenging problem in visual tracking is the robust estimation of object scale in complex scenes. Currently, there exist two popular strategies to tackle this problem. One is the analytical method [2, 27] which moves and analytically deforms a template to minimize the difference between the template and an image region, and the other exhaustively searches for the proper scale among several discretized scales [1, 6]. In order to handle large variations of scale effectively and efficiently, Danelljan *et al.* [6] applied the idea of kernelized correlation filter to a pyramid representation of candidate bounding boxes to speed up exhaustive searches greatly. Nevertheless, their approach had to construct the pyramid in advance as the search method does, thus was in fact still based on several discrete scales. And in order to estimate object scales accurately, the pyramid had to include 33 layers. In this paper, we will present a novel method, which is based on fast feature pyramids [8] and optimization technique to efficiently determine the optimal scale in the continuous scale space. As far as we know, our scale estimation method is the first one which explores the truly minimal number of layers of feature pyramid and avoids constructing the pyramid before searching for proper scales. Moreover, our optimal scale estimation is generic, and can be incorporated into any tracker which does not contain inherent scale estimation.

In summary, the **main contributions** of our work includes two aspects. 1) A multi-kernel correlation filter (MKCF) based tracker is proposed. In fact, the MKCF can be accepted as a general framework of correlation filter in the sense that it embraces both strengths of multiple channels and multiple kernels. 2) An optimal and efficient algorithm is developed to determine the scale of target object. In this paper, the power law of image scaling [30, 8] is introduced into visual tracking community for the first time in order to quickly determine the object scale.

The remainder of this paper is organized as follows. In

Sec.2, we briefly overview the related work. In Sec.3, the general correlation filter framework of multiple kernels and multiple channels is derived. And the optimal and efficient algorithm to search for proper scales is presented. Our tracking algorithm, MKCF, is then described in detail. To understand the MKCF more clearly, Sec.4 provides details of our implementation. Experimental results and comparison with other state-of-the-art approaches are presented in Sec.5. Sec.6 summarizes our work.

2. Related Work

Multiple kernel learning (MKL) aims at simultaneously learning a kernel and the associated predictor in supervised learning settings. Rakotomamonjy *et al.* [28] proposed an efficient algorithm, named SimpleMKL, for solving the MKL problem through reduced gradient descent in a primal formulation. Varma and Ray [33] extended the MKL formulation in [28] by introducing an additional constraint on combinational coefficients and applied it to object classification. Vedaldi *et al.* [34] and Gehler and Nowozin [11] applied MKL based approaches to object detection and classification. Cortes *et al.* [5] studied the problem of learning kernels of the same family with an L_2 regularization for ridge regression (RR) [29]. In this paper, we extend the MKL formulation in [28] to RR, and present a novel multi-kernel RR approach.

In recent years, Bolme *et al.* [4] proposed an extension of traditional correlation filters [19] referred to as Minimum Output Sum of Squared Error (MOSSE) filter. The original MOSSE was expressed in the Fourier domain. In fact, it is easy to observe that the expression of MOSSE in the spatial domain is just the ridge regression [29] with a linear kernel. Therefore, Henriques *et al.* [15] proposed the kernelized correlation filter (KCF) by introducing the kernel trick into ridge regression [29]. The generalizations of MOSSE and KCF to multiple channels have also been proposed [3, 10, 14]. Danelljan *et al.* [7] further extended the KCF to multiple channels, one per adaptive color attribute. Henriques *et al.* [17] utilized the circulant structure of Gram matrix to speed up the training of pose detectors in the Fourier domain. It is noted that all these works [4, 15, 3, 10, 14, 7] can only employ a single kernel. In this paper, we propose a brand new multi-kernel correlation filter which is able to fully take advantage of invariance-discriminative power spectrums of various features.

Ruderman and Bialek [30] explored how the statistics of natural images behave as a function of the scale at which an image ensemble is captured. And their discovery implies that the ratio of the statistics of two image sets under two different scales is approximately the power function of the ratio of the two scales. Dollár [8] extended this power law to a pair of images of the same scene, and pointed out that the feature ratio of two images under different scales is also

approximately the power function of the ratio of the two scales. In this paper, we utilize Dollár's power law to speed up the determination of the object scale during tracking. To the best of our knowledge, it is the first time to make use of the power law in the visual tracking community.

3. Multi-kernel Correlation Filter Based Tracking

3.1. Training Multi-kernel Correlation Filter

The training goal of ridge regression [29] is to find function $f(x)$ which minimizes the squared error over training samples x_i 's and their regression targets y_i 's, *i.e.*,

$$\min_f \frac{1}{2} \sum_{i=0}^{l-1} (f(x_i) - y_i)^2 + \lambda \|f\|_K^2, \quad (1)$$

where l is the number of samples, f lies in a bounded convex subset of a Reproducing Kernel Hilbert Space defined by a positive definite kernel function $k(\cdot, \cdot)$, and $\lambda > 0$ is the regularization parameter. By means of the Representer Theorem [31], the solution f to the Tikhonov regularization problem can be expressed as

$$f(x) = \sum_{i=0}^{l-1} \alpha_i k(x_i, x). \quad (2)$$

Then, $\|f\|_K^2 = \alpha^T K \alpha$, where $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_{l-1})^T$, and K is the positive semi-definite kernel matrix with $K_{ij} = k(x_i, x_j)$ as its elements. Problem (1) becomes to find α , *i.e.*,

$$\min_{\alpha} \frac{1}{2} \|\mathbf{y} - K\alpha\|_2^2 + \frac{\lambda}{2} \alpha^T K \alpha, \quad (3)$$

where $\mathbf{y} = (y_0, y_1, \dots, y_{l-1})^T$.

