

# Efficient High-Resolution Template Matching with Vector Quantized Nearest Neighbour Fields

Ankit Gupta<sup>a,\*</sup>, Ida-Maria Sintorn<sup>a,b</sup>

<sup>a</sup>*Department of Information Technology, Uppsala University, Uppsala, 752 37, Sweden*

<sup>b</sup>*Vironova AB, Stockholm, 113 30, Sweden*

---

## Abstract

Template matching is a fundamental problem in computer vision and has applications in various fields, such as object detection, image registration, and object tracking. The current state-of-the-art methods rely on nearest-neighbour (NN) matching in which the query feature space is converted to NN space by representing each query pixel with its NN in the template pixels. The NN-based methods have been shown to perform better in occlusions, changes in appearance, illumination variations, and non-rigid transformations. However, NN matching scales poorly with high-resolution data and high feature dimensions. In this work, we present an NN-based template-matching method which efficiently reduces the NN computations and introduces filtering in the NN fields to consider deformations. A vector quantization step first represents the template with  $k$  features, and then filtering compares the template and query distributions over the  $k$  features. We show that state-of-the-art performance was achieved in low-resolution data, and our method outperforms previous methods at higher resolution showing the robustness and scalability of the approach.

*Keywords:* template matching, vector quantized nearest neighbour field (VQ-NNF), object detection, high-resolution template matching.

---



---

\*Corresponding Author.

E-mail addresses: ankit.gupta@it.uu.se, ida.sintorn@it.uu.se

## 1. Introduction

Template matching refers to locating a small template image,  $T$  in a larger image  $I$ . It is a fundamental problem in computer vision. It has applications in various fields such as object detection [33, 26, 31, 37], object tracking [40, 43], document information identification [34], counterfeit detection [28] and image registration [46, 25]. It also has a central role in deep learning due to the increasing demand for annotated data. Template matching has, for example, been used as a tool for human-in-the-loop data annotation frameworks [10, 37, 15, 24, 11] as it offers fast detection at a relatively low cost of labour and resources. Here, we refer to human-in-the-loop data annotation as object detection on an image set where one or multiple instances of different classes can be present. Template matching allows users to find similar objects quickly without expensive classifier training.

Traditional template matching approaches such as sum-of-squared distance (SSD), sum-of-absolute distance (SAD), and normalized cross-correlation (NCC) are very efficient. However, they consider every pixel pair in the sliding subwindow of query image  $I$  and  $T$ , making them vulnerable to occlusion and transformations. More recently, nearest neighbour field (NNF) based matching approaches [29, 36, 35, 22] have been suggested and shown to overcome these shortcomings, making them state-of-the-art in the field. NNFs constitute a general and non-parametric framework for generating correspondences between the sub-regions of images and have been successfully used for motion-tracking in videos [51, 4], optical flow algorithms [7], and structural image editing [3]. NNF-based template matching approaches rely only on the subset of “good” matches between the template and the query subwindow. This makes them more robust against complex non-rigid transformations, occlusions, and background clutter. They achieve this by matching between two point sets, namely, the template point set and the query point set. A similarity measure is defined between the sets based on the matching statistics, for example, bi-directional matches [9] or unique matches [36].

Due to the powerful representation capability, features from pre-trained deep learning models are attractive for template matching. However, the approximate nearest neighbour search methods used in the algorithms mentioned above scale poorly with both the feature dimensions as seen in [2, 22, 23] and the template sizes (more feature points). Both speed and recall of the approximate nearest neighbour methods are affected negatively by increasing either the number of data points or the feature dimensions. Emphasizing

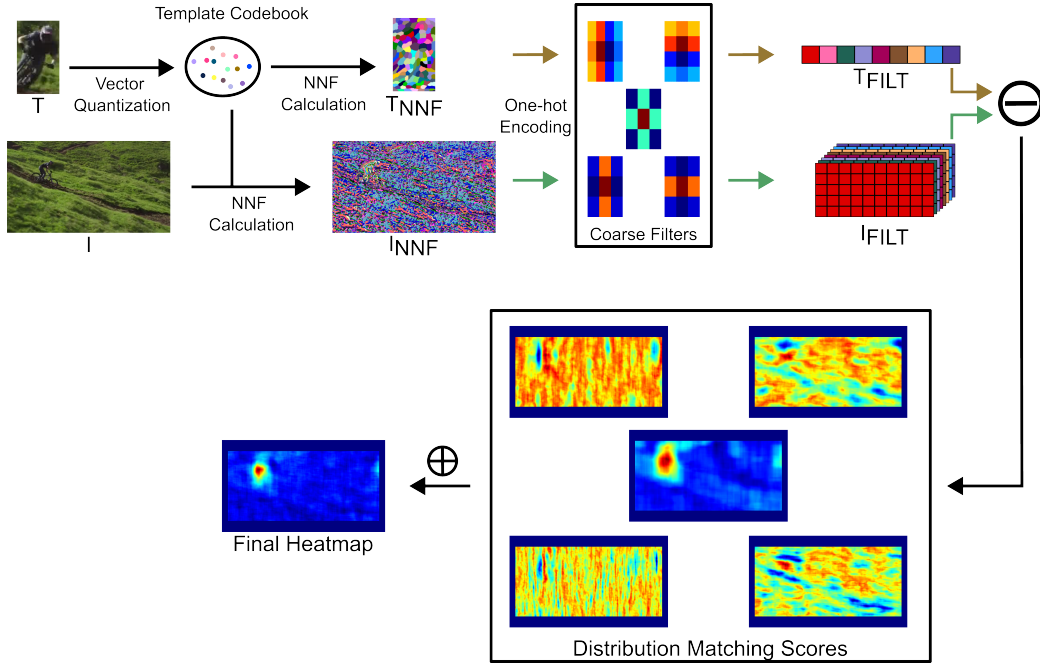


Figure 1: Overview of the proposed template matching method. First, a  $k$ -sized vector-quantized codebook is constructed using the template points. The codebook is used to generate the NNF label image for the template and the query image. Multiple template representations are generated using coarse filters and compared with the filtered query distribution. The matching scores from different filters are combined to get the final heatmap.

speed, the methods use PCA to reduce the feature dimensions, which reduces the features’ representation capacity.

An inadvertent effect of using too many points in the NNF representation is the difficulty in considering the deformation implied in the NNF in the similarity score. To model the deformation in the NNF, previous methods rely on a deformation measure for each pixel in the query subwindow, penalising large relative distances between the template point location and its NN match location. This relies on the strong assumption of the similar orientation of the template and query subwindow, which makes these methods incapable of handling significant rotational or deformation changes. It also requires more computations as the relative distance has to be calculated for each point in the query sub-window.

NNF-based approaches represent the points in the query image with the

nearest neighbour of the points in the template image. This formulation reduces the pixel representation space in relation to the template. For example, in an 8-bit RGB image, the possible pixel values can be from a set size of  $|255|^3$ , but representing the image with the NNF of a template size  $|w \times h|$  will reduce the representation set to  $wh$  thus greatly lowering the set of possible values. It can be viewed as a form of vector quantization [17] where the codebook is defined by the template pixels, and the query image pixels are then quantized with respect to the codebook. However, as mentioned above, the codebook is still too large for practical implementations and contains redundant information, as the representation of nearby pixels would be very similar.

In this work, we introduce a fast and flexible template-matching approach that reduces the NN calculations in previous approaches and better utilizes the NNF. We reduce the template codebook significantly to  $k$  points obtained by vector quantization of the template points. The points in the codebook represent the  $k$  major patterns in the template. This changes the computational complexity in the matching step to depend on only  $k$  points instead of the number of template pixels, greatly reducing the NN computation complexity and NNF creation time. The difference in the pixel distribution in the template and query subwindow among these  $k$  patterns is then used as the similarity measure in the matching. Since the number of NNF labels is small, simple coarse filters can be used to model the deformation in the NNF instead of considering pixel-wise distances. Our major contributions are:

- We present a simple and fast NNF-based template matching method that greatly reduces computational costs in the NNF creation while maintaining performance.
- We introduce filtering in the NNF space to model the deformation using coarse filters and show that state-of-the-art performance can be achieved with simple gaussian and haar filters.
- We show that our method’s quantitative performance and run-time scale better with the image resolution than previous approaches.

## 2. Related Work

Traditional methods in template matching, such as SSD, SAD, and NCC, work well in cases where only translational variance is present. They are very



sensitive to deformations and non-rigid transformations. Classical algorithms that use SSD and SAD as the dissimilarity measure are comprehensively reviewed in [30].

Different approaches have been proposed to model affine transformations between the template and query window [18, 38, 20, 48, 14]. More recently, [16] proposed an approach based on template feature co-occurrence matrix statistics for the matching algorithm reaching state-of-the-art performance on the BBS datasets. In [19], a method is proposed that provides a quadratic improvement in the search complexity while also being robust to partial occlusion by formulating the template matching as a consensus set maximization problem, i.e., finding the transformation where the maximum number of pixels between the template and query window are co-visible. In [45], an adaptive radial ring code histogram (ARRCH) image descriptor is used, which is robust against large-scale and rotation variations. The ARRCH descriptor is created by generating a histogram of the stable pixels at different concentric rings in the template. In [44], a superpixel region binary descriptor (SRBD) is suggested to construct a multilevel semantic fusion vector used in the matching. First, the template is divided into superpixels by the KD-SLIC [1] algorithm, whereafter, a region binary vector is constructed by describing the dominant superpixel orientation. The method is robust against large deformations, occlusions, noise and illumination variations; however, the computation complexity is  $O(360 \cdot K \cdot |I|)$  where  $K$  is the number of superpixels. This makes it challenging for real-time use. In [8], a quality-aware template matching approach is proposed, which can also be implemented as a trainable layer in a deep learning model.

