

1 **An accurate and distraction-free vision-based structural displacement measurement**
2 **method integrating Siamese network based tracker and correlation-based template**
3 **matching**

4 Yan Xu^a, Jian Zhang^{a*}, James Brownjohn^b

5 ^a Jiangsu Key Laboratory of Engineering Mechanics, Department of Engineering Mechanics,
6 Southeast University, Nanjing 210096, China.

7 ^b Vibration Engineering section, College of Engineering, Mathematics and Physical Sciences,
8 University of Exeter, Exeter EX4 4QF, UK.

9 * corresponding author: jian@seu.edu.cn

10 **ABSTRACT**

11 Vision-based displacement measurement receives increasing attention on non-contact bridge
12 monitoring while it faces challenges in long-time field applications due to the presence of
13 environmental variations. To overcome this issue, this study proposes a novel distraction-free
14 displacement measurement approach by integrating deep learning-based Siamese tracker with
15 correlation-based template matching. The Siamese tracker used applies deep feature
16 representations and learned similarity measures for image matching and also considers adaptive
17 template update with time. Since the estimated bounding boxes by the Siamese tracker have
18 size changes within frame sequences, a correction step is added to remove the centroid drifts
19 between the template and the predicted target regions using correlation-based template
20 matching. The proposed method is validated first in an indoor test and then implemented in
21 monitoring tests on a short-span footbridge and a long-span road bridge, demonstrating its
22 potential to handle challenging scenarios including partial occlusion, illumination changes,
23 background variations and shade effects.

24 **KEYWORDS**

25 Displacement measurement, vision-based method, Siamese network, template matching,
26 background variations.

27 1 INTRODUCTION

28 Bridge displacement is a significant metric for bridge condition assessment and of great interest
29 to bridge owners. The displacement data collected during the normal operation reflect the
30 bridge serviceability condition while the data collected during controlled vehicle load testing
31 are useful for the estimation of load carrying capacity. For flexible bridges, the displacement
32 data also carry short-time dynamic performance induced by the wind or traffic.

33 Vision-based measurement receives increasing attention in bridge displacement monitoring due
34 to its advantages of non-contact, easy installation and cost effective, etc. Existing studies have
35 demonstrated the potential application on structural condition assessment including system
36 identification [1–3], finite element model calibration [4], damage detection [5] and bridge WIM
37 system [6].

38 *1.1 Review of existing vision-based displacement measurement methods*

39 Vision-based displacement measurement is the process of localising target patterns in image
40 sequences and converting the computed target motions in image plane to true structural
41 displacement via a projection relationship. Thus, target region localisation is the key component
42 with a few variants of methods available including correlation-based template matching, optical
43 flow estimation and sparse keypoint matching, etc.

44 Correlation-based template matching is an area-based image matching method that works by
45 searching in a new frame for an area most closely resembling a predefined template. The
46 similarity measure is usually applied to image intensity values of grayscale images over a
47 rectangle area with default resolution in pixel level. Interpolation schemes are added to refine
48 the resolution to sub-pixel level. The method has been widely applied in structural monitoring
49 from the earliest work on the Humber Bridge and Second Severn Crossing in 1990s [7,8] to
50 recent displacement measurement applications on a railway bridge [9], a long-span bridge [10]
51 and a high-rise building [11], etc. Instead of considering consistent two-dimensional rigid
52 motion over the target area, the digital image correlation is an extension mostly used in
53 experimental mechanics under large deformation defining a shape distortion function of the
54 tracking area [12]. It was implemented in a short-span railway bridge monitoring exercise [13]
55 but the large deformation assumption is usually unnecessary for bridge measurement purposes.

56 Correlation-based template matching method is sensitive to illumination variation, partial target
57 occlusion, partial shading, and background disturbance, etc. and thus is often difficult to
58 guarantee robust performance over a long time in outdoor field environmental conditions [14].

59 The classic optical flow estimation detects motions or flows of all pixels in an image resulting
60 from brightness pattern shift [15]. The apparent velocity of movement is computed by
61 variational approaches by minimising energy based on brightness constancy and spatial
62 smoothness [16]. The measured results inherently contain sub-pixel resolution and the method
63 was implemented for field monitoring tests on a footbridge [17] and bridge stay-cable vibration
64 measurement [18,19]. Another popular variant is phase-based optical flow estimation based on
65 local phase constancy assumption proposed by Fleet and Jepson [20]. The method mainly
66 focuses on the application of system identification, i.e. extracting modal frequencies and mode
67 shapes in laboratory tests [21,22] and identifying modal frequencies of high-rise tower
68 buildings [23]. In one recent work, Dong etc. [24] proposed a deep learning based full field
69 optical flow approach for displacement monitoring on a grandstand structure.

70 Sparse keypoint matching techniques apply the process of transforming an image into a
71 collection of sparse feature representations and then finding their correspondences among
72 image sequences using a suitable distance measure. The common descriptors for feature
73 representations are scale-invariant feature transform [25], speeded up robust features [26] and
74 Oriented FAST and Rotated BRIEF [27], etc. and the distance measures are usually the
75 Euclidean distance in feature space [28] for float-point based descriptors and the Hamming
76 distance [29] for binary descriptors. Khuc and Catbas [30,31] applied the FREAK and SIFT
77 methods for deformation measurement in a stadium structure and a railway bridge, and Ehrhart
78 and Lienhart [17,32] adopted the ORB method for deformation measurement in a short-span
79 footbridge. Three review works [14,33,34] summarised challenges faced by vision-based
80 displacement approaches, mainly concerning robustness with respect to environmental
81 variations and camera mounting instability. These limitations could impose measurement
82 uncertainty, especially for continuous long-time tests. Concerning tracking robustness under
83 environmental variations (e.g. target pattern, illumination and background variations), deep

84 learning (DL) based techniques could be a potential solution, learning patterns in visual inputs
85 and improving prediction performance using big data and plentiful computing resources.

86 *1.2 Review of deep learning based target tracking methods*

87 Computer vision techniques have been widely used in structural inspection and monitoring
88 applications, including surface defect detection [36], 3D reconstruction of structural geometry
89 [37], strain [38] and displacement monitoring, etc. Vision-based methods for outdoor
90 applications are highly susceptible to uncontrolled environmental conditions, such as lighting
91 variations, shadows, atmospheric interference and wind gusts [39]. In the task of image-based
92 defect detection, the research focus has recently moved from the earlier image processing
93 methods (e.g. edge and boundary detection, background subtraction and thresholding, etc.) to
94 DL techniques [39]. Most of them use a typical CNN or its variations to classify defeat images
95 pre-trained on a large dataset [41] and successfully applied on images of different lighting
96 conditions and viewing angles [42]. This could provide as a hint for the displacement
97 monitoring task.