It has been shown that using multiple kernels instead of a single one can improve the algorithms' discrimination [22, 33]. Given the base kernels, k_m where $m = 1, 2, \dots, M$, a usual approach is to consider $k(x_i, x_j)$ be a convex combination of base kernels, *i.e.*, $k(x_i, x_j) = \sum_{m=1}^M d_m k_m(x_i, x_j)$, where $k(x_i, x_j) = (k_1(x_i, x_j), k_2(x_i, x_j), \dots, k_M(x_i, x_j))^T$, $\mathbf{d} = (d_1, d_2, \dots, d_M)^T$, $\sum_{m=1}^M d_m = 1$, and $d_m \geq 0$. Hence we have

$$K = \sum_{m=1}^M d_m K_m, \quad (4)$$

where K_m is the m^{th} base kernel matrix with $K_{ij}^m = k_m(x_i, x_j)$ as its elements. Substituting K in Eq. (4) for that in (3) and introducing the additional constraint on the sum of d_m 's, we obtain the object function $F(\mathbf{d})$ of multiple kernel version of ridge regression problem (3) as follows.

low.

$$\min_{\mathbf{d}} F(\mathbf{d}) = \frac{1}{2} \|\mathbf{y} - \sum_{m=1}^M d_m^2 K_m\|_2^2 + \frac{\lambda}{2} \sum_{m=1}^M d_m^2 - 1, \quad (5)$$

where $\lambda = 10^{-3}$ and $\lambda = 10^{-2}$ in our current experiments. In order to ensure that all combination coefficients are positive, d_m^2 's, instead of d_m 's, are used in Eq. (5), *i.e.*, $K = \sum_{m=1}^M d_m^2 K_m$. It is noted that Eq. (5) is equivalent to a constrained multiple kernel optimization problem with the third item of Eq. (5) as its constraint and λ as its Lagrangian multiplier, and its optimal solution can be expressed as

$$f(x) = \sum_{i=0}^{l-1} \alpha_i d^2 k(x_i, x), \quad (6)$$

where $\mathbf{d}^2 = (d_1^2, d_2^2, \dots, d_M^2)^T$.

It should be pointed out that Eq. (5) is also the Lagrangian function with λ as its Lagrangian multiplier without a second power on $\sum_{m=1}^M d_m^2 - 1$. In this case, $\min_{\mathbf{d}} F(\mathbf{d})$ is a linear programming (LP) problem w.r.t. \mathbf{d}^2 , given λ . Since the optimal solution of LP problem will always be at the vertex of linear feasible region, the optimal \mathbf{d}^2 must be a unit vector. This means that the combination of multiple kernels will be discarded and only one kernel left. This case does not meet our goal of exploring multiple kernels simultaneously to improve the tracking performance. Therefore, a second power on $\sum_{m=1}^M d_m^2 - 1$ is necessary for Eq. (5).

Theorem 1 Let $\{K_m\}$ be positive semi-definite. Then, (a) given \mathbf{d} , $F(\mathbf{d})$ is convex w.r.t. \mathbf{d}^2 ; (b) given \mathbf{d}^2 , $F(\mathbf{d})$ is convex w.r.t. \mathbf{d} .

$F(\mathbf{d})$ is a differentiable function, we can find a minimizer simply by taking its gradients w.r.t. \mathbf{d} and then solving them. To solve (5), we let $F(\mathbf{d}) = 0$, and achieve that

$$\mathbf{d} = \left(\sum_{m=1}^M d_m^2 K_m + \lambda \mathbf{I} \right)^{-1} \mathbf{y}, \quad (7)$$

where \mathbf{I} is the $l \times l$ identity matrix. We propose two approaches to find out \mathbf{d}^2 . One is an iteration approach, and the other is analytic. The iteration approach employs the gradient descent method to achieve the minimum step by step, *i.e.*,

$$\begin{aligned} \mathbf{d}_{t+1}^2 &= \mathbf{d}_t^2 + \eta_{d,t} \mathbf{d}_t^2 F(\mathbf{d}_t) \\ &= \mathbf{d}_t^2 + \frac{\eta_{d,t}}{2} (\mathbf{d}_t^2 \mathbf{B} - \mathbf{c}), \end{aligned} \quad (8)$$

where B is an $M \times M$ matrix with elements $b_{mn} = K_{mn} + 2$, $K_{mn} = K_m K_n + K_n K_m$, c is an M dimensional vector with $c_n = K_n(2y - 1) + 2$ as its elements, $d_{t,t} > 0$ is the optimal step length in the t^{th} iteration [9]. Theorem 1 (a) ensures that such iteration will converge to the minimum. Or, the optimization process (8) will terminate at the boundary of region $d^2 = 0$ whenever any of the component of d_{t+1}^2 is less than or equal to 0.

The analytic approach finds out d^2 through solving the system of equations $dF(\cdot, d) = 0$. Specifically, it is easy to derive from $dF(\cdot, d) = 0$ that

$$d_n \sum_{m=1}^M d_m^2 K_{mn} + 2 \sum_{m=1}^M d_m^2 = d_n c_n, \quad (9)$$

where d_n is an element of d . Therefore,

$$d_n = 0, \text{ or } \sum_{m=1}^M d_m^2 K_{mn} + 2 \sum_{m=1}^M d_m^2 = c_n. \quad (10)$$