Methods relying on NNF-based similarity measures have shown great promise due to their robustness against occlusion and deformations. In [9], the Best-Buddies Similarity (BBS) measure was introduced based on the properties of the nearest-neighbour (NN) matches between the features of the template and the query image. This bidirectional matching focuses only on the relevant corresponding features, thus making it robust against deformations and occlusion, outperforming previous state-of-the-art. This approach is slow and hard to use in practice as the computational complexity for BBS calculation is  $O(|I| \cdot |T| \cdot |I|)$  (where  $|I|$  denotes the size of the query image and  $|T|$  denotes the size of the template). In [36], inspired by the idea of matching objects for texture synthesis using the patch diversity [13], diversity similarity (DIS) and deformable diversity similarity (DDIS) measures were proposed that only rely on the NN matching in a single direction. This

approach is faster than BBS; however, with the computational complexity of  $O(|I| \cdot |T|)$ , it still becomes time-consuming for larger image and template sizes.

To reduce the computational complexity of the methods proposed above while retaining the performance, [35] proposed image-based unpopularity (IWU) and deformable image-based unpopularity (DIWU). They use the diversity within the query image NNF instead of the subwindow of NNF, reducing the complexity to  $O(I)$ . The authors used “unpopularity” in meaningfully the same way as “diversity” used in [36]. In [21], a majority neighbour similarity and annulus projection transformation were used to provide a fast template matching method which is also robust to different challenges. Recently, [22] proposed a global-aware diversity-based method, GAD, which combines the IWU and DIS scores to propose a parameter-free algorithm, also with  $O(I)$  complexity. An inverse NNF-based method was proposed in [23] where instead of searching the NNs of query points, the diversity of the NNs of template points is used. However, deformability is not considered, which limits the method’s usability. A scale-adaptive NNF-based method was proposed in [50]. It extends DIS through statistical analysis, which makes it robust to outliers. However, the method performs calculations over multiple scales, which is time-consuming. In [49], authors use bi-directional NN calculations to make the diversity similarity measure robust against scaling, rotation, and illumination changes at the expense of added computational complexity.

### 3. Method

We define the input to our template matching method as a template  $\tau$  of size  $w \times h$  and a query image  $I$  of size  $W \times H$ . The goal is to find the  $w \times h$  subwindow  $q$  in  $I$  most similar to  $\tau$ . To do this, a similarity score  $S_{q_i}$  is calculated for each query subwindow  $q_i$  using the sliding window procedure, and the subwindow with the highest score is our desired output.

Our template-matching approach consists of three main steps. An overview of the pipeline is shown in Fig. 1, and each step is described below.

#### 3.1. Template Vector Quantization

Our method uses K-means clustering of the template features to get the  $k$  quantized vectors for the codebook. This reduces the number of points from  $w \times h$  to  $k$  (where  $k \ll wh$ ) for the NN computations in the next step.

Furthermore, it reduces the redundant information present in the template representation. The  $k$  cluster centers,  $C = \{c_i\}_{i=1}^k$ , are then used for the new template representation for the NNF calculations. We next represent each pixel in the template with the cluster it belongs to and form a template label map,  $L_\tau$ , as shown in Fig. 1. The NNF for the query image  $I$  is calculated using the codebook vectors and represented similarly as  $L_I$ .

### 3.2. Feature Vector Construction

Since the template point set is now greatly reduced in size, the similarity measures used in previous state-of-the-art methods can't be used. They rely on the diversity of the matches in the query subwindow [36] or whole image [35] and depend on all the pixels in the template to provide good discrimination between a potential match and background in the query image. Instead, we represent the template and the query subwindow with a distribution histogram of the cluster labels and use the difference between the two as our similarity measure. As the number of clusters is small, this process can be done efficiently by first converting the label image into a one-hot encoded (binary) image with  $k$  channels and then using the integral image of each channel to reduce the operations per query subwindow. The integral image representation enables rapid summation over rectangular image subregions and can be evaluated in constant time for any rectangular size. The integral image  $ii$  at each location  $(x, y)$  of the  $k$  dimensional one-hot image  $\iota : \mathbb{Z}^2 \rightarrow \{0, 1\}^k$  is defined as:

$$\iota(x, y) = \sum_{x' \leq x, y' \leq y} \iota(x', y'). \quad (1)$$

In an integral image  $\iota$ , the sum of the values in a rectangular region  $A$  defined by a top-left point  $(x_1, y_1)$ , top-right point  $(x_2, y_1)$ , bottom-left point  $(x_1, y_2)$ , and bottom-right point  $(x_2, y_2)$  is calculated as:

$$sum(A) = \iota(x_2, y_2) + \iota(x_1, y_1) - \iota(x_1, y_2) - \iota(x_2, y_1) \quad (2)$$

Thus, a  $k$  dimension histogram vector can be calculated efficiently to represent the cluster labels for a window size of  $w \times h$  in the query image.

As evident by the limitations of the histogram-based methods, this formulation removes the orientation information of the template. Furthermore, it does not explicitly incorporate spatial information, meaning it does not give any spatial preferences to labels at different locations in the template.

To address these limitations, we employed coarse filters and compared the filter response of the template and the query subwindow to calculate the similarity score. Coarse filters operate on larger regions of an image, such as blocks or segments, rather than operating on individual pixels, combining information from these regions. Fig. 1 shows a few examples of coarse filters used in this work.

### 3.2.1. Coarse filters with dilated convolutions

As the weight in a coarse filter is the same for a set (region) of pixels, dilated convolution can be used to implement the filters and work on the integral images efficiently. Dilated convolutions are often used for multi-scale context aggregation in convolutional neural networks (CNNs) for semantic segmentation[47, 41, 5, 6]. Dilated convolutions work similarly to regular convolutions but allow for skipping pixels during convolution. The number of pixels to be skipped is referred to as the dilation rate and is used to increase the receptive field of the convolution while keeping the same computational complexity. Thus, coarse filters implemented using dilated convolutions and made to work on the integral images allow for fast computation of filter responses over a relatively large window size.

**Dilated Convolutions:** Let  $f : ([-h_f, h_f] \times [-w_f, w_f]) \cap \mathbb{Z}^2$  be a discrete filter of size  $r_h \times r_w$ , where  $r_w = 2w_f + 1$ ,  $r_h = 2h_f + 1$ . The filter response is then calculated using dilated convolution as:

$$(\iota_{*(d_x, d_y)} f)(x, y) = \frac{1}{(r_w \cdot d_x \cdot r_h \cdot d_y)} \sum_{m=-h_f}^{h_f} \sum_{k=-w_f}^{w_f} \iota(x+d_x \cdot k, y+d_y \cdot m) \cdot f(k, m), \quad (3)$$

where  $d_x$  and  $d_y$  are the dilation rates of the kernel in the x and y direction, respectively. The receptive field of the filter is  $(d_x \cdot (r_w - 1) + 1, d_y \cdot (r_h - 1) + 1)$  and hence can be modified using the dilation rate and the kernel size. For example, the receptive field of  $(w, h)$  can be achieved by a convolutional kernel of size  $(3, 3)$  by setting the dilation rate to  $((w - 1)/2, (h - 1)/2)$ . Thus the filter response over a larger area can be achieved while requiring only relatively few (9) operations for each query sub-window. This, however, comes at the cost of large bin size  $(w/3, h/3)$  for the label aggregation and, thus, lower deformation granularity. This means that the filtering operation can not capture label distribution shifts in the NNF within the bin size.

The granularity of the filter response for a fixed receptive field can be modulated in the following ways (as also illustrated in Fig. 2):

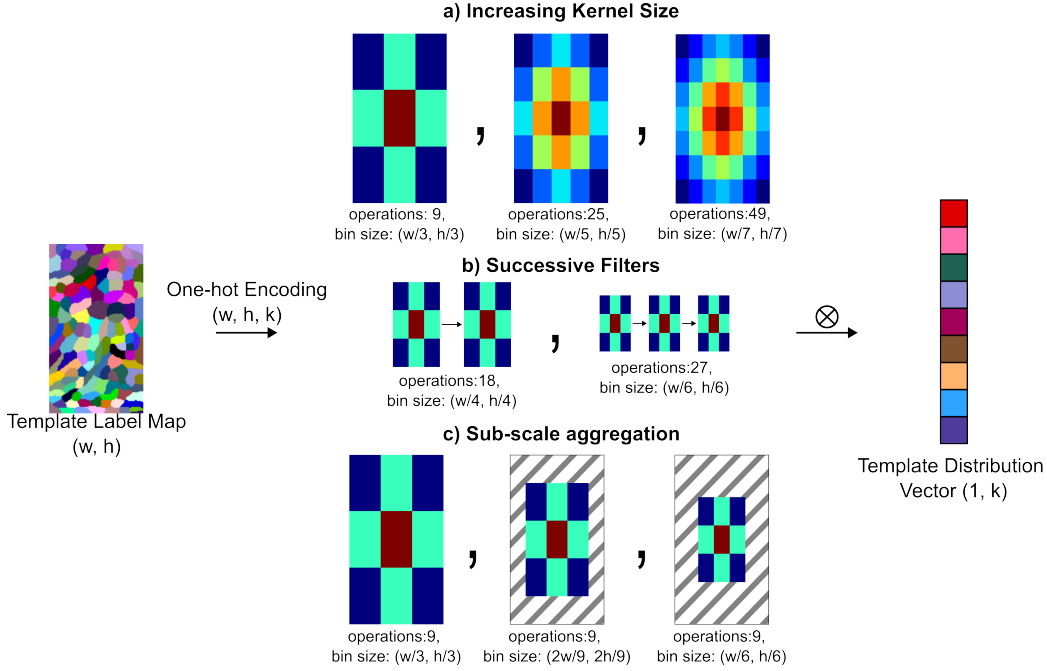


Figure 2: Different ways of modulating the bin size for label distribution aggregation by a) increasing kernel size, b) successive filtering and c) sub-scale aggregation. The grey-shaded region in c) shows the region not considered in the filtering.