98 DL techniques employ multiple, deep layers of neurons that capture underlying pattern
99 representations from a dataset, enabling them to learn richer abstractions of inputs. The classic
100 feature matching by SIFT, SURF and ORB, etc., are describing the sparse and salient key-
101 points by their local image gradient variables while the DL-based retrieval models for feature
102 representation usually compute hierarchical layer-wise representations, capturing increasingly
103 complex image characteristics. The DL-based feature retrieval models have validated to
104 outperform SIFT-like detectors [35], particularly in cases where SIFT contains many outliers
105 or cannot match a sufficient number of feature points. Besides the successful applications in
106 automatic visual defect detection, DL techniques have been applied in structural health
107 monitoring applications such as data anomaly detection in long-time monitoring data [43] and
108 computer vision-based vibration measurement and modal frequency identification [44], etc.

109 Convolutional neural networks (CNNs) are used primarily in computer vision applications.
110 The task of target localization using the end-to-end learning framework is generalised as a
111 classification problem where the decision boundary is obtained by online learning of a
112 discriminative classifier using image patches from the target object and the background. One

113 popular tracking framework is based on ‘Siamese networks’ following the template matching
114 concept. A Siamese network consists of two-branch CNNs with tied parameters. It implicitly
115 encodes the object template and the search region to deep feature representations in another
116 space and then fuses them with a specific tensor to predicts their similarity. The earlier work
117 by Bertinetto et al. [45] first proposed a Siamese tracker (SiamFC), followed by a few extension
118 works. CFNet [46] adds a correlation filter to the template branch and makes the Siamese
119 network shallower but more efficient. However, they are lack of bounding box regression
120 requiring multi-scale test of high computation efforts. The SiamRPN tracker [47] introduces
121 the region proposal network after the Siamese network and performs joint classification and
122 regression for tracking. The DaSiamRPN tracker [48] further introduces a distractor-aware
123 module and improves the discrimination power of the model. These trackers do not consider
124 updating the template image which is inadequate for long-term tracking in presence of
125 appearance changes, fast motion, or occlusion. Zhang et al. [49] proposed a convolutional
126 neural network (CNN) architecture, UpdateNet to learn an adaptive target template update
127 strategy given the initial template, the accumulated template and the template of the current
128 frame. The UpdateNet architecture is general and can be integrated into all existing Siamese
129 trackers. To adapt to the target’s scale and aspect ratio changes, the predicted bounding box is
130 designed to have size changes among frame sequences instead of a fixed size as in traditional
131 template matching. It has the advantage over the traditional template matching method to be
132 robust to occlusion, lighting variation and pattern changes, etc. However, the predicted
133 bounding box centroid might deviate from the template centre in the initial frame due to the
134 predicted size changes. Therefore, it is infeasible to directly apply the method on the structural
135 displacement measurement.

136 *1.3 Purpose of this study*

137 The field applications undergo environmental variations (e.g. illumination conditions, shadow,
138 partial occlusion and other variations), making it challenging for a vision system to achieve a
139 robust and accurate displacement measurement over a long time. To overcome these challenges,
140 this study proposes a novel distraction-free target tracking approach by integrating deep
141 learning-based Siamese tracker with traditional correlation-based template matching. There are

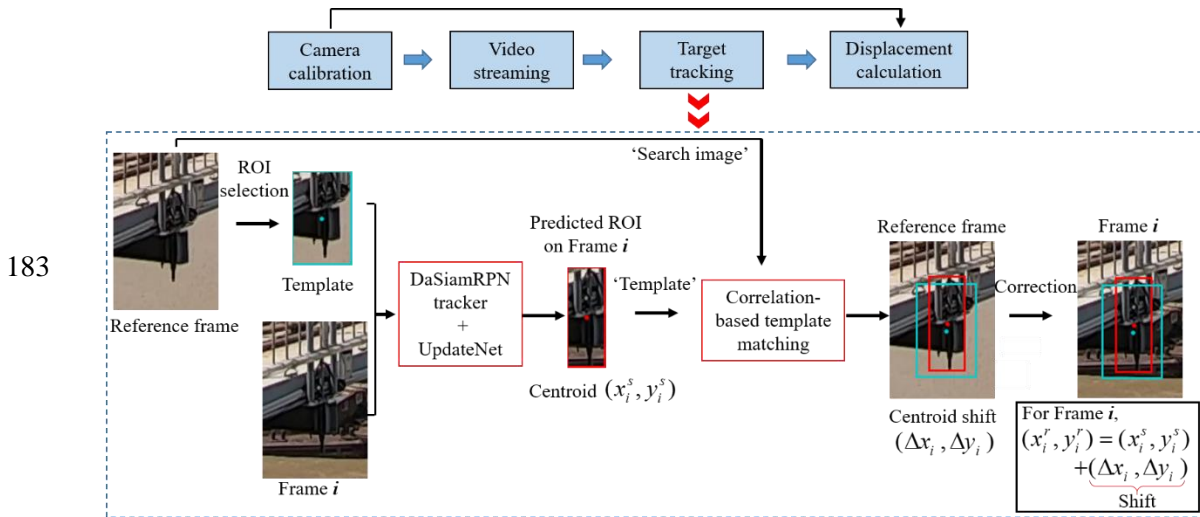
142 a few variations of Siamese trackers for template matching and the DaSiamRPN tracker
143 integrated with the UpdateNet for adaptive template updating is adopted in this work, which is
144 robust in challenging scenarios over long-term monitoring. Different from the fixed target size
145 setting in template matching, the Siamese tracker include a bounding-box regression layer to
146 predict target localisation, which consists of four regression coefficients, two-directional
147 position translation and size scaling of the bounding box. Since our task for structural
148 displacement measurement is based on the quantification of image translations of the target
149 region, a correction step is added to remove the centroid drifts between the template and the
150 predicted bounding boxes due to image size changes using correlation-based template matching.
151 The proposed method is validated first in a laboratory test and then implemented in monitoring
152 tests on a short-span footbridge and a long-span road bridge, demonstrating its potential to
153 handle challenging scenarios including occlusion, illumination and background changes.
154 To that end, section 2 introduces the basic principles for the DaSiamRPN tracker and the
155 UpdateNet template update scheme as well as our proposed method for structural displacement
156 measurement. Section 3 provides the information of an indoor validation test considering the
157 scenarios of occlusion and illumination changes. Section 4 and 5 demonstrates two field
158 monitoring tests on a short-span footbridge and a long-span road bridge in presence of
159 background changes and shade influence, respectively.