Suppose $d_n = 0$ only if $n \in S_0 \subseteq S_a \subseteq \{1, 2, \dots, M\}$, and $\bar{S}_0 = S_a \setminus S_0$, where the cardinality of S_0 equals N , and $0 \leq N < M$. Then Eq. (10) can be expressed as

$$\left(\sum_{m \in \bar{S}_0} K_{mn} + 2 \right) d_m^2 = c_n, \quad (11)$$

where $n \in \bar{S}_0$. This is a system of linear equations w.r.t. d_m^2 's, and can be briefly expressed as

$$d_p^2 = c_p B_p^{-1}, \quad (12)$$

where B_p is a $(M - N) \times (M - N)$ matrix with elements $b_{mn} = K_{mn} + 2$, d_p^2 and c_p are two column vectors with d_m^2 's and c_n 's as elements, respectively, $m, n \in \bar{S}_0$. If B_p is not invertible, the generalized inverse B_p^+ will be calculated instead of B_p^{-1} . If any element of $c_p B_p^{-1}$ is less than 0, there does not exist any solution d_p^2 for Eq. (12) in such S_0 and \bar{S}_0 . Otherwise, according to Theorem 1, alternately evaluating Eqs. (7) and (12) will make $F(\cdot, d)$ converge to a point $(\cdot, s(d_p^2 - d_0))$, where d_p^2 is evaluated via Eq. (12) and initially $d_p^2 = (1/M, 1/M, \dots, 1/M)$, d_0 is the vector with d_n , $n \in S_0$, as its elements, $d_p^2 - d_0$ is the union set of all components of both d_p^2 and d_0 , $s(v)$ sorts the elements of v ascendingly according to their subscripts, and generates a new column vector d^2 . For example, suppose $S_0 = \{2\}$, $\bar{S}_0 = \{1, 3\}$, $d_p^2 = \{d_1^2, d_3^2\} = \{0.4, 0.6\}$, and $d_0 = \{d_2\} = \{0\}$. Then $d^2 = s(v) = s(d_p^2 - d_0) = \{d_1^2, d_2, d_3^2\} = \{0.4, 0, 0.6\}$. In practice, a satisfactory convergence (\cdot, d^2) can be achieved after a couple of iterations of Eqs. (7) and (12).

In fact, N can take any value from 0 to $M - 1$, and given N , there are C_M^N different S_0 's. To distinguish these S_0 's one another, we introduce $S_0^{N,c}$ to represent the S_0 whose

cardinality is N and whose elements are the c^{th} combination of C_M^N ones. The C_M^N combinations consist of all possible N elements of S_a . It is clear that $c = 1, 2, \dots, C_M^N$. Through evaluating the convergent point by means of the above process for every $S_0^{N,c}$, at most $\sum_{N=0}^{M-1} C_M^N$ minimizers of $F(\cdot, d)$ will be found. The minimizer that makes $F(\cdot, d)$ minimal among the minimizers will be accepted as the optimal solution of Eq. (10).

It is interesting to examine the optimization of $F(\cdot, d)$ in the case that only two complementary features are included, i.e. $M = 2$, given \cdot . Because $(d_1^2 + d_2^2 - 1)^2$ is one of the three items to optimize in Eq. (5), it is advantage for the optimization that $d_1^2 = 1$ if $d_2^2 = 0$, or vice versa. Therefore, there are three combinations of the values of d_1^2 and d_2^2 , i.e., $(1, 0)$, $(0, 1)$, and (v_1, v_2) , where (v_1, v_2) is obtained through solving Eq. (12) with $N = 0$. This means that the optimal solution of Eq. (10) will be selected from three candidates: only employing one of two complementary features, and using the linear combination of both features. And the final solution that minimizes $F(\cdot, d)$ will be among the three candidates of d^2 , and will be applied to the construction of object appearance model in the current frame.

It is clear that the analytic approach to solving d is preferred if M is small enough. In our current implementation, we found that the analytic approach is more efficient than the iteration one because $M = 2$. The former spent about half of the time the later spent. Of course, the iteration approach may be more efficient than the analytic one if M is large. And that how many kernels will be applied depends on the tradeoff between the performance and computational burden.

3.2. Fast Evaluation in Training

All correlation filter based tracking algorithms [4, 15, 17, 7] consider the training samples, x_i , to be generated through cyclically shifting a base sample. Therefore, the optimization of $F(\cdot, d)$ described in the preceding section can be speeded up by means of fast Fourier transform (FFT).

Because kernel matrices K_m s are circulant [15], and the inverses, products, and sums of circulant matrices are still circulant [12], K_{mn} is circulant. Denote k_{mn} to be the first row of K_{mn} , and k_m the first column of K_m . It is clear that $k_m = k_m$ because K_m is symmetric. The evaluation of Eq. (7) can be accelerated by

$$= F_1^{-1} \frac{F_1(y)}{F_1 \sum_{m=1}^M d_m^2 k_m + 1}, \quad (13)$$

and that of Eq. (12) can be done by

$$k_{mn} = F_1^{-1} (F_1(k_n) - F_1(k_m)) + F_1^{-1} (F_1(k_m) - F_1(k_n)), \quad (14a)$$

$$b_{mn} = F_1^{-1}(F_1(k_{mn}) - F_1(\cdot)) + 2, \quad (14b)$$

$$c_m = (-y) F_1^{-1}(F_1(k_m) - F_1(\cdot)) - 1. \quad (14c)$$

It is clear that the linear combination of circulant matrices is still circulant. Therefore, $\sum_{m=1}^M d_m^2 K_m$ is circulant, and its first row is $\sum_{m=1}^M d_m^2 k_m$. The evaluation of $F(\cdot, d)$ can then be accelerated as follows.

$$F_2(\cdot, d) = \frac{1}{2} y - F_1^{-1} \left(F_1 \left(\sum_{m=1}^M d_m^2 k_m \right) - F_1(\cdot) \right) + \frac{1}{2} F_1^{-1} \left(F_1 \left(\sum_{m=1}^M d_m^2 k_m \right) - F_1(\cdot) \right). \quad (15)$$

It is noted that only the first two items of $F(\cdot, d)$ are necessary for the evaluation of $F(\cdot, d)$, as its third item is equivalent to a constraint.