1. Increasing the filter kernel size. Since the bin size is inversely proportional to the kernel size for a fixed receptive field, increasing the kernel size increases the granularity of the filtering. However, the computations required for each window increase quadratically with the kernel size, which might be undesirable.
2. Successive filtering. Applying the filters multiple times increases the granularity exponentially with the dilation rate. At the same time, the number of computations grows linearly.
3. Sub-scale aggregation. Reducing the receptive field results in finer filter bins and, thus, higher granularity. This enables variable deformation penalties for selective template sub-regions. For example, a higher deformation penalty can be enforced in a template subregion by reducing the receptive field without increasing the number of computations.

**Filter Modification for Integral Image:** To get the filter response over the integral image, the coarse filter must first be converted into a dilation

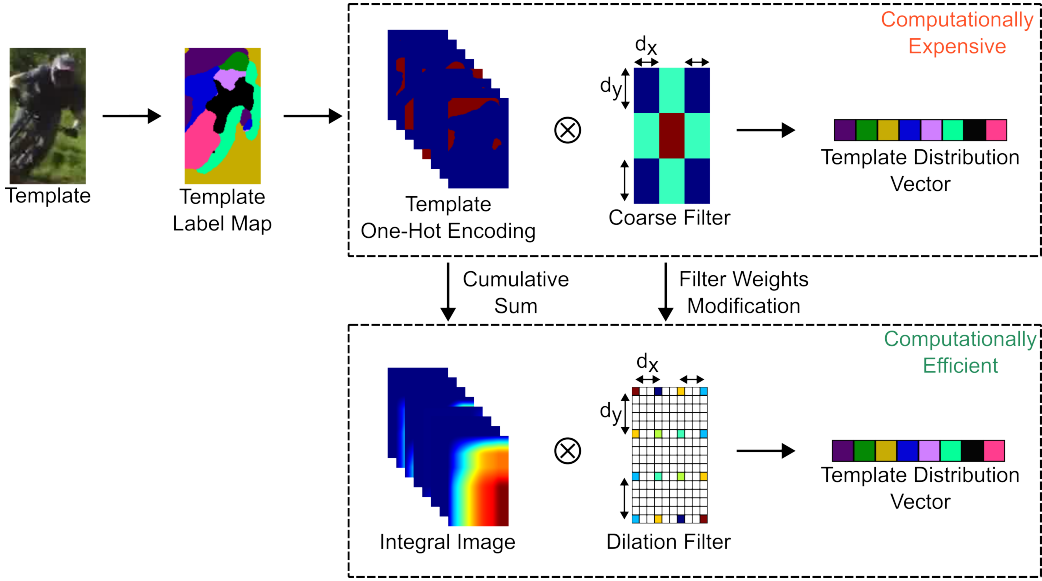


Figure 3: Process to efficiently calculate the coarse filter response over the label map. The one-hot encoded image is converted to an integral image, and the coarse filter is converted into a dilation filter to reduce the computational complexity. A small  $k = 8$  was chosen here for illustration purposes.

filter. Let  $f_i$  be a coarse filter of kernel size  $(r_w, r_h)$  with each weight repeated over the rectangular bins of size  $(d_x, d_y)$ . The dilation filter  $f_d$  would then have the kernel size of  $(r_w + 1, r_h + 1)$  and the filter weights of the dilation filter can be calculated using Algorithm 1. Thus, the number of calculations per window has been reduced from  $r_w \cdot d_x \cdot r_h \cdot d_y$  to  $(r_w + 1) \cdot (r_h + 1)$  making the filtering process efficient. Fig. 3 shows the process of calculating the filter response with this approach.

### 3.2.2. Template Representation

Multiple different filters can be used in the aforementioned ways to model spatial and orientation information of the template label distribution without much computational overhead. Here, we used the combination of coarse gaussian and rectangular haar-like filters to construct multiple template representation vectors from the response. A gaussian filter is the simplest way to impart spatial information as it assigns more weight to central regions and gradually decreases the weight as the distance from the center increases. Haar-like filters were introduced in [39] for rapid object detection using a

---

**Algorithm 1:** Modifying dilation filter weights for region aggregation on the integral image

---

**Result:** Dilation filter for integral image,  $f_d$  of size  $(r_w + 1, r_h + 1)$   
Original filter  $f_i$  of size  $(r_w, r_h)$ ;  
Integral Multiplier,  $m = [[1, -1], [-1, 1]]$ ;  
Initialize  $f_d$  with zeros;  
**for**  $i$  **in**  $r_w$  **do**  
    **for**  $j$  **in**  $r_h$  **do**  
         $f_d(i : i + 1, j : j + 1) += f_i(i, j) * m$   
    **end**  
**end**

---

boosted cascade of a large set of simple filters. Haar filters can be used to encode the orientation efficiently. This technique of combining Gaussian and Haar-like filter responses was chosen due to its ability to encode both spatial and orientation information with minimal computational overhead. More complex filters can also be used to generate richer template representations; however, we show that state-of-the-art (SOTA) performance can be achieved using these simple filters at different subscales (the c-option in Fig. 2). Since SOTA was achieved using the subscale approach, the main experiment did not consider alternative approaches of increasing the kernel size or successive filtering.

The template representation vectors obtained from the filter responses,  $R(\tau)$ , can be defined as a collection of filtered responses:

$$R(\tau) = \{L_{\tau,s}^f\}, \forall s \in S, f \in \{f_i\}_{i=1}^n$$

where,  $L_{\tau,s}^f = \iota_\tau *_{((w-1)/(s \cdot (r_w-1)), (h-1)/(s \cdot (r_h-1)))} f$

$S = \{1, 2, \dots, s\}$  defines the set of all levels of sub-scale aggregation considered for the template, i.e., for a scale  $s$ , the dilation rate was changed such that the receptive field of the filter was  $(w/s, h/s)$ ; and  $f$  represents one of the  $n$  different filters used for feature vector construction. The number of template representation vectors would then be  $n \cdot S$ . Since the number of operations for each filter is significantly less than the template window size  $w \times h$ , multiple filters can efficiently define different template representations.

### 3.3. Similarity Score

The NN of the  $k$  cluster centers is calculated for each pixel in the query image and converted into a label image  $L_I$ . The distribution similarity score at each sub-window  $q_i \in Q$  in the query image is then defined as

$$\text{sim}(q_i, \tau) = \sum_{s \in S, \sigma_i \in \{\sigma_i\}_{i=0}^n} -w_{s, \sigma_i} \cdot |L_{q, s}^{\sigma_i} - L_{\tau, s}^{\sigma_i}| \quad (4)$$

where  $w_{s, \sigma_i}$  is the weight assigned to the difference in the representations of filter  $\sigma_i$  at scale  $s$ .

## 4. Computational Complexity

**K-means computation:** The first step in the method is to calculate the  $k$  cluster centers of the  $w \times h \triangleq l$  template points. The k-means clustering algorithm complexity is  $O(klt)$  where  $t$  is the number of iterations. The complexity would be  $O(l)$  as  $k, t \ll l$ .

**NN Search:** The NN calculation is done by calculating the euclidean distance of all the points in the image  $W \times H \triangleq L$  to the  $k$  cluster centers. The complexity of this step is  $O(kL)$ .

**Filter Responses:** The indices of the NN are converted into a one-hot tensor of size  $W \times H \times k$  and an integral image is constructed in  $O(L)$ . Different coarse filters at different scales can be applied to aggregate the distribution. Each filtering operation with the dilated filter size  $(r_h + 1) \times (r_w + 1)$  has a time complexity of  $O(((r_h + 1) \cdot (r_w + 1) \cdot k \cdot L))$ . Multiple scale representations are then computed in  $O(s \cdot (r_h + 1) \cdot (r_w + 1) \cdot k \cdot L)$  where  $s$  is the maximum number of scales considered. Similarly, adding  $n$  different filters would increase the complexity by  $O(n)$ . Here  $r_h \cdot r_w$  can be ignored as  $r_h \cdot r_w \ll k, L$ . Hence, the overall complexity of this step would be  $O(nskL)$ .

**Similarity Score Calculation:** The distribution similarity score is calculated by subtracting the  $n \cdot s$  template distribution vectors from the query distribution vector in  $O(nskL)$ .

**Target Localization:** Maxima location requires another sweep through the image which takes another  $O(L)$  computations.

Overall, the complexity of our method would be  $O(l) + O(kL) + O(nskL) + O(nskL) + O(L) \approx O(nskL)$  where  $1 < s \cdot n \cdot k < l \ll L$ . Furthermore, the operations are easily parallelizable in a GPU to speed up the computations.