160 2 PROPOSED METHOD

161 A vision-based system comprises camera devices, a computer with video processing software
162 and accessories like a tripod. The video processing procedure could fit into a four-component
163 framework in Figure 1, i.e. camera calibration, video streaming, target tracking and
164 displacement calculation.

165 Camera calibration is to determine the projection transformation between the 2D image
166 coordinate system and the 3D structural coordinate system. The projection transformation could
167 be either scaling factor determined by the camera-to-target distance or planar homography
168 calibrated based on a few planar point correspondences. Target tracking is critical in the video
169 processing procedure to locate the target regions in the image plane through tracking methods.
170 The structural displacement could be easily derived from the outputs of the previous two steps.

171 Target tracking method is the main factor to influence the measurement accuracy and
 172 robustness. Although there are a few variants of target tracking methods with field validations,
 173 it is still challenging for a long-term monitoring campaign in the presence of environmental
 174 variations. To overcome this limitation, a novel target tracking approach is proposed in this
 175 study by integrating deep learning-based Siamese tracker with traditional correlation-based
 176 template matching. The main flowchart is shown in Figure 1. With the region of interest (ROI)
 177 selected in the reference frame, a Siamese tracker (DaSiamRPN + UpdateNet) is employed to
 178 predict the bounding box on the current frame (Frame i). Since the predicted bounding box has
 179 size changes compared with the template, the output of the Siamese tracker is corrected using
 180 correlation-based template matching for target position refinement. The principals of the
 181 DaSiamRPN and UpdateNet are introduced in Section 2.1 and Section 2.2 provides the
 182 framework of our proposed method.



184 Figure 1 Flowchart of the proposed target tracking method.

185 2.1 Siamese tracker for object tracking

186 The Siamese tracker selected for initial template matching is the DaSiamRPN tracker integrated
 187 with the UpdateNet. This is because the DaSiamRPN architecture considers a distractor-aware
 188 learning module and a local-to-global search region strategy, making it robust in challenging
 189 scenarios like occlusion, illumination change and other variations. Also, the UpdateNet is
 190 integrated to learn an adaptive target template update strategy to meet the long-term monitoring

191 requirement. The basic principles of the DaSiamRPN and the UpdateNet are briefly introduced
192 here and more details could refer to [48,49].

193 **DaSiamRPN**

194 The DaSiamRPN for object tracking is based on the SiamRPNBIG architecture extended from
195 the SiamRPN. The difference with the initial SiamRPN architecture [47] is that the image
196 dimension of the search area in the current frame is expanded from 255×255 to 271×271 and
197 the number of channels for CNN feature maps from 256 to 512 as shown in the Siamese branch.

198 As shown in Figure 2, it consists of a Siamese subnetwork for feature extraction and a region
199 proposal subnetwork for proposal extraction. Specifically, the Siamese network used is the
200 modified AlexNet [50], where the groups from the second and fourth convolutional layers
201 (conv2 and conv4) are removed. The Siamese feature extraction subnetwork consists of two
202 branches, the template branch with the target patch in the template frame as input and the
203 detection branch using the search region in the current frame as input. The two branches share
204 parameters in CNN so that the two patches are implicitly encoded by the same transformation
205 which is suitable for the subsequent tasks. The Region Proposal Network (RPN) subnetwork
206 consists of two branches, the foreground-background classification branch and bounding box
207 regression branch. The pair-wise correlation between the template feature map and current
208 feature map is firstly computed on both branches with the template feature maps used as kernels.
209 In the classification branch, the computed correlation features are passed through a softmax
210 layer to derive the classification scores representing negative and positive activation of each
211 anchor at corresponding locations on original map. In the regression branch, a linear regression
212 layer is employed to predict four regression coefficients representing the position and size
213 changes of the bounding box to refine the coordinates of the positive anchors. To train the RPN
214 network, the training loss used is a multi-task loss combined by the cross-entropy loss for the
215 classification branch and the smooth L1 loss with normalized coordinates for the regression
216 branch, respectively.

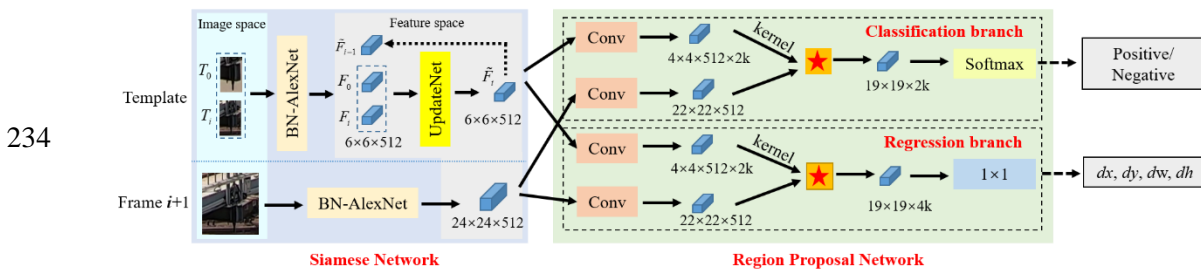
217 Since high quality training data is crucial for the success of end-to-end representation learning,
218 the DaSiamRPN framework includes a series of strategies to improve the generalization of the
219 learned features and eliminate the imbalanced distribution of the training data. One is to expand

220 the categories of positive pairs by introducing existing large-scale detection datasets and data
 221 augmentation techniques. Besides, diverse semantic negative pairs consisting of labelled targets
 222 both in the same and different categories are added in the training process.

223 In the online tracking process, a distractor-aware module is designed which can effectively
 224 transfer the general embedding to the current video domain. Distractors in context of the target
 225 are selected in each frame by the non-maximum suppression to generate a distractor set. Instead
 226 of directly using the cross correlation response between the template and the proposal with
 227 highest score in the embedding space as the similarity metric, this response is subtracted by the
 228 weighted sum of the cross correlation between the template and the distractor set to make full
 229 use of the negative label information.