3.3. Fast Detection

3.3.1 Determine Central Location

According to Eq. (6), the multiple kernel correlation filter evaluates the responses of all test samples $z_j = P^j z$, $j = 0, 1, \dots, l-1$, in the current frame l as

$$y^j(z) = \sum_{m=1}^M d_m^2 \sum_{i=0}^{l-1} i k_m(z_j, \bar{x}_i), \quad (16)$$

where z is the base test sample, P is the permutation matrix [15, 16]. $\bar{x}_i = P^i x$, $i = 0, 1, \dots, l-1$, \bar{x} is the weighted average of located samples in several earlier frames (i.e., \bar{x}^{new} in Sec. 3.3.4). Because $k_m(\cdot)$, $m = 1, 2, \dots, M$, is permutation-matrix-invariant, the response map, $y(z)$, of all virtual samples generated by z can be evaluated as

$$y(z) = (y^0(z), \dots, y^{l-1}(z)) = \sum_{m=1}^M d_m^2 C(\bar{k}_m), \quad (17)$$

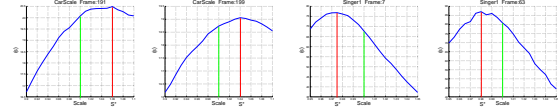
where $\bar{k}_m = (\bar{k}_{m,0}, \dots, \bar{k}_{m,l-1})$, $\bar{k}_{m,i} = k_m(z, P^i \bar{x})$, and $C(\bar{k}_m)$ is the circulant matrix with \bar{k}_m as its first row. Therefore, the response map can be accelerated as follows.

$$y(z) = \sum_{m=1}^M d_m^2 F_1^{-1} \left(F_1(\bar{k}_m) - F_1(\cdot) \right). \quad (18)$$

The element of $y(z)$ which takes the maximal value is accepted as the optimal object location in l .

3.3.2 Identify Optimal Scale

Suppose that $(l)_{\max}(y(z))$ is the center location of the test patch whose response is maximal among all test patches



(a)



(b)

Figure 2. (a) $y(s)$ usually has a dominant mode, which is assumed to be the optimal object scale s in current frame. 0.618 method [24] is employed to search for s . The object scale s of the last frame is supposed to be 1. (b) $y(s)$ usually also has a dominant mode. Its curve around s , however, is sometime a little bit too flat in comparison with (a)'s. See Sec. 3.3.2 for details. The figure is best viewed with high resolution display.

in l , $D(\cdot, s)$ is an l 's patch, whose center is $(l)_{\max}(y(z))$, with scale s , s_x is the scale of patch x , and $R(D, s)$ denotes the image patch D re-sampled by s . Let

$$(s) = \text{PSR}(y(R(D((l), s), s_x))), \quad (19)$$

where PSR is the peak to sidelobe ratio [4] in the response map $y(\cdot)$. Intuitively, the evaluation of (s) consists of four steps. The first is to extract a patch $D((l), s)$, and the second to re-sample $D((l), s)$ by scale s_x . The re-sampled $D((l), s)$ will be used as base sample then. The third is to generate the response map $y(\cdot)$ by using the base sample and Eq. (17), and the last to evaluate the PSR of $y(\cdot)$. It is noted that s_x is the template scale in the last frame. We experimentally observed that (s) usually possesses a dominant mode. Fig. 2 (a) shows several (s) 's generated in our experiments on sequences *CarScale* and *Singer1*. Therefore, it is reasonable to assume that the correct scale, s , of the object will maximize PSR, i.e., $s = \arg \max_s (s)$. Then the optimization technique is employed to seek out s efficiently in our tracker.

In fact, the optimal solution, s , of (s) can easily be sought out through exact line search because (s) is a univariate function. In our current implementation, we employed the golden section method, also called 0.618 method [24], to optimize (s) . This is because the 0.618 method only evaluates (s) once per iteration. Specifically, suppose that the aspect ratio of the target object is constant in the whole sequence, the object scale is 1 and the length and width of the located bounding box are l_0 and w_0 , respectively, in the last frame, and the initial interval of s for searching is set to be $s = [0.9, 1.1]$. Then, $\|s\| = 0.2$, where $\|s\|$ is the length of interval s , the maximum of length and width of the maximal candidate bounding boxes are $1.1 \cdot \max(l_0, w_0)$, the minimum of those are $0.9 \cdot \max(l_0, w_0)$, and their difference is $0.2 \cdot \max(l_0, w_0) =$

$\|s\| \max(l_o, w_o)$. As iteration proceeds, $\|s\|$ gets smaller and smaller until $\|s\| \max(l_o, w_o) < 1$ pixel. When the iteration stops, s is obtained and returned.

In each iteration, the length of searching interval is narrowed to 0.618 times the one in the last round. In other words, the computational complexity of 0.618 method is $O(\log(\max(l_o, w_o)))$. In comparison, the traditional exhaustive method has to search each layer of pyramid, and requires $O(\max(l_o, w_o))$ to find out the optimal scale.

Another possible choice is to use $f(\cdot)$ of Eq. (6), instead of PSR, to find out the optimal object scale. That is,

$$o(s) = y^0(R(D(I, s), s_{\bar{x}})). \quad (20)$$

We observed experimentally that, similar to (s) , $o(s)$ also usually possesses a dominant mode. Fig. 2 (b) shows several $o(s)$'s generated in our experiments in the same sequences and frames as used in Fig. 2 (a). Nevertheless, $o(s)$ may become a little bit too flat around the optimal scale, therefore, is not as robust as (s) . In fact, the performance will decrease at least 1% on the visual tracking benchmark [36] if Eq. (20), rather than Eq. (19), is utilized. Consequently, Eq. (19) is adopted to identify the optimal scale in our tracker.