## 5. Experiments

**Datasets:** We evaluated our method on three datasets with different image sizes and challenges. The first dataset, BBS, a subset of the Online Object Tracking benchmark [42], consists of 90 video sequences with challenges such as complex deformation, occlusions, scale differences etc. It was acquired by sampling frames from [42] at different intervals. Three random pairs of frames with constant frame differences,  $dF = \{25, 50, 100\}$ , were extracted from each sequence. This resulted in three sub-datasets, namely, BBS25, BBS50, and BBS100, consisting of 270, 270, and 252 images, respectively. The sampling was repeated 5 times to extract more robust statistics. The image size in the dataset was relatively small, 320x480 or 320x240 pixels.

The TinyTLP dataset is based on a shortened version of the object-tracking dataset Track Long and Prosper (TLP) [27]. This dataset contains 50 video clips with 600 frames of size 1280x720 pixels. To reduce redundancy, we followed the protocol used in [35] and sampled 50 frames [1, 11, ..., 491] from each video and then used the 100th next frame as the query pair. The dataset created consists of 2500 template-query pairs.

The TLPattr dataset [27] is a collection of 91 short clips of different durations focusing on six different challenge attributes in the TLP dataset, namely, illumination variation, fast motion, background clutter, out-of-view or full occlusions, scale variations, and partial occlusions. The sequences in TLPattr are selected such that only one of the abovementioned challenge attributes is present in a sequence. There are 15 sequences belonging to a particular attribute, except for scale variation, which has 16 sequences. We randomly selected 15 images from each sequence as templates and chose the 100th next frame as the query image. Overall, the dataset consists of 1351 template-query pairs.

**Implementation:** We implemented our algorithm in `python` programming language using `PyTorch` [32] library to efficiently utilize GPU parallelization. All experiments were performed on a `Intel(R) Core(TM) i7-5930K` CPU and a `GeForce GTX TITAN X` GPU, respectively. Following a similar feature extraction strategy as in [36, 35], we investigated and compared two feature descriptors in our experiments: colour and deep features. The colour features were extracted by taking the 3x3 overlapping patches of RGB values. The deep features were extracted from a pre-trained ResNet-18 [12]. This differs from [36, 35] where the VGG model is used. ResNet-18 is a smaller and more parameter-efficient alternative with slightly lower performance on the

ImageNet dataset (69.7%, 89.1% vs 72.4%, 90.8% top-1 and top-5 accuracy) but with a significantly lower number of parameters (11.7 million vs 143.7 million). Here, the features were extracted from the bottleneck and consecutive stages in the model, concatenated and resized to the original input. The highest feature dimension considered in our experiments was 512, for which the output from the `conv1` (64), `conv2_x` (64), `conv3_x` (128), and `conv4_x` (256) layers were concatenated and resized to the original image size.

We used a fast GPU PyTorch-based clustering library<sup>1</sup> to perform K-means clustering to find  $k$  cluster centers to represent the template. The label aggregation was done using coarse filters implemented as dilated convolutional kernels. In our experiments, we used a  $3 \times 3$  gaussian kernel with a sigma of 2 and dilation of  $(w/3, h/3)$  to give equal weights to all the bins. This was intentional; we wanted to model the histogram comparison without any spatial location preferences and treat the performance as the baseline performance. For modelling the orientation, 2-rectangle and 3-rectangle haar filters were used by multiplying the gaussian kernel with the filter weights. The filter weights were modified according to Algorithm 1 to work on the integral images to compute the distribution aggregation in the query image efficiently. Multiple template vectors were generated and stored for comparison. We then found the exact NN for each feature vector in the query image in the  $k$  template cluster centers based on the L2 distance measure. The NNF was then converted to one-hot encoding to calculate the feature vector in each query sub-window.

**Quantitative Evaluation:** We evaluated the performance of the different representation schemes in our method using Intersection over Union (IoU) between the estimated window,  $\tau_x$  from the algorithm and the ground truth,  $\tau_{GT}$ . It is defined as follows:

$$IoU(\tau_x, \tau_{GT}) = \frac{\tau_x \cap \tau_{GT}}{\tau_x \cup \tau_{GT}}. \quad (5)$$

The performance for each dataset is shown as the mean IoU (MIOU), i.e. the mean of IoUs of all the samples, and Success Rate (SR), which is defined as the fraction of the samples in the dataset where  $IoU > 0.5$ .

**Performance Comparison:** We compare our method with the official code

---

<sup>1</sup>The k-means clustering library is available at: [https://github.com/DeMoriarty/fast\\_pytorch\\_kmeans](https://github.com/DeMoriarty/fast_pytorch_kmeans)

releases of DDIS [36] and DIWU [35] as they are the state-of-the-art as well as the closest approaches to ours. The reported time for our method includes the clustering, NN computation and similarity score calculation time. For DDIS and DIWU, the reported time includes the NN computation time and the score calculation time. The computation time of our method with DDIS and DIWU is not directly comparable as the computations for our method are performed on a GPU, and the official implementation of these methods uses a CPU for NNF computation. However, the relative runtime difference between the methods while using RGB vs deep features highlights the challenges in using DDIS and DIWU effectively, especially at higher resolutions.

### 5.1. Design Choices Evaluation

To quantify the effects of different design choices, we first evaluate the performance of our method on the BBS dataset for the codebook size ( $k$ ), multiple distribution aggregation scales, and w/wo the haar filters. We considered RGB and deep features with feature dimensions,  $d = \{27, 512\}$ , number of clusters,  $k = \{4, 8, 16, 32, 64, 128\}$ , total scales,  $S = \{1, 2, 3\}$ , and the inclusion of haar 2-rectangle and 3-rectangle filters, namely, haar-2x, haar-2y, haar-3x, and haar-3y. The weight of the similarity score for the final map from different scales is empirically set to  $1/s$ , 1.0 for the gaussian and 0.25 for each haar filter. Fig. 4 shows the performance on the BBS25, BBS50, and BBS100 datasets and Figure 5 shows the analysis on TinyTLP and TLPattr datasets. The following trends can be observed from Fig. 4 and 5.

**Increasing the codebook size boosts performance.** This is the general trend as shown within each plot (increasing number of clusters on the x-axis of MIOU plots) and between plots (increasing number of features from 27 to 512). With ResNet features of dimension 512 (4 second row, 5 bottom row, first and third columns) and the cluster size set to 128 (right-most points within the MIOU plots), the base method (blue dots), meaning scale of one and no haar filters, reaches close to the performance of DDIS and DIWU for BBS and TLP datasets. This shows the benefits of our method in greatly reducing the NN calculations while managing similar performance with just histogram comparison. Furthermore, even without including any shape information (i.e., the base method), our method shows the competitiveness of simple histogram matching in the NNFs with the similarity measures used in DDIS and DIWU. A similar trend of performance increase with the codebook size was observed with other intermediate feature dimensions ( $D=64, 128, 256$ ) of ResNet features.

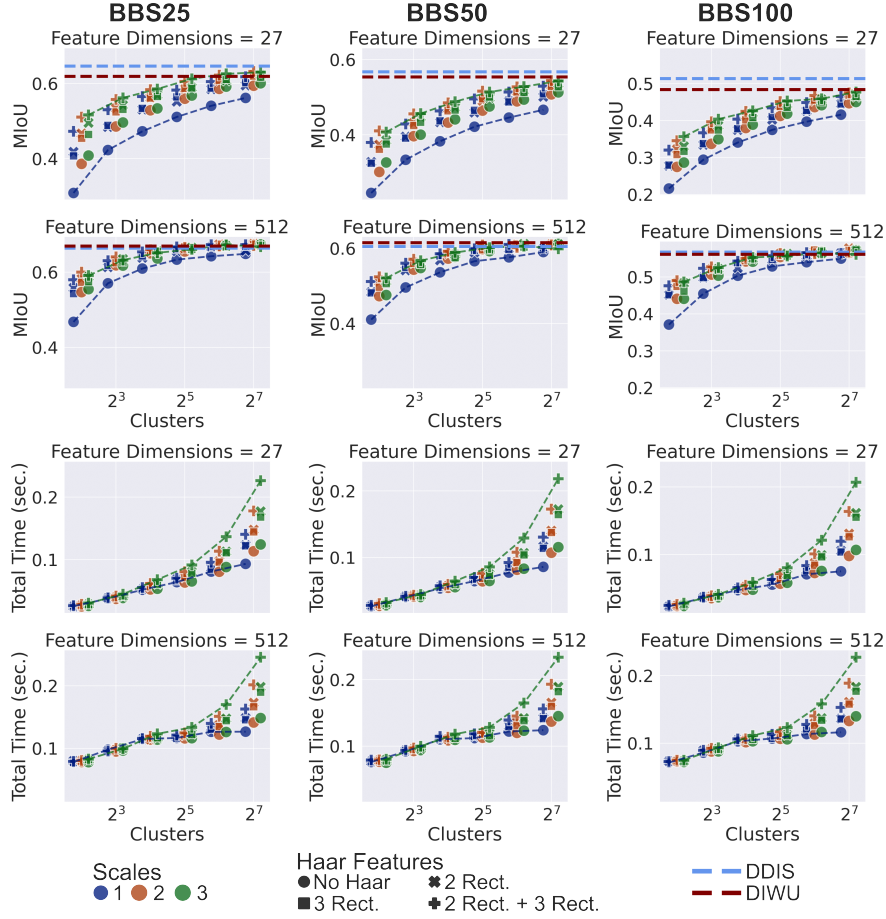


Figure 4: Average MIoU (top two rows) and Total time (bottom two rows) performance of our approach with different hyper-parameters on the BBS25, BBS50, and BBS100 datasets. The different scales shown in different colours are shifted slightly for better visualization. DDIS and DIWU performances are shown with the dashed lines in the MIoU plots. The lines in the plots highlight the trends for the base configuration without spatial information (blue line) and the configuration with the most filters and scales (green line) for each cluster size. The x-axis for all the plots is shown in the log scale.