230 To adapt the long-term tracking application which might include severe out-of-view or full
 231 occlusion, detection scores are taken as a metric indicating the tracking quality and an iterative
 232 local-to-global search strategy is designed to re-detect the target during failure cases.

233



235 Figure 2 Structure of DaSiamRPN with the SiamRPNBIG architecture integrated with the UpdateNet
 236 for template updating.

237 UpdateNet

238 In the target tracking step, the template is the basis to find the best candidate region in the new
 239 frame and hence a good template is crucial for robust object tracking. However, the
 240 effectiveness of the template cannot be guaranteed for long-term object tracking on site in
 241 presence of illumination, occlusion and background variations. Therefore, the template must be
 242 updated iteratively to improve the matching efficiency for robust long-term object tracking.

243 The UpdateNet framework is a CNN which aims to estimate the optimal template \tilde{F}_i for
 244 tracking the next frame ($i+1$) given the initial template (T_0), the template of the current frame

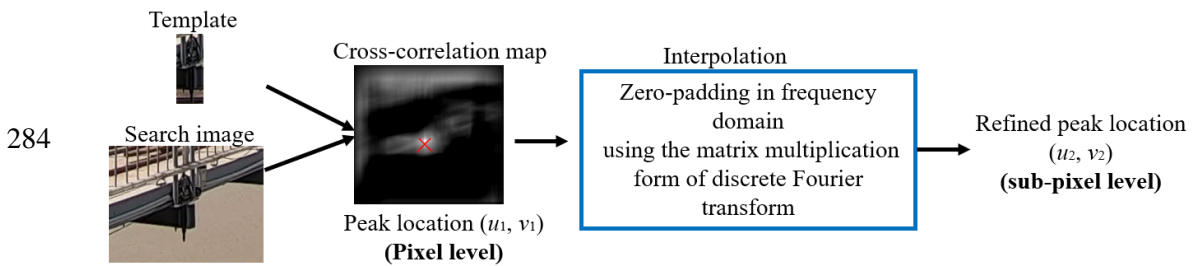
245 (T_i) and the last accumulated template in feature space (\tilde{F}_{i-1}). Specifically, the UpdateNet is a
 246 two-layer fully convolutional neural network: one $1 \times 1 \times 3 \cdot C \times 96$ convolutional layer, followed
 247 by a ReLU and a second convolutional layer of dimensions $1 \times 1 \times 96 \times C$ where $C = 512$. As
 248 shown in Figure 2, the input of the UpdateNet is the initial template F_0 , the last accumulated
 249 template \tilde{F}_{i-1} and the template of the current frame F_i in the embedding feature space by a
 250 fixed fully convolutional network and the output is updated accumulated template of the current
 251 frame \tilde{F}_i . The training process is by minimizing the Euclidean distance between the updated
 252 template and the ground-truth template of the next frame. A multi-stage training approach is
 253 employed, and the first stage involves template updating using the standard linear update. In
 254 the posterior stages, the UpdateNet model trained in the previous stage is applied to get
 255 accumulated templates and for object location predictions. The UpdateNet which is compact
 256 can easily be integrated into existing Siamese trackers and here it is employed together with the
 257 DaSiamRPN.

258 The Siamese tracker is end-to-end offline trained with large-scale image pairs and then online-
 259 tracking as a local one-shot detection task. For the application, the pre-trained model used in
 260 the DaSiamRPN is the SiamRPNBIG model provided by the authors [51] trained on VID [52],
 261 Youtube-BB [53], COCO Detection [54] and ImageNet Detection [52] datasets with data
 262 augmentation. For the UpdateNet, the pre-trained model is derived using a three-stage training
 263 on the VOT2018 dataset [55] with the achieved expected average overlap of 0.403.

264 2.2 Proposed method

265 The Siamese tracker has the advantage of localising target regions in challenging scenarios like
 266 occlusion, illumination change and other variations. The output of bounding box predictions
 267 are scale and aspect ratio varied to adapt to the target pattern changes among frame sequences.
 268 Therefore, it is infeasible to directly employ the image coordinate changes of the predicted
 269 bounding box centroids for structural displacement measurement. To solve this issue, an
 270 additional step is supplemented to calculate the centroid drift between the initial template frame
 271 and the predicted bounding box regions due to region size changes.

272 Considering that the traditional correlation-based template matching achieves sub-pixel
 273 accuracy in ideal scenarios, it is adopted here to quantify the centroid drifts of the bounding
 274 boxes following the Siamese tracker. The correlation-based template matching is the process
 275 of searching in the current frame for an area most closely resembling a predefined template in
 276 the initial frame. A target region is selected as the template that is a subset image in the reference
 277 frame. A matching criterion is defined to evaluate the similarity degree between the template
 278 and the new frame and the criterion used is zero-mean normalised cross correlation coefficient
 279 (ZNCC). The target location in the new frame corresponds to the peak location in the similarity
 280 matrix that has resolution at pixel level. Subpixel interpolation schemes [9] are required to
 281 refine the tracking results to sub-pixel level and the interpolation method used in this study is
 282 zero-padding in frequency domain using the matrix multiplication form of discrete Fourier
 283 transform [56].



285 Figure 3 Framework of correlation-based template matching.

286 Different from general template matching applications, the ‘template’ used as the reference
 287 here is the image subset in the current frame predicted by the Siamese tracker as the target
 288 region. As shown in in Figure 1, the ‘search image’ which has a larger image size than the
 289 template is a subset image cropped from the initial reference frame. Regarding the two image
 290 subsets are the predicted resembling area with the highest similarity score by the Siamese
 291 tracker, the search range of the peak location in the correlation map could be limited to a small
 292 value (e.g. 10 pixels from the centroid). This constraint could effectively avoid the drift error
 293 due to apparent target pattern changes in challenging scenarios.

294 The output of target centroid location in the current frame is the Siamese output corrected by
 295 the estimated centroid drift by the correlation-based template matching.

296 To convert to the structural displacement, the projection relationship between the structural
297 coordinate system and the image plane is pre-computed by the scaling factor using camera-to-
298 target distance or by the planar homography matrix derived from a few planar point
299 correspondences. It is noted that the tracked target is planar and coplanar with the computed
300 2D displacement.