3.3.3 Fast Feature Scaling

Refer to Eq. (19) again; $R(D(I, s), s_{\bar{x}})$ will re-sample $D(I, s)$ by scale $s_{\bar{x}}$. In order to accelerate the evaluation of $R(D(I, s), s_{\bar{x}})$, rather than rescaling image patch $D(I, s)$ and then extracting its feature in every iteration, we applied Dollár's power law [8] to rescale the feature directly. It has been pointed out that the processing time can be reduced several times in this way [8]. Specifically, if X_m is the m^{th} feature of image patch X , then Dollár's power law states that

$$\begin{aligned} R(D(I, s), s_{\bar{x}})_m \\ R(D(I, s)_m, s_{\bar{x}}) \cdot \frac{s_{\bar{x}}^{-m}}{s}, \end{aligned} \quad (21)$$

where m is a feature related constant, $R(\cdot)_m$ and $D(\cdot)_m$ are the m^{th} features of patches $R(\cdot)$ and $D(\cdot)$, respectively.

To fully take this advantage to speed up searching for the optimal scale of the object, we extract in practice the features of an image only once. And all the feature scaling operations discussed in Sec. 3.3.2 will be performed on feature channels.

3.3.4 Updating Filter

In our tracker, the object appearance model is (\cdot, d^2, \bar{x}) . We adopt the following formulation to update \bar{x} .

$$\begin{aligned} \bar{x}_m^{\text{new}} = & (1 - m)R(\bar{x}_m, s) \cdot \frac{s}{s_{\bar{x}}}^{-m} \\ & + mD(I, s)_m, \end{aligned} \quad (22)$$

where m is the learning rate, s is the optimal object scale in the current frame (Sec. 3.3.2), $m = 1, 2, \dots, M$. It is clear that \bar{x}_m^{new} is the weighted sum of historical templates and the m^{th} channel feature of object bounding box in the current frame, where the historical template has to be fast scaled to the current scale s by means of Dollár's power law [8]. $\bar{x}^{\text{new}} = (\bar{x}_1^{\text{new}}, \dots, \bar{x}_M^{\text{new}})$. \bar{x}^{new} and $(d^{\text{new}})^2$ are evaluated by using \bar{x}^{new} and Eqs. (12), (13), (14), and (15). It is noted that $k_m(x, x)$ will be evaluation by only using x_m and x_m .

The whole process of our tracker is summarized in Alg.1.

4. Implementation Details

Each of color and HOG features uses a kernel, *i.e.*, $M = 2$. The color scheme proposed in [7] is adopted in our tracker. To enhance the robustness against object deformation and speed up tracking, the HOG feature has only six gradient orientations and the cell size is 4×4 . Gaussian kernel is used for both features with $\text{color} = 0.4$ and $\text{HOG} = 0.5$, which ensures that all K_m s are positive definite [25]. According to Dollár *et al.* [8], the resizing coefficients m in Eq.(21) are set as follows: $\text{color} = 0$ and $\text{HOG} = 0.078$.

In order to reduce high-frequency noise in the frequency domain resulted from the large discontinuity between opposite edges of a cyclic-extended image patch, all feature patches (*e.g.* x_m and $D(\max(y(z)), s)_m$) are banded with a sine window for a sine window can reduce values near the borders to zero, and eliminate discontinuities.

In spite of the high efficiency of our continuous scale estimation, introducing scale processing inevitably increases the time complexity of our tracker. It is generally observed that, in most cases, the variation of target scale is much slower than that of its location. Therefore, it is superfluous to execute scale estimation as frequently as location. In our current experiments, the ratio of the number of scale estimations to that of locations is 0.5.

5. Experimental Results

Our tracker MKCF was implemented in MATLAB. The experiments were performed on a PC with Intel Core i5 3.20GHz CPU and 16GB RAM. We compared our MKCF to other 4 state-of-the-art trackers, Struck [13], DSST [6], CN2 [7], and KCF [16] with the visual tracking benchmark [36] which includes 50 image sequences. All parameter values of MKCF were kept consistent across all experimental comparisons. It is noted that DSST is the winner of VOT2014 challenge.¹ The mean fps of MKCF over the 50 sequences is about 15.

The performance of our tracker were quantitatively evaluated with popular criteria used in [1, 15, 36, 7, 6, 16], *i.e.*,

¹www.votchallenge.net/vot2014/results.html

Algorithm 1 Multi-kernel Correlation Filter Based Tracker

- **Input:** Frame t , $t = 0, 1, 2, \dots$, initial object patch x_0 of $l_1 \times l_2 \times c$, where c is the number of feature channels, $l_1 \times l_2$ Gaussian-shaped regression target y .
- **Output:** optimal locations $l_1, l_2 \dots$ in subsequent frames.
- **Initialization:** $t = 0$, $\bar{x} = x_t$.
- **Training** (Secs. 3.1 and 3.2):
 1. Generate virtual training set based on \bar{x} , $v_F =$.
 2. for $N = 0$ to $M - 1$
 - for $c = 1$ to C_M^N
 - Calculate (\cdot, d^2) by using \bar{x} and Eqs. (12), (13), (14), and (15).
 - If $v_F > F(\cdot, d^2)$
 - $v_F = F(\cdot, d^2)$, $(\cdot_o, d_o^2) = (\cdot, d^2)$.
- **Location** (Sec. 3.3):
 3. Determine the object location by using (\cdot_o, d_o^2) and Eq. (18) in frame $t + 1$.
 4. Determine the object scale by using Eqs. (19) and (21).
 5. Update \bar{x} by using Eqs. (22).
 6. $t = t + 1$, go to 1.