### Multiple scale aggregation and haar filters boost the performance.

As seen in Fig. 4(top two rows) and Fig. 5 (first and third columns), the sub-scale aggregation and haar filters have clear gains in performance. This boost grows, as expected, smaller as cluster size and feature dimensions grow larger, as can be seen from the trend lines (blue vs green). This suggests that adding simple shape information using haar filters and sub-scale aggregation

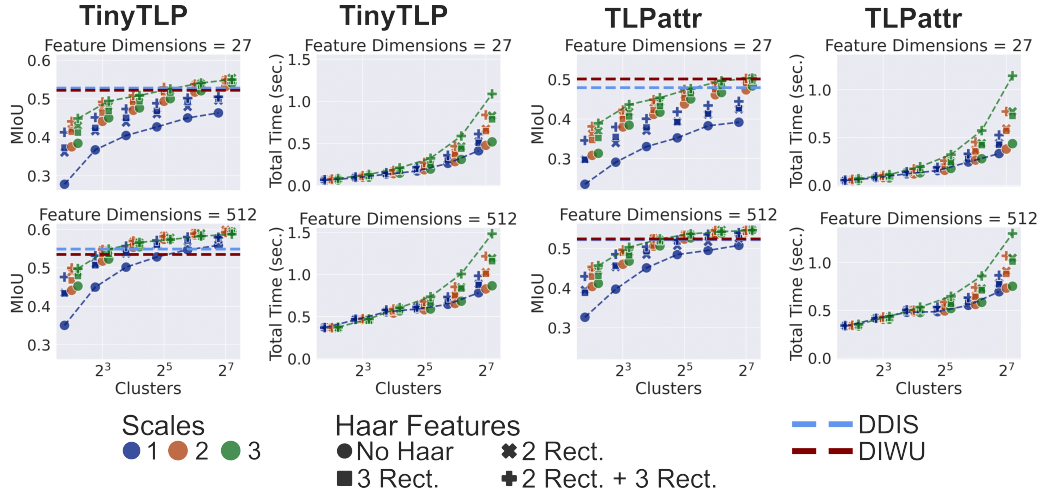


Figure 5: Average MIoU (top two rows) and Total time (bottom two rows) performance of our approach with different hyper-parameters on the TinyTLP and TLPattr datasets. The different scales shown in different colours are shifted slightly for better visualization. DDIS and DIWU performances are shown with the dashed lines in the MIoU plots. The lines in the plots highlight the trends for the base configuration without spatial information (blue line) and the configuration with the most filters and scales (green line) for each cluster size. The x-axis for all the plots is shown in the log scale.

can increase the distribution matching specificity at lower feature dimensions and cluster sizes. It also enables the method to reach the performance of DDIS and DIWU at high dimensions with smaller cluster sizes. The gains diminish much faster with increasing the codebook size for deep features than colour features. This might be attributed to the high feature-matching specificity of the deep features compared with the colour features.

**The execution time increases with cluster size and the number of filters used.** Fig. 4 (bottom two rows) and Fig. 5 (second and fourth column) show the total time taken by the different configurations of our method. The runtimes of colour and deep features differ by a constant  $\sim 0.1$  sec. for BBS datasets and  $\sim 0.5$  sec. for TLP datasets for the same configurations. The trendlines in the time plots (green and blue lines) show the runtime scales worse when using multiple filters (green) than just increasing the codebook size (blue).

The experiments show that multiple hyperparameter configurations can reach the performance of DDIS and DIWU. Here, for the final results, we chose the overall best configuration on BBS25, BBS50, and BBS100 for com-

Table 1: Performance results for BBS datasets. RGB features are marked with (C), and ResNet features are denoted by (D). The best results in each column are written in bold.

Method	BBS25		BBS50		BBS100		Total		
	SR	MIoU	SR	MIoU	SR	MIoU	SR	MIoU	Time
DDIS (C)	<b>0.781</b>	<b>0.638</b>	<b>0.695</b>	<b>0.567</b>	<b>0.594</b>	<b>0.493</b>	<b>0.690</b>	<b>0.566</b>	1.336
DIWU (C)	0.771	0.624	0.684	0.554	0.584	0.485	0.680	0.554	0.349
Ours best (C)	0.764	0.632	0.657	0.544	0.552	0.475	0.658	0.550	<b>0.171</b>
DDIS (D)	0.813	0.667	0.748	<b>0.613</b>	0.682	0.569	0.748	0.617	6.337
DIWU (D)	<b>0.821</b>	0.665	<b>0.751</b>	0.611	<b>0.690</b>	0.571	<b>0.754</b>	0.616	5.883
Ours best (D)	0.813	<b>0.675</b>	0.732	0.611	<b>0.690</b>	<b>0.577</b>	0.745	<b>0.621</b>	<b>0.167</b>

Table 2: Performance results for TinyTLP and TLPattr datasets. RGB features are marked with (C), and ResNet features are denoted by (D). The best results in each column are written in bold.

Method	TinyTLP			TLPattr		
	SR	MIoU	Time	SR	MIoU	Time
DDIS (C)	0.607	0.528	22.565	0.557	0.479	25.486
DIWU (C)	0.613	0.522	2.508	<b>0.600</b>	0.501	1.948
Ours best (C)	<b>0.638</b>	<b>0.554</b>	<b>0.935</b>	0.597	<b>0.505</b>	<b>0.772</b>
DDIS (D)	0.631	0.549	49.686	0.617	0.522	41.362
DIWU (D)	0.619	0.535	33.899	0.628	0.524	34.063
Ours best (D)	<b>0.694</b>	<b>0.598</b>	<b>1.040</b>	<b>0.655</b>	<b>0.549</b>	<b>0.911</b>

parison. We found that in BBS dataset, for RGB features, the best overall performance was achieved at  $k = 128$ ,  $S = 2$ , and using both haar 2-rectangle and 3-rectangle filters. For deep features, the best configuration was found to be at  $k = 128$ ,  $S = 2$ , and using only haar 2-rectangle filters. Our best configuration in the TLP datasets for RGB features was with  $k = 128$ ,  $S = 3$ , and the best configuration for the deep features was with  $k = 128$ ,  $S = 2$ , both with haar 2-rectangle filters. Hence, we fixed  $k = 128$  for the rest of the experiments and showcased the best configuration results.

## 5.2. Quantitative Results

The results for the BBS25, BBS50, and BBS100 datasets are shown in Tab. 1. The images in the datasets are relatively small (320x480 or 320x240). Our method lags slightly behind DDIS and DIWU for RGB features; however, for the deep features, the performance is marginally better than that of others. The runtime of our best method for the deep features is approximately the same as the RGB features, while the runtime increased around

4.7 times for DDIS (0.349 to 6.337 sec) and 16.8 times for DIWU (0.349 to 5.883). This further shows that increasing the feature dimensions in our method provides better performance without compromising the speed.

The results for the TinyTLP and TLPattr datasets are presented in Tab. 2. The images in the datasets are of higher resolution (1280x720), which makes the datasets suitable for showcasing the advantages provided by our method compared with DDIS and DIWU. Our method outperforms both DDIS and DIWU for RGB and deep features. For RGB features, our best configuration outperforms the DDIS by approx. 2.6% and 2.6%, and DIWU by approx. 2.2% and 0.4% in TinyTLP and TLPattr datasets, respectively. For deep features, our best configuration outperforms the DDIS by approx. 4.9% and 2.7%, and DIWU by approx. 6.3% and 2.5% in TinyTLP and TLPattr datasets, respectively. Here, similar to the BBS performance, the RGB and deep features runtimes of our method are marginally different for both datasets. In contrast, for DDIS, the runtime increased by 2.2 (22.5 to 49.7 secs) and 1.6 (25.5 to 41.4 secs) times for TinyTLP and TLPattr, respectively, and for DIWU, the runtime increased by 13.6 (2.5 to 33.9 secs) and 17.5 times (1.9 to 34.1 secs) for TinyTLP and TLPattr, respectively.

### 5.3. Qualitative Results

Following the quantitative results in the previous section, we also show the qualitative comparison of our method with DDIS and DIWU. Fig. 6 and 7 show the results of the methods on the BBS and TLPattr datasets, respectively. Although the heatmaps appear noisier than the compared methods, our method successfully finds the template in the query image. Furthermore, our method seems robust in cases where objects similar to the template are present, while DDIS and DIWU fail to detect the template successfully.

### 5.4. Ablation Study

We perform ablation studies to assess our method’s robustness under rotation, scaling, attributes, and dimensionality reduction challenges.

*Note:* Even though all the Python random seeds were fixed in the experiments, the inherent system randomness might cause the evaluation metric results to differ in different experimental runs. For fair evaluation, each ablation experiment is conducted within the same script run to minimize such randomness within the experiment. This also showcases the spread of the performance of our method and indicates the possible randomness when using it in real settings.