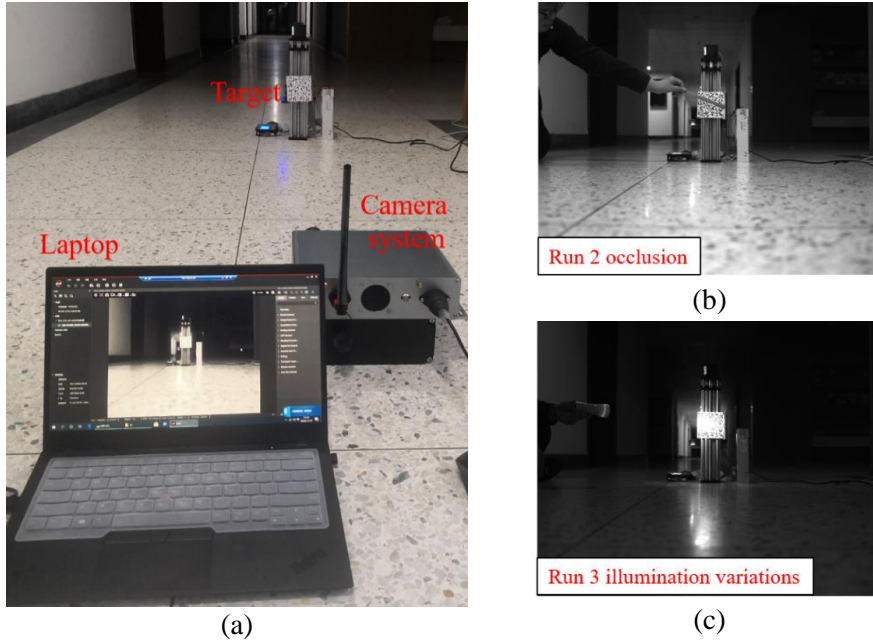
301 3 LABORATORY VALIDATION

302 To validate the effectiveness of the proposed method for structural displacement measurement
303 on challenging scenarios, an indoor test of reciprocating motions triggered by a linear actuator
304 was conducted considering two cases including partial occlusion and pattern variations due to
305 illumination. Section 3.1 and 3.2 describe the test setup and the results obtained, respectively.

306 *3.1 Test configuration*

307 A linear actuator was implemented to generate reciprocating motion with the amplitude of 8
308 cm. The test configuration is shown in Figure 4. A target of speckle patterns with the dimension
309 of 10 cm by 10 cm was attached to the centre of the cover tube.

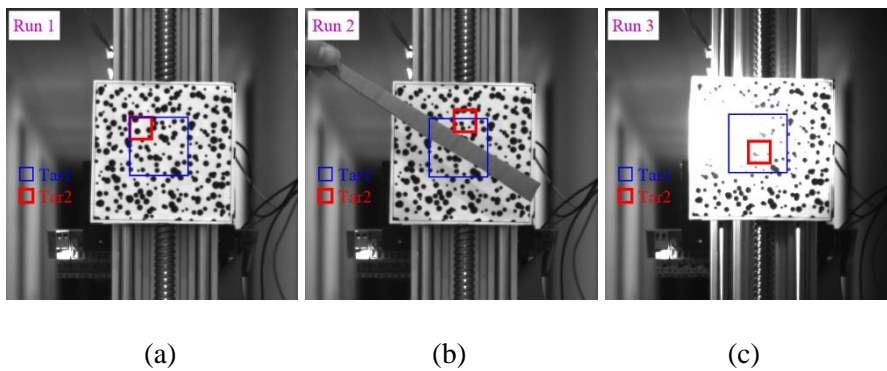
310 The video acquisition device used is an industry camera (Hikvision MV-CA050-20UM) with
311 the focal length of 16 mm which is arranged 1.77 m away from the target. The acquired images
312 are grayscale. Three runs of tests were conducted and the differences of test conditions occurred
313 in Run 2 and Run 3 are the presence of partially occlusion of the target and illumination changes
314 caused by an adjacent lamp, respectively shown in Figure 4 (b) and (c).



315 Figure 4 Test configuration (a) and sample frames captured by the camera in Run 2 (b) and Run
 316 3 (c).

317 3.2 Test Results

318 To demonstrate the working performance of the traditional correlation-based template matching
 319 methods (denoted as CM in the following study), two targets on the speckle patterns are selected
 320 for analysis with the dimensions of 30 pixels by 30 pixels and 80 pixels by 80 pixels,
 321 respectively. Figure 5(a) shows the template frame captured in Run 1 and the image was
 322 cropped for better visualisation.



324 (a) (b) (c)
 325 Figure 5 Selected target regions in the template frame in Run 1 (a), one sample frame at 14.2 s
 326 in Run 2 (b) and one frame at 19.5 s in Run 3(c).

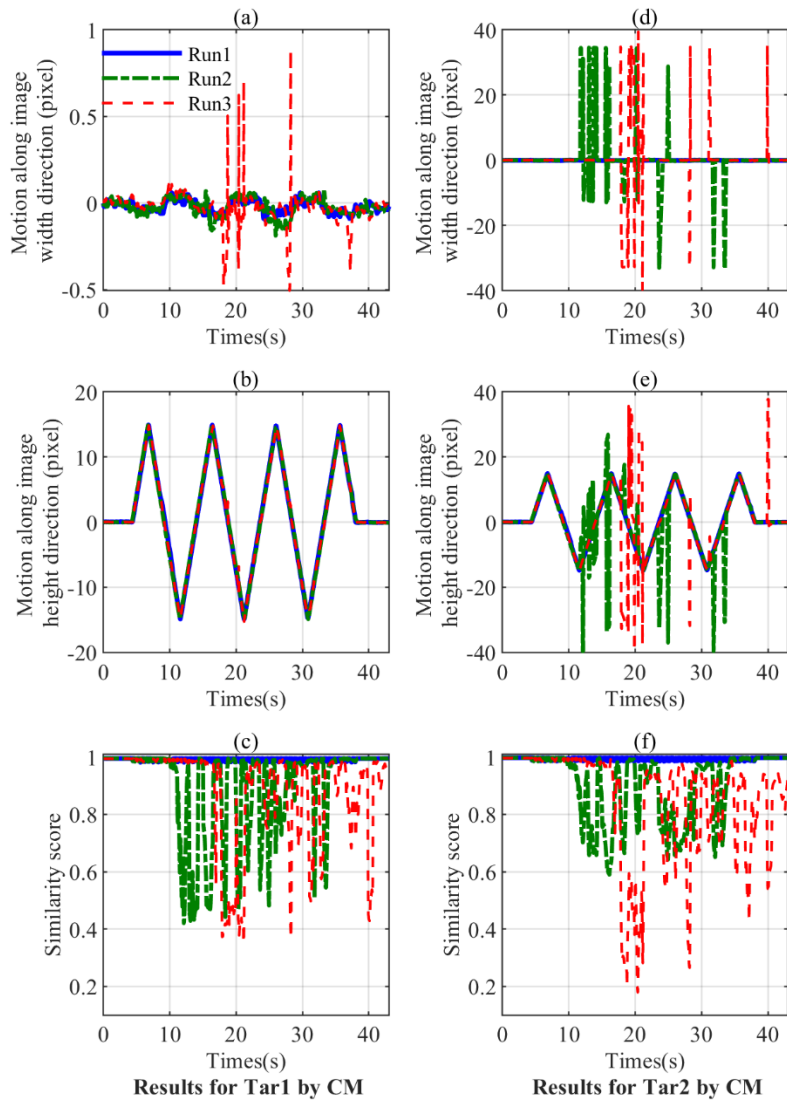
327 Tracked motions by the CM in the three runs are indicated in Figure 6. The left column of the
 328 figure presents the tracked motions along the image width and height direction and the
 329 similarity scores between the predicted target region and the template image for the target