center error, distance precision, precision plot, overlap ratio, overlap precision, and success plot. Center error is calculated as the average Euclidean distance between the centers of located objects and their ground truths in a sequence. Distance precision is the percentage of frames where the objects are located within the center errors in 0 to t_c pixels, where $t_c = 20$, and the precision plot is simply a curve of the distance precisions with t_c changing from 0 to 50 pixels. Overlap ratio is defined as the average ratio of intersection and union of the estimated bounding box and ground truth in a sequence, overlap precision as the percentage of frames where the overlap ratio exceeds t_o in a sequence, where $t_o = 0.5$. And the success plot is simply a curve of overlap precisions with t_o changing from 0 to 1.

5.1. Fast Feature Scaling vs. Traditional One

In our approach, we extract the feature from a patch only once in a frame, and then scaling the patch's feature is approximately implemented by directly scaling the extracted feature channels. Such a way will significantly reduce the computational cost in extracting features [8]. While saving the processing time, it is also desirable not to lose much in tracking accuracy. Fig. 3 includes the average precision and success plots of the two versions, MKCF and MKCFN, of our tracker over the 28 sequences annotated with scale variation [36]. MKCF is just described in Alg. 1, while

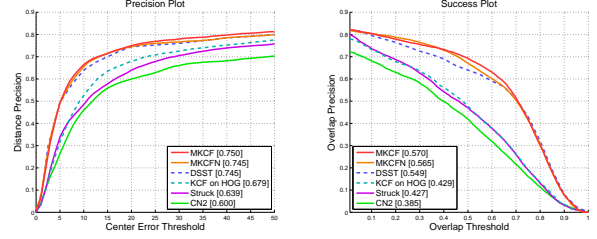


Figure 3. The average success plot of two versions, MKCF and MKCFN, of our tracker, DSST [6], KCF [16], Struck [13], and CN2 [7] over 28 sequences annotated with scale variation [36]. See Secs. 5.1 and 5.2 for details. The area under curve (AUC) of three trackers are reported in the legend.

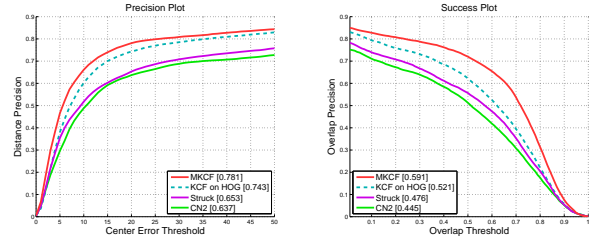


Figure 4. The average precision and success plots of our MKCF, Struck [13], KCF [16], and CN2 [7] over 50 sequences [36]. The mean distance precision scores and AUCs of each tracker are also reported in the legend. The figure is best viewed in color.

MKCFN is a variant of MKCF of which the fast Dollár's law [8] is replaced by a traditional one. The traditional scheme scales the patch's feature by first scaling original patches and then calculating their features in the scaled original channels. Except for the scaling schemes pointed out above, MKCF and MKCFN are exactly the same in terms of implementation and parameter setting. Surprisingly, Fig. 3 indicates that using the fast Dollár's law scheme generates higher distance and overlap precisions than using the traditional one, and the AUC of MKCF is also larger than that of MKCFN. In contrary to the conclusion in object detection domain that approximately scaling features with Dollár's law will slightly reduce the detection accuracy [8], approximate feature extraction scheme indeed improves the robustness and accuracy of tracking. The experimental comparison reveals the *contrary affects* of the approximate scheme on object detection and visual tracking.

5.2. Comparison to State-of-the-art Trackers

Fig. 4 shows the average precision and success plots of our tracker MKCF, Struck [13], KCF [16], and CN2 [7] over the 50 sequences. The mean distance precisions and AUCs are also included in it. It is seen that our tracker outperforms other three ones.

We also compared our MKCF to Struck, KCF (on HOG), and CN2 over the 50 sequences with respect to the 11 anno-

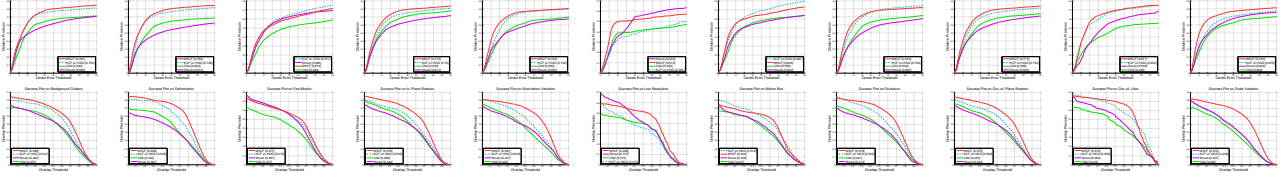


Figure 5. The average precision plots of our MKCF and Struck [13], KCF [16], and CN2 [7] on 11 attributes [36]. The mean distance precision scores and AUCs of each tracker are also reported in legends. The figure is best viewed with high resolution display.