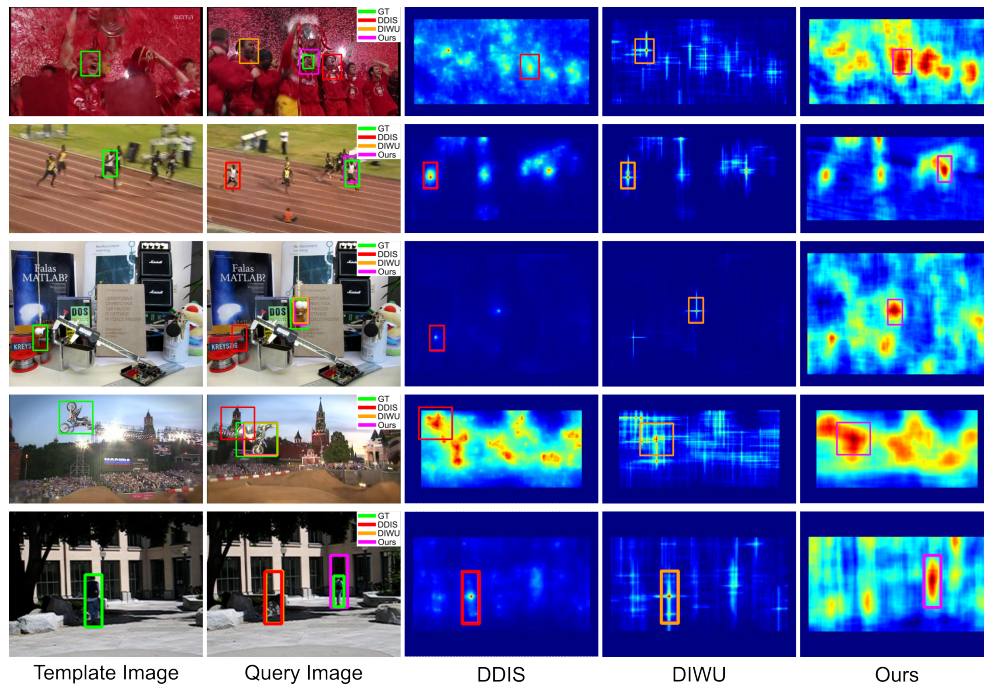


Figure 6: Qualitative results for the BBS dataset. Images were chosen to highlight the differences between the methods. The first column shows the template image with the template marked in green. The second column shows the query image with the results of DDIS (in red), DIWU (in orange), and our method (in pink). The third, fourth, and fifth columns show the heatmaps of the DDIS, DIWU, and our method, respectively, with the predicted bounding box marked in the same colours as the second column.

#### 5.4.1. Rotation Experiments

We constructed new datasets from BBS and TinyTLP by rotating the query images in the dataset by an angle  $\theta$  while keeping the same template image and box same. The query box is also rotated at the same angle and treated as the new ground truth. The prediction box dimensions are chosen to be the same as the rotated template box dimensions for better comparison. We considered  $\theta \in \{60^\circ, 120^\circ, 180^\circ\}$  and compared the results with the original dataset performance without rotation. We chose to show the performance of our method on RGB features and highlight the results of all the spatial filters on our best-performing aggregation scale. Deep features were not considered as the runtime of DDIS and DIWU with deep features is quite high.

Fig. 8 shows the performance of the rotation experiment on the datasets.



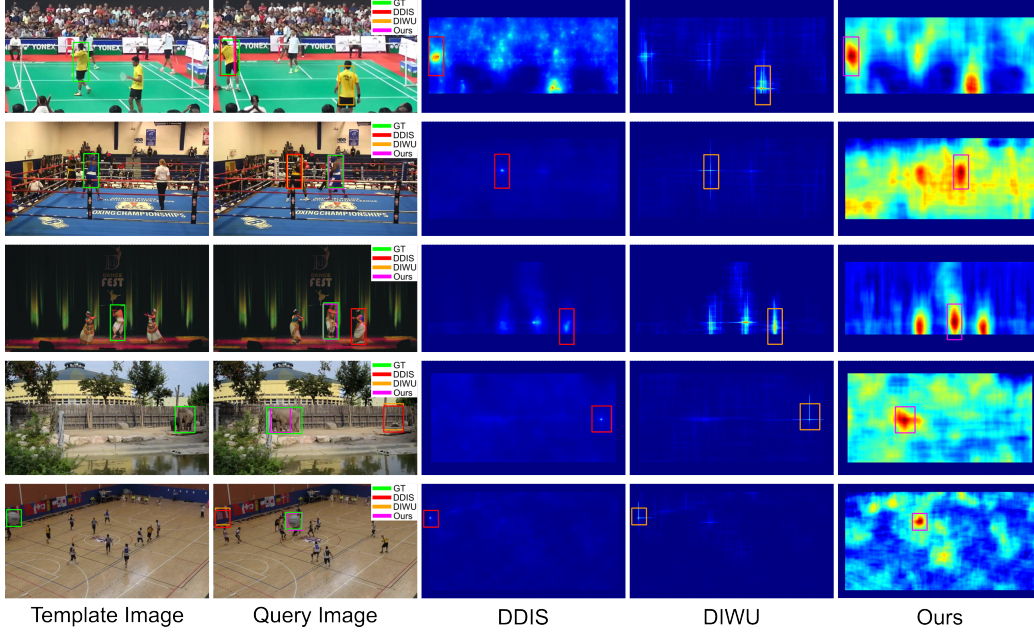


Figure 7: Qualitative results for the TLPattr dataset. Images were chosen to highlight the differences between the methods. The first column shows the template image with the template marked in green. The second column shows the query image with the results of DDIS (in red), DIWU (in orange), and our method (in pink). The third, fourth, and fifth columns show the heatmaps of the DDIS, DIWU, and our method, respectively, with the predicted bounding box marked in the same colours as the second column.

DDIS and DIWU performance drops sharply with increased rotation. As expected, DDIS and DIWU are not robust against large rotations as their deformation penalty relies on the pixel distance. However, the performances of the different configurations of our method do not degrade as sharply. For BBS datasets, the performance of our method at  $\theta = 60^\circ$  is close to that of DDIS and DIWU; however, our method outperforms them at higher rotational deformations,  $\theta = 120^\circ, 180^\circ$ . For the TinyTLP dataset, our method outperforms DDIS and DIWU for every rotation angle. DDIS and DIWU perform well in the original dataset, however, with a higher penalty for larger deformations, DDIS and DIWU are not robust against large rotations. Our methods perform similarly to DDIS and DIWU for the original dataset, but the performance doesn't degrade as severely when rotations are present. The performance of our methods increases back up at  $\theta = 180^\circ$ , which can be attributed to the filters used. The similarity score from the gaussian filter

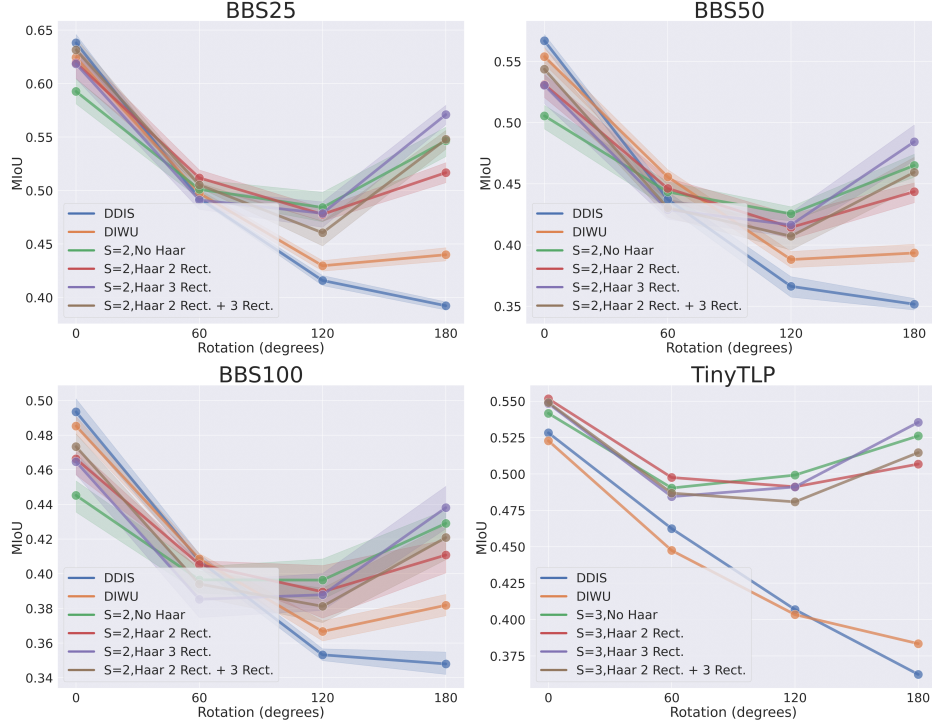


Figure 8: MIoU performance on the rotation experiments on the BBS25, BBS50, BBS100, and TinyTLP datasets.

representations would be high for  $\theta = 180^\circ$  as the filters are symmetric.

#### 5.4.2. Scale Experiments

We compare the effect on performance and runtime with respect to the image resolution on the TLPattr dataset. The original resolution is  $1280 \times 720$ ; hence, we downsample the images with different scale factors in the  $[1/4, 1]$  range. The template, query images, and ground truths are downsampled with the scale factor. Scale factors above one are not chosen as the image resolution is high enough to showcase the runtime variations, and upsampling doesn't add any meaningful information. The experiment also reflects the robustness of the performance if downsampling of the images is preferred to speed up the process.