330 region Tar1. The motion along the image width direction is expected to be zeros. For Tar1, the
331 root mean squared errors (RMSE) for three runs are 0.03 pixel, 0.05 pixel and 0.10 pixel while
332 the maximum deviations in three runs are 0.08 pixel, 0.17 pixel and 0.88 pixel. About the
333 motion along the image height direction, the result measured in the ideal case (Run 1) is taken
334 as the reference and the cross-correlation coefficients for the measurement in Run 2 and Run 3
335 both reach 99.9%. Due to the occurrence of partial occlusion in Run 2 and illumination
336 variations in Run 3, the similarity scores experience sharp decreases in some time periods
337 shown in Figure 6(c). Generally, the CM provides a reliable measurement for the larger target
338 Tar1 in presence of partial target pattern changes.

339 The right column of Figure 6 presents the measured results for the smaller target Tar2. Taking
340 the measurement in Run 1 as the reference, the measurement in Run 2 and Run 3 includes some
341 sharp large deviations with the amplitude over 20 pixels in both the image width and height
342 directions. Sample frames taken from Run 2 and Run 3 shown in Figure 5(b) and (c) indicate
343 that these apparent measurement deviations correspond to tracking failures when large pattern
344 changes occur on most of the target region.

345 Analysis results for both two targets demonstrate that the CM is not robust to severe target
346 pattern changes over time and setting a threshold on similarity score measures is necessary to
347 remove the measurement of low confidence or tracking failures. Also, it is better to select a
348 large target region with stable patterns and consistent motions for tracking while it might be not
349 available in field tests.

350

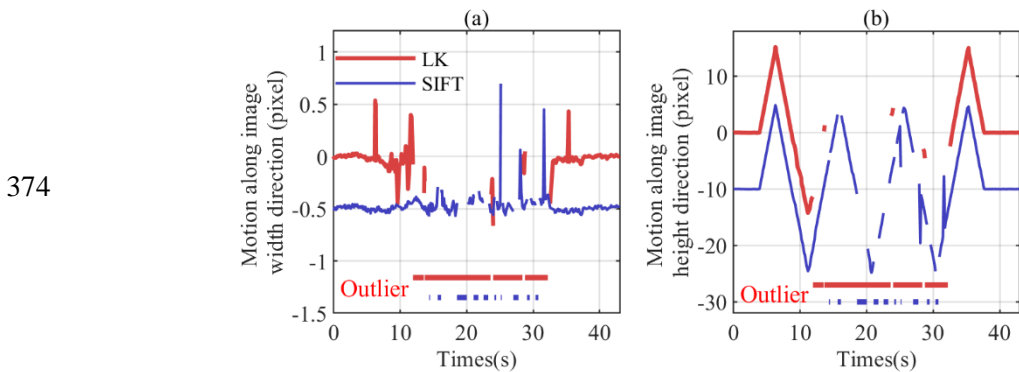


351 Figure 6 Image motions along two image directions and similarity scores using correlation-
352 based template matching for the target Tar1 (left column, a-c) and the target Tar2 (right
353 column, d-f).

354 As a comparison, two other classic target localisation methods, i.e. Lucas- Kanade (LK) optical
355 flow and the SIFT matching were implemented to predict the motions of the larger target Tar1
356 in Run 3. For the LK optical flow estimation, the feature points in the template image was
357 extracted first and then refined to sub-pixel level. Then the optical flow for the extracted sparse
358 features was estimated in the subsequent frames using the iterative Lucas-Kanade method with
359 pyramids. The outliers in feature point correspondences were filtered by apply RANSAC

360 geometric transformation estimation. The remaining feature points are evaluated by the re-
 361 projection error of the 2D translation to further check the geometric consistency. The averaged
 362 image coordinate movement between feature point correspondences were taken as the image
 363 motion of the target region. For the SIFT matching, the key-points were detected for the
 364 computation of their descriptors in the template and the subsequent frames independently. The
 365 FLANN (Fast Library for Approximate Nearest Neighbours) match was implemented to find
 366 nearest neighbours in two sets of descriptors. The detected matches were sorted by their
 367 distance with the first 50% closest matches kept. The sorted key-points were post-processed for
 368 outlier removal similar to the process in the LK.

369 The results given in Figure 7 are the averaged image motions of sparse key-points after
 370 RANSAC geometric verification. It shows that the measurements contain many outliers due to
 371 insufficient number of matched key-points or large re-projection error after RANSAC
 372 geometric transformation. Thus, these two methods are not robust for patterns in severe varying
 373 lighting conditions.



375 Figure 7 Image motions along image width (a) and height directions (b) for the target Tar1 in
 376 Run 3 measured by Lucas-Kanade (LK) optical flow and the SIFT matching.