	Boy	Card	CarScale	Couple	Crossing	David	Dog1	Doll	Dudek	Fleetface	Freeman1	Freeman3	Freeman4	Girl
MKCF	97.2	100	89.7	45	100	99.6	100	100	98.8	63.6	61	38.7	31.1	83
DSST [6]	100	100	84.5	10.7	100	100	100	99.6	98.1	66.5	36.8	31.3	41.7	24.2
KCF [16]	99.2	36.4	44.4	24.3	95	62.2	65.1	55.2	97.6	66.9	16.3	29.1	18.4	74.2
CN2 [7]	95.3	27.6	44.8	10.7	96.7	58.8	65.3	72.8	96.1	58.7	15	33	17.3	46.4
Struck [13]	99.7	39.8	43.3	54.3	94.2	23.6	65.3	68.8	98.0	66.6	21.8	20.0	15.9	98.0
	Ironman	Lenning	Liquor	Matrix	mRolling	Shaking	Singer1	Skating1	Skiing	Soccer	Trellis	Walking	Walking2	Woman
MKCF	13.3	27.2	99.4	19	7.3	96.2	100	63.5	4.9	13.5	96.7	99.8	99.8	93.8
DSST [6]	13.3	26.9	40.9	18	6.7	100	100	54.8	4.9	52.8	96.8	99.8	100	93.3
KCF [16]	15.1	44.2	99	13	7.9	1.4	27.6	36.3	7.4	39.3	84	51.5	38	93.6
CN2 [7]	13.3	29.1	20.4	1	7.3	67.4	27.6	37.3	8.9	48.2	65.9	45.9	38.4	24.5
Struck [13]	4.8	64.1	40.6	12.0	15.9	16.7	29.9	37.0	3.7	15.6	78.4	56.6	43.4	93.5

Table 1. Overlap precision in percent on the 28 sequences annotated with scale variation. The best scores are shown in bold.

tated attributions [36]. Fig. 5 reports the average precision and success plots, the mean distance precisions, and AUCs. It is seen in Fig. 5 that our MKCF almost consistently outperforms other three trackers, except for the following three cases. 1) MKCF is inferior to KCF (on HOG) over the 12 sequences of motion blur, but superior to all other trackers if the center error threshold $t_c = 11$ or the overlap threshold $t_o = 0.44$. 2) Over the 17 sequences of fast motion, MKCF is inferior to KCF (on HOG) in the precision plot. 3) MKCF is inferior to Struck over 4 sequences of low resolution at distance precision of $t_c = 20$, but superior to all other trackers if $t_c = 17$. It has been accepted through experiments that the representation with HOG-based feature performs poorly over low resolution sequences [6, 16]. By combining HOG and its complementary feature color through multi-kernel technique, our MKCF improves its performance on low resolution sequences by a large margin.

Table 1 lists a per-sequence comparison of our MKCF to DSST [6], KCF, Struck, and CN2 in overlap precision. It is easy to check out that MKCF outperforms all other trackers on 14 sequences, is superior to DSST on 13 out of the 28 sequences, and has a similar accuracy to DSST on 7 sequences. Fig. 3 shows the average distance and success plots of MKCF and DSST over 28 sequences annotated with scale variation. The mean distance precision and AUC are also reported. We do not show the performance curves of other state-of-the-art trackers in Fig. 3, because DSST outperforms them consistently on the 28 sequences [6]. It is clear from Table 1 and Fig. 3 that MKCF outperforms DSST on the benchmark.

It is interesting to notice that the difference of performance between our MKCF and DSST is made by the difference of their curves at around 0.5. This means that MKCF is more robust than DSST in difficult frames where the lo-

cations and scales of object are easily bias away from the correct ones a little bit more.

6. Conclusion

A novel tracking algorithm, MKCF, has been presented in this paper. To construct multi-feature appearance models, we proposed an approach to fusing features of different types by means of multiple kernel learning. Instead of building traditional image pyramids in advance, we employed the optimization technique to search for the correct object scale in the continuous scale space efficiently. The whole tracking algorithm is accelerated by FFT. Extensive experiments on the benchmark have shown that our algorithm outperforms the state-of-the-art algorithms.

A. Proof of Theorem 1

(a) $d^2 F(\cdot, d) = \frac{1}{2}(d^2 B - c)$, $\frac{1}{2} d^2 F(\cdot, d) = \frac{1}{2} B$, where B is an $M \times M$ matrix with elements $b_{mn} = K_{mn} + 2$, $K_{mn} = K_m K_n + K_n K_m$, c is an $M \times 1$ vector with $c_n = K_n(2y - \cdot) + 2$ as its elements. All K_m s are circulant (Theorem 1, [15]) and positive definite, $K_m = U \frac{1}{2} U^T$ [12], $K_m K_n = U \frac{1}{2} U^T$, where $\frac{1}{2} = \frac{1}{2} \frac{1}{2}$, K_m s are positive definite, $b_{mn} > 0$ as 0 , $d^2, d^2 B d^2 > 0$ if $d^2 = 0$. Therefore, $F(\cdot, d)$ is convex w.r.t. d^2 .

(b) $\frac{1}{2} F(\cdot, d) = K(K + I)$. K is positive definite and circulant, K is positive definite and circulant. $K = U \frac{1}{2} U^T$ [12], $K + I = U(\frac{1}{2} + I)U^T$, $K + I$ is positive definite as > 0 , $K(K + I) = U(\frac{1}{2} + \frac{1}{2})U^T$. $K(K + I)$ is positive definite as > 0 . Therefore, $F(\cdot, d)$ is convex w.r.t. d .