Figure 9 show our method's MIoU and runtime performance along with DDIS and DIWU. As can be seen from the RGB feature accuracy plot, our methods' performance peaked at a lower resolution (between 0.5-0.75) and

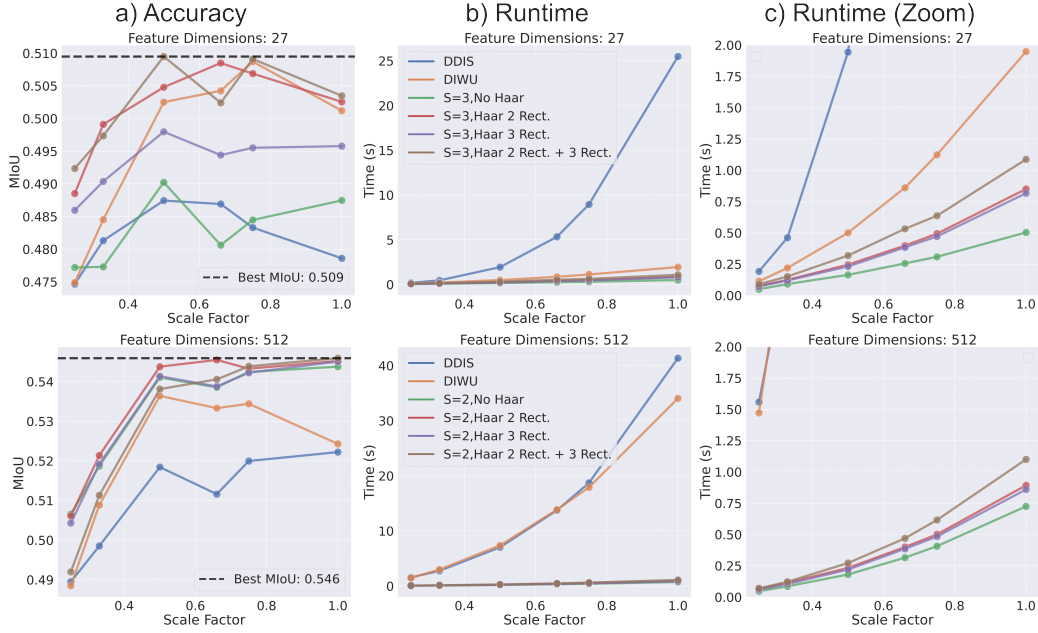


Figure 9: MIoU and Runtime performance at different scale factors on the TLPattr dataset. The x-axis shows the relative scale factor with the image of size  $1280 \times 720$ . The y-axis corresponds to a) The mean IoU of the algorithms with the RGB features (top) and deep features (below), b) the runtime of the algorithms and c) the zoomed-in runtime with the y-axis clipped to two seconds.

either remained the same or decreased. We hypothesize that RGB features provide limited benefits at higher resolutions. DIWU performs similarly to our method with haar filters, while DDIS performs similarly without haar filters. For deep features, all configurations of our method perform better than DDIS and DIWU. Furthermore, after the scale factor of 0.5, our methods show only fractional change in performance, which indicates that images can be processed at a lower resolution while maintaining performance.

The runtime analysis is shown in 9 b), and c). As can be seen from the figures, all configurations of our method scale better with the image resolution than DDIS and DIWU. For RGB features, the DIWU scales similarly to our methods with the resolution, but with deep features, both DDIS and DIWU scale equally poor with resolution.

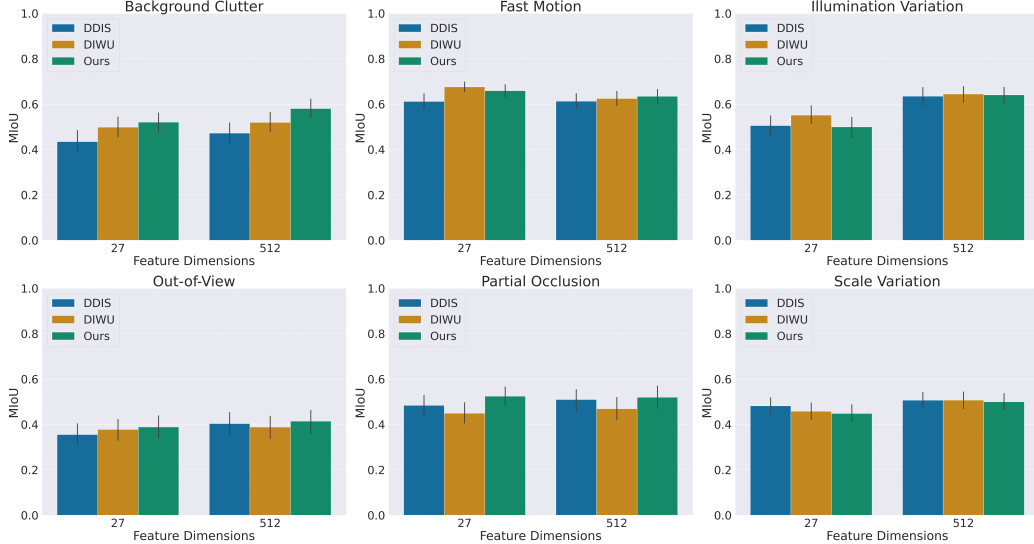


Figure 10: The MIoU performances of our best method with the DDIS and DIWU for the different challenge attributes present in the TLPattr dataset. DDIS, DIWU and our performance are shown in blue, yellow and green bars, respectively.

#### 5.4.3. Attribute Experiments

We evaluate the performance of our method on the various challenge attributes of the TLPattr dataset and again compare it with DDIS and DIWU. The attribute-wise performance for both RGB and deep features is shown in Fig. 10.

Our best method performs comparably to DDIS and DIWU for most challenge attributes. DDIS and DIWU rely on a smaller subset of good matches, making them robust to partial occlusions and background clutter. In contrast, our method relies on the distribution of all the matches. Our method outperforms them in both categories, which shows that distribution matching might be better at a higher resolution than NN diversity-based scoring. For the rest of the categories, our performance is similar to or slightly better than DDIS and DIWU.

#### 5.4.4. Dimensionality Reduction

Although our method makes the application of high-dimensional features relatively fast over the GPU, using a lower-dimensional representation for the NNF creation might still be desirable for lower memory and shorter runtime requirements. DDIS and DIWU use PCA dimensionality reduction

Table 3: Dimensionality reduction results for the BBS and TLP datasets. The best results in each column are written in bold.

Feature Dim.	BBS Datasets				TLP Datasets	
	BBS25	BBS50	BBS100	BBS(Total)	TinyTLP	TLPAttr
MIoU(D=512)	<b>0.677</b>	<b>0.614</b>	<b>0.574</b>	<b>0.621</b>	0.596	<b>0.546</b>
MIoU(D=18)	<b>0.677</b>	0.613	0.570	0.620	<b>0.597</b>	<b>0.546</b>
MIoU(D=9)	0.670	0.602	0.567	0.613	0.595	0.545
Time(D=512)	0.161	0.159	0.164	0.161	0.957	0.888
Time(D=18)	0.120	<b>0.119</b>	0.123	0.121	<b>0.660</b>	<b>0.677</b>
Time(D=9)	<b>0.118</b>	<b>0.119</b>	<b>0.120</b>	<b>0.119</b>	0.678	0.683

to reduce the computational time of the methods. The methods reduce the dimensionality of the features to  $D = 9$  first to speed up the NN search. Here, similarly, we explore the effects of dimensionality reduction by reducing the feature dimensions of our deep model from  $D = 512$  to  $D = 9$  and 18 and measure the relative change in performance on all the datasets. Only the results of our best-performing configuration on the datasets are shown.

Tab. 3 shows the MIoU and runtime of the dimensionality reduction on all the datasets. The overall performance of our method on the BBS datasets is reduced by approx. 1% of the original when using the same dimensions ( $D = 9$ ) as used in the DDIS and DIWU. The runtime of our method, however, was reduced by approx. 25%. For  $D = 18$ , a similar performance to the original was achieved while reducing the runtime by a similar margin.

For TLP datasets, no significant performance reduction was noticed for  $D = 9, 18$ . The runtime for  $D = 9$  was reduced by approx. 29% and 23% for TinyTLP and TLPAttr datasets, respectively. This further shows the robustness of the similarity measure on high-resolution datasets, even with dimensionality reduction. We also noticed that the runtime of  $D = 18$  was lower than that of  $D = 9$ . This can be due to the K-means algorithm being optimized for faster runtime of higher dimensional features with more data points.

## 6. Conclusion

We present a fast and robust template-matching framework using vector quantized NNFs, outperforming the previous state-of-the-art methods on high-resolution images while greatly reducing the NN computation cost. Our method efficiently uses the NN matching paradigm, is easily parallelizable to

be implemented on a GPU, and scales well with the image size. Furthermore, the similarity measure is comparable across images, unlike the other fast NNF-based approaches like DDIS and GAD, where the scores are based on image statistics. This can be useful in speeding up the process in applications of template matching where multiple images are considered, e.g., annotation applications and object tracking. We showcase the robustness of the method through exhaustive ablation studies. We show that state-of-the-art performance can be reached with simple coarse filters. This work paves the way for more robust and flexible template matching with more complex and custom-designed filters in the NNFs. The aforementioned advantages make this approach suitable for human-in-the-loop systems where the parameters can be adjusted, and the system’s performance can be validated quickly by humans and applied to a large set of images.

## 7. Acknowledgement

The work was supported in part by the Swedish Foundation for Strategic Research under Grant BD15-0008SB16-0046 and in part by the European Research Council under Grant ERC-2015-CoG 683810.

## References

- [1] Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S., 2012. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence* 34, 2274–2282.
- [2] Aumüller, M., Bernhardsson, E., Faithfull, A., 2020. Ann-benchmarks: A benchmarking tool for approximate nearest neighbor algorithms. *Information Systems* 87, 101374.
- [3] Barnes, C., Shechtman, E., Finkelstein, A., Goldman, D.B., 2009. Patch-match: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.* 28, 24.
- [4] Ben-Zrihem, N., Zelnik-Manor, L., 2015. Approximate nearest neighbor fields in video, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5233–5242.