377 The DaSiamRPN+UpdateNet and our proposed method (denoted as Siam and Siam+CM in the
 378 following study) were also implemented to analyse the video data by tracking the smaller target
 379 Tar2 and the results are shown in Figure 8. For the results measured by the Siam (left row), in
 380 Run 1, the measured motion in image width direction has a RMSE of 0.18 pixel and maximum
 381 deviation of 0.52 pixel while the motion in image height direction reaches a high similarity with
 382 results by the CM with the cross-correlation coefficient reaching 99.8%. In Run 2 and Run 3,
 383 the predicted bounding boxes deviate from the actual position with box size changes after 11s

384 and 18 s when the target patterns experience apparent variations as shown in Figure 9. The
 385 predicted bounding boxes have apparent size changes and shift to adjacent resembling area with
 386 salient patterns.

387 For the results measured by the Siam+CM (right row), the measured motion in image width
 388 direction has the RMSEs of 0.04 pixel, 0.04 pixel and 0.06 pixel for the three runs while the
 389 maximum deviations are 0.10 pixel, 0.09 pixel and 0.46 pixel, respectively. Compared with the
 390 measurement by the CM in Run 1, the motions in image height direction in the three runs are
 391 of high similarity with the cross-correlation coefficient over 99.9%.

392 The 2D structural displacement for the target region Tar2 by Siam+CM is shown in Figure 10.
 393 The amplitudes in the vertical direction are 8.0 cm consistent with the test settings. The
 394 horizontal movement has a similar shape with the amplitude of 0.1 mm which is the leakage of
 395 vertical components caused by the error in the definition of the structural coordinate system.

396 The measurement speed of the proposed method is provided in Table 1. The Computer
 397 hardware specifications are the CPU, Intel Core i9-10900K (10 cores, 20 siblings) and the GPU,
 398 NVIDIA GeForce RTX 2080 Ti (11G). The Programming language is Python in Linux
 399 environment. By making use of the multiple processing cores, the measurement speed for CM
 400 could be increased from 26.4 FPS (frame per second) to 166.4 FPS. The proposed method
 401 Siam+CM which runs sequentially the Siam and CM process, reaches the tracking speed of
 402 56.5 FPS which is sufficient for most real-time bridge monitoring applications.

403 Table 1 Measurement speed of the proposed method

Frame resolution (pixel)		2592×2048
Template dimension (pixel)		30×30
Search image dimension (pixel)		110×110
Measurement speed (FPS)	Siam	85.5
	CM	26.4
	CM (multi)	166.4
	Siam+CM (multi)	56.5

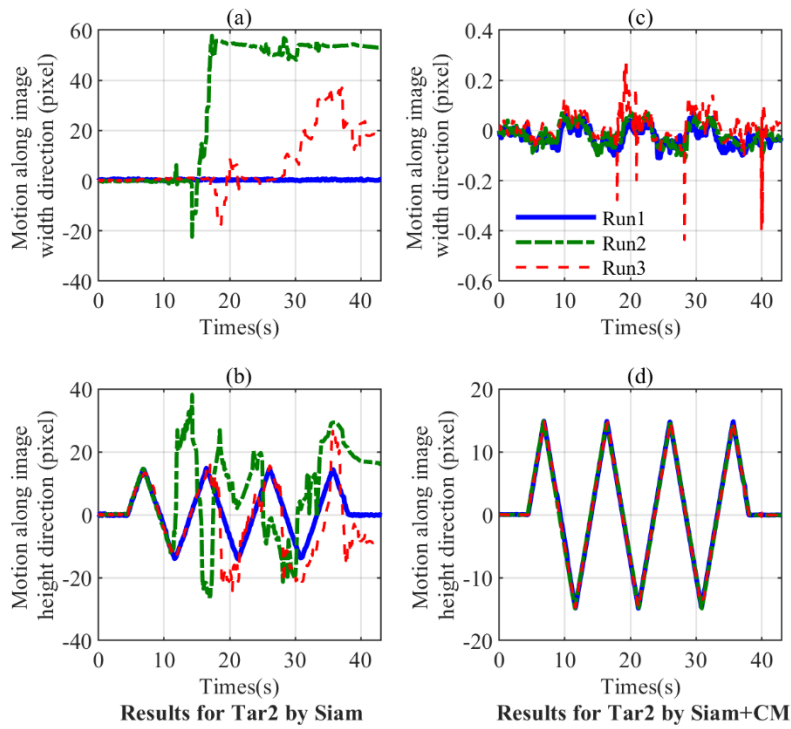
404 Observations indicate that the Siamese tracker (DaSiamRPN+UpdateNet) localises the target
 405 regions by size changes and shift to adjacent regions with high similarity in presence of target
 406 pattern changes. By supplementing the correlation-based template matching as a followed
 407 correction step, it reaches a robust and accurate measurement in challenging scenarios. The
 408 measurement results are evaluated by the actuator records and the root mean square (RMS)

409 errors are summarised in Table 2. The maximum RMS errors in three runs are 0.05 mm in
 410 horizontal direction and 0.158 mm in vertical direction. The measurement accuracy in cases of
 411 pattern occlusion and varying illumination are similar to that in ideal case.

412 Table 2 Measurement error of the proposed method in three runs

Run No.	Horizontal direction	Vertical direction	
	RMS error (mm)	RMS error (mm)	Correlation coefficients
Run 1	0.048	0.138	99.95%
Run 2	0.044	0.089	99.98%
Run 3	0.050	0.158	99.92%

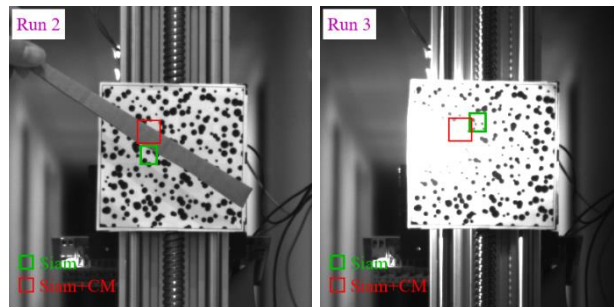
413



414

415 Figure 8 Image motions along two image directions for the target Tar2 using the
 416 DaSiamRPN+UpdateNet (left column, a-b) and our proposed method (right column, c-d).