References

- [1] B. Babenko, M.-H. Yang, and S. Belongie. Robust object tracking with online multiple instance learning. *IEEE T-PAMI*, 33(No.8):1619–1632, 2011. **2, 6**
- [2] S. Baker and I. Matthews. Lucas-kanade 20 years on: A unifying framework. *IJCV*, Vol.56(No.3):221–255, 2004. **2**
- [3] V. Boddeti, T. Kanade, and B. Kumar. Correlation filters for object alignment. In *CVPR*, 2013. **2**
- [4] D. Bolme, R. Beveridge, B. Draper, and Y. Lui. Visual object tracking using adaptive correlation filters. In *CVPR*, 2010. **1, 2, 4, 5**
- [5] C. Cortes, M. Mohri, and A. Rostamizadeh. ℓ_2 regularization for learning kernels. In *UAI*, 2009. **2**
- [6] M. Danelljan, G. Hager, F. Khan, and M. Felsberg. Accurate scale estimation for robust visual tracking. In *BMVC*, 2014. **1, 2, 6, 7, 8**
- [7] M. Danelljan, F. Khan, M. Felsberg, and J. van de Weijer. Adaptive color attributes for real-time visual tracking. In *CVPR*, 2014. **1, 2, 4, 6, 7, 8**
- [8] P. Dollár, R. Appel, S. Belongie, and P. Perona. Fast feature pyramids for object detection. *IEEE T-PAMI*, Vol.36:1532–1545, 2014. **2, 6, 7**
- [9] R. Fletcher. *Practical Methods of Optimization*. John Wiley & Sons, 2nd edition, 1987. **4**
- [10] H. Galoogahi, T. Sim, and S. Lucey. Multi-channel correlation filters. In *ICCV*, 2013. **2**
- [11] P. Gehler and S. Nowozin. On feature combination for multiclass object classification. In *ICCV*, 2009. **1, 2**
- [12] R. Gray. *Toeplitz and Circulant Matrices: A review*. Now Publishers Inc., 2006. **4, 8**
- [13] S. Hare, A. Saffari, and P. Torr. Struck: Structured output tracking with kernels. In *ICCV*, 2011. **1, 6, 7, 8**
- [14] J. Henriques, J. Carreira, R. Caseiro, and J. Batista. Beyond hard negative mining: Efficient detector learning via block-circulant decomposition. In *ICCV*, 2013. **2**
- [15] J. Henriques, R. Caseiro, P. Martins, and J. Batista. Exploiting the circulant structure of tracking-by-detection with kernels. In *ECCV*, 2012. **1, 2, 4, 5, 6, 8**
- [16] J. Henriques, R. Caseiro, P. Martins, and J. Batista. High-speed tracking with kernelized correlation filters. *IEEE T-PAMI*, Vol.37(No.3):pp.583–596, 2015. **1, 2, 5, 6, 7, 8**
- [17] J. Henriques, P. Martins, R. Caseiro, and J. Batista. Fast training of pose detectors in the fourier domain. In *NIPS*, 2014. **2, 4**
- [18] Z. Hong, C. Wang, X. Mei, D. Prokhorov, and D. Tao. Tracking using multilevel quantizations. In *ECCV*, 2014. **1**
- [19] B. Kumar, A. Mahalanobis, and R. Juday. *Correlation Pattern Recognition*. Cambridge University Press, 2005. **2**
- [20] J. Kwon, J. Roh, K.-M. Lee, and L. Van Gool. Robust visual tracking with double bounding box model. In *ECCV*, 2014. **1**
- [21] X. Lan, A. Ma, and P. Yuen. Multi-cue visual tracking using robust feature-level fusion based on joint sparse representation. In *CVPR*, 2014. **1**
- [22] G. Lanckriet, T. De Bie, N. Cristianini, M. Jordan, and W. Noble. A statistical framework for genomic data fusion. *Bioinformatics*, 20:2626–2635, 2004. **3**
- [23] D. Lee, J.-Y. Sim, and C.-S. Kim. Visual tracking using pertinent patch selection and masking. In *CVPR*, 2014. **1**
- [24] D. Luenberger. *Introduction to Linear and Nonlinear Programming*. Addison-Wesley Reading, MA, 1973. **5**
- [25] C. Micchelle. Interpolation of scattered data: Distance matrices and conditionally positive definite functions. *Constructive Appr.*, Vol.2:pp.11–22, 1986. **6**
- [26] H. Nam, S. Hong, and B. Han. Online graph-based tracking. In *ECCV*, 2014. **1**
- [27] S. Oron, A. Bar-Hillel, and S. Avidan. Extended lucas-kanade tracking. In *ECCV*, 2014. **2**
- [28] A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. SimpleMKL. *JMLR*, 9:2491–2521, 2008. **2**
- [29] R. Rifkin, G. Yeo, and T. Poggio. Regularized least-squares classification. *Nato Science Series Sub Series III: Computer and Systems Sciences*, pp.131–154.:2003, 190. **2, 3**
- [30] D. Ruderman. The statistics of natural images. *Network: Computation in Neural Systems*, 5:517–548, 1994. **2**
- [31] B. Schölkopf and A. Smola. *Learning with Kernels*. MIT press Cambridge, MA, 2002. **3**
- [32] M. Tang and X. Peng. Robust tracking with discriminative ranking lists. *IEEE T-IP*, Vol.21(No.7):3273–3281, 2012. **1**
- [33] M. Varma and D. Ray. Learning the discriminative power-invariance trade-off. In *ICCV*, 2007. **1, 2, 3**
- [34] A. Vedaldi, V. Gulshan, M. Varma, and A. Zisserman. Multiple kernels for object detection. In *ICCV*, 2009. **2**
- [35] Y. Wu, G. Blasch, G. Chen, L. Bai, and H. Ling. Multiple source data fusion via sparse representation for robust visual tracking. In *FUSION*, 2011. **1**
- [36] Y. Wu, J. Lim, and M.-H. Yang. Online object tracking - a benchmark. In *CVPR*, 2013. **1, 2, 6, 7, 8**
- [37] F. Yang, H. Lu, and M. Yang. Robust visual tracking via multiple kernel boosting with affinity constraints. *IEEE T-CSVT*, 24:242–254, 2014. **1**