- [5] Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L., 2017a. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* 40, 834–848.
- [6] Chen, L.C., Papandreou, G., Schroff, F., Adam, H., 2017b. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587* .
- [7] Chen, Z., Jin, H., Lin, Z., Cohen, S., Wu, Y., 2013. Large displacement optical flow from nearest neighbor fields, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2443–2450.
- [8] Cheng, J., Wu, Y., AbdAlmageed, W., Natarajan, P., 2019. Qatm: Quality-aware template matching for deep learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11553–11562.
- [9] Dekel, T., Oron, S., Rubinstein, M., Avidan, S., Freeman, W.T., 2015. Best-buddies similarity for robust template matching, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2021–2029.
- [10] Diete, A., Sztyler, T., Stuckenschmidt, H., 2017. A smart data annotation tool for multi-sensor activity recognition, in: *2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*, IEEE. pp. 111–116.
- [11] Gupta, A., Sabirsh, A., Wählby, C., Sintorn, I.M., 2022. Simsearch: A human-in-the-loop learning framework for fast detection of regions of interest in microscopy images. *IEEE journal of biomedical and health informatics* 26, 4079–4089.
- [12] He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- [13] Jamriška, O., Fišer, J., Asente, P., Lu, J., Shechtman, E., Šykora, D., 2015. Lazyfluids: Appearance transfer for fluid animations. *ACM Transactions on Graphics (TOG)* 34, 1–10.

- [14] Jia, D., Cao, J., Song, W.d., Tang, X.l., Zhu, H., 2016. Colour fast (cfast) match: fast affine template matching for colour images. *Electronics Letters* 52, 1220–1221.
- [15] Jin, J., Dauwels, J., Cash, S., Westover, M.B., 2014. Spikegui: Software for rapid interictal discharge annotation via template matching and on-line machine learning, in: 2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, IEEE. pp. 4435–4438.
- [16] Kat, R., Jevnisek, R., Avidan, S., 2018. Matching pixels using co-occurrence statistics, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1751–1759.
- [17] Khan, M.A., Smith, M.J., 2005. 5.3 - fundamentals of vector quantization, in: BOVIK, A. (Ed.), *Handbook of Image and Video Processing (Second Edition)*. second edition ed.. Academic Press, Burlington. Communications, Networking and Multimedia, pp. 673–688. URL: <https://www.sciencedirect.com/science/article/pii/B9780121197926501030>, doi:<https://doi.org/10.1016/B978-012119792-6/50103-0>.
- [18] Kim, H.Y., De Araújo, S.A., 2007. Grayscale template-matching invariant to rotation, scale, translation, brightness and contrast, in: *Advances in Image and Video Technology: Second Pacific Rim Symposium, PSIVT 2007 Santiago, Chile, December 17-19, 2007 Proceedings 2*, Springer. pp. 100–113.
- [19] Korman, S., Milam, M., Soatto, S., 2018. Oatm: Occlusion aware template matching by consensus set maximization, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2675–2683.
- [20] Korman, S., Reichman, D., Tsur, G., Avidan, S., 2013. Fast-match: Fast affine template matching, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2331–2338.
- [21] Lai, J., Lei, L., Deng, K., Yan, R., Ruan, Y., Jinyun, Z., 2020. Fast and robust template matching with majority neighbour similarity and annulus projection transformation. *Pattern recognition* 98, 107029.



- [22] Lan, Y., Wu, X., Li, Y., 2021. Gad: A global-aware diversity-based template matching method. *IEEE Transactions on Instrumentation and Measurement* 71, 1–13.
- [23] Lan, Y., Xiang, X., Zhang, H., Qi, S., 2020. A highly efficient and robust method for nnf-based template matching, in: *2020 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, IEEE. pp. 1–6.
- [24] Lebbink, M.N., Geerts, W.J., van der Krift, T.P., Bouwhuis, M., Hertzberger, L.O., Verkleij, A.J., Koster, A.J., 2007. Template matching as a tool for annotation of tomograms of stained biological structures. *Journal of Structural Biology* 158, 327–335.
- [25] Li, L., Han, L., Ding, M., Cao, H., Hu, H., 2021. A deep learning semantic template matching framework for remote sensing image registration. *ISPRS Journal of Photogrammetry and Remote Sensing* 181, 205–217.
- [26] Lucas, B.A., Himes, B.A., Xue, L., Grant, T., Mahamid, J., Grigorieff, N., 2021. Locating macromolecular assemblies in cells by 2d template matching with cistem. *Elife* 10, e68946.
- [27] Moudgil, A., Gandhi, V., 2018. Long-term visual object tracking benchmark, in: *asian conference on computer vision*, Springer. pp. 629–645.
- [28] Nagaraj, P., Muneeswaran, V., Muthamil Sudar, K., Hammed, S., Lokesh, D.L., Samara Simha Reddy, V., 2022. An exemplary template matching techniques for counterfeit currency detection, in: *Second International Conference on Image Processing and Capsule Networks: ICIPCN 2021 2*, Springer. pp. 370–378.
- [29] Oron, S., Dekel, T., Xue, T., Freeman, W.T., Avidan, S., 2017. Best-buddies similarity—robust template matching using mutual nearest neighbors. *IEEE transactions on pattern analysis and machine intelligence* 40, 1799–1813.
- [30] Ouyang, W., Tombari, F., Mattoccia, S., Di Stefano, L., Cham, W.K., 2011. Performance evaluation of full search equivalent pattern matching algorithms. *IEEE transactions on pattern analysis and machine intelligence* 34, 127–143.

- [31] Park, K., Patten, T., Prankl, J., Vincze, M., 2019. Multi-task template matching for object detection, segmentation and pose estimation using depth images, in: 2019 International Conference on Robotics and Automation (ICRA), IEEE. pp. 7207–7213.
- [32] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S., 2019. Pytorch: An imperative style, high-performance deep learning library, in: Advances in Neural Information Processing Systems 32. Curran Associates, Inc., pp. 8024–8035. URL: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [33] Puranic, A., Deepak, K., Umadevi, V., 2016. Vehicle number plate recognition system: a literature review and implementation using template matching. International Journal of Computer Applications 134, 12–16.
- [34] Sun, Y., Mao, X., Hong, S., Xu, W., Gui, G., 2019. Template matching-based method for intelligent invoice information identification. IEEE access 7, 28392–28401.
- [35] Talker, L., Moses, Y., Shimshoni, I., 2018. Efficient sliding window computation for nn-based template matching, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 404–418.
- [36] Talmi, I., Mechrez, R., Zelnik-Manor, L., 2017. Template matching with deformable diversity similarity, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 175–183.
- [37] Thomas, L.S., Gehrig, J., 2020. Multi-template matching: a versatile tool for object-localization in microscopy images. BMC bioinformatics 21, 1–8.
- [38] Tian, Y., Narasimhan, S.G., 2012. Globally optimal estimation of non-rigid image distortion. International journal of computer vision 98, 279–302.

- [39] Viola, P., Jones, M., 2001. Rapid object detection using a boosted cascade of simple features, in: Proceedings of the 2001 IEEE computer society conference on computer vision and pattern recognition. CVPR 2001, Ieee. pp. I–I.
- [40] Wang, S., Wang, H., Zhou, Y., Liu, J., Dai, P., Du, X., Wahab, M.A., 2021. Automatic laser profile recognition and fast tracking for structured light measurement using deep learning and template matching. *Measurement* 169, 108362.
- [41] Wei, Y., Xiao, H., Shi, H., Jie, Z., Feng, J., Huang, T.S., 2018. Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7268–7277.
- [42] Wu, Y., Lim, J., Yang, M.H., 2013. Online object tracking: A benchmark, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2411–2418.
- [43] Xu, Y., Zhang, J., Brownjohn, J., 2021. An accurate and distraction-free vision-based structural displacement measurement method integrating siamese network based tracker and correlation-based template matching. *Measurement* 179, 109506.
- [44] Yang, H., Huang, C., Wang, F., Song, K., Yin, Z., 2019a. Robust semantic template matching using a superpixel region binary descriptor. *IEEE transactions on image processing* 28, 3061–3074.
- [45] Yang, H., Huang, C., Wang, F., Song, K., Zheng, S., Yin, Z., 2019b. Large-scale and rotation-invariant template matching using adaptive radial ring code histograms. *Pattern Recognition* 91, 345–356.
- [46] Ye, Y., Bruzzone, L., Shan, J., Bovolo, F., Zhu, Q., 2019. Fast and robust matching for multimodal remote sensing image registration. *IEEE Transactions on Geoscience and Remote Sensing* 57, 9059–9070.
- [47] Yu, F., Koltun, V., 2015. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122* .
- [48] Zhang, C., Akashi, T., 2015. Fast affine template matching over galois field., in: BMVC, pp. 121–1.

- [49] Zhang, Y., Zhang, C., Akashi, T., 2022. Ds-sri: Diversity similarity measure against scaling, rotation, and illumination change for robust template matching. *IET Image Processing* 16, 2738–2751.
- [50] Zhang, Z., Yang, X., Jia, X., 2020. Scale-adaptive nn-based similarity for robust template matching. *IEEE Transactions on Instrumentation and Measurement* 70, 1–9.
- [51] Zhou, T., Lu, Y., Lv, F., Di, H., Zhao, Q., Zhang, J., 2015. Abrupt motion tracking via nearest neighbor field driven stochastic sampling. *Neurocomputing* 165, 350–360.