417



418

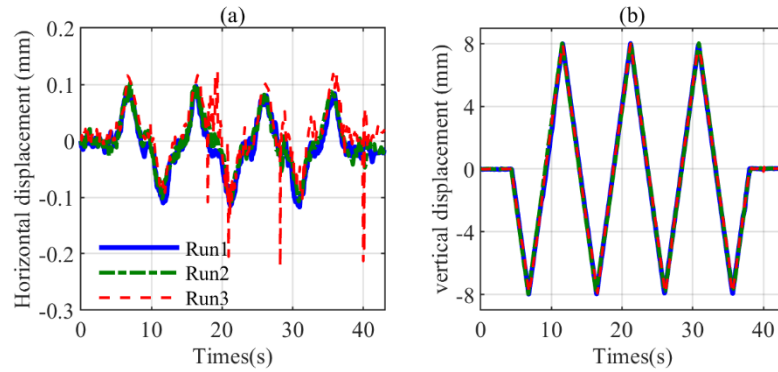
(a)

(b)

419 Figure 9 Demonstration of the predicted bounding boxes by the DaSiamRPN+UpdateNet and
420 our proposed method on the target Tar2: (a) cropped frame in Run 2 at 14.2 s; and (b) cropped
421 frame in Run 3 at 33.0 s.

422

423



424 Figure 10 Displacement time histories along horizontal (a) and vertical directions (b) on the
425 target Tar2 using our proposed method.

426 4 FIELD TEST ON A SHORT-SPAN FOOTBRIDGE

427 This section describes a case study of using the proposed target tracking method for measuring
428 displacement of a suspension footbridge during bridge rehabilitation status. The recorded video
429 streams include frequent background variations due to coal barge passage. Section 4.1 and 4.2
430 describe the test setup and the results obtained, respectively.

431 4.1 Bridge and test information

432 The tested footbridge as shown in Figure 11 (a) is a canal water overpass in Huai'an, China with
433 the span length of 115.7 m. The camera used for video acquisition was a GoPro Hero7 which
434 was mounted on a tripod in one platform of the east tower. The acquired videos were outputted
435 for post-processing using the proposed method.

436 One sample frame is given in Figure 11 (b) and the tracked target is the deck area connecting
437 the third vertical hanger, which is approximately 4 meters away from the east tower. The
438 projection transformation used is the planar homography matrix which is calibrated using the
439 six planar point correspondences (marked in Figure 11 (b)) between the image coordinates and
440 the structural coordinates. The known dimensions from the design drawing is the distance
441 between two adjacent vertical hangers (2.0 m) and the distance from the top parapet connectors

442 to the bottom hanger connectors (1.41 m). Cargo barges frequently pass through the canal
443 during the day time as shown in Figure 11 (c) which causes apparent background variations in
444 the deck target region. The video frames were converted to grayscale before performing target
445 tracking process.



447 (a)



448 (b)

449 (c)

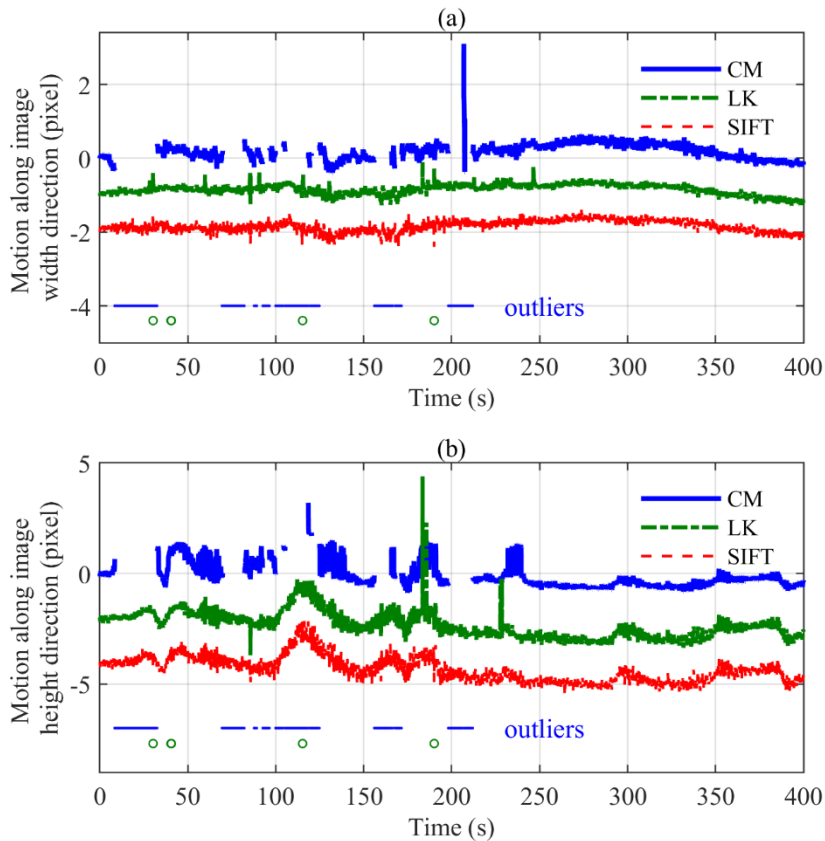
450 Figure 11 Pictures of the tested suspension footbridge (a), one sample frame with annotations
451 of target regions and control points (b), and another video frame recorded when a cargo barge
452 passed through (c).

453 4.2 Measurement results

454 To demonstrate the effectiveness of the proposed distraction-free displacement measurement
455 method, a 400-second video stream is truncated for analysis. In the beginning 250 seconds, the
456 target patterns experienced apparent background variations due to the passage of a series of
457 cargo barges tied together with mooring cables.

458 Three classic target localisation methods (CM, LK and SIFT) were implemented to predict the
459 motions with the results shown in Figure 12. For the CM, the tracked results with the similarity
460 scores lower than 0.7 are taken as tracking failure and removed as outliers while there some
461 still some local sharp peaks in the first 250 s. For the LK and SIFT, the outliers are removed by
462 applying threshold (0.1 pixel) on re-projection error after RANSAC geometric transformation.

463 The measurements by the LK and SIFT are of high similarity with the correlation coefficient of
 464 97.1%. The CM results after 250 seconds have high consistency with the outputs by the other
 465 two methods. It shows that the LK and the SIFT are robust to the partial background distractions
 466 while the CM fails to acquire reliable results. This is because the CM is area-based matching
 467 which utilises the image intensity information of the whole rectangle target region while the
 468 LK and the SIFT only track the sparse and salient key-points within the rectangle target region.



469

470 Figure 12 Image motions along two image directions (a-b) using the three classic methods
 471 (CM, LK and SIFT) for the deck target.

472 Besides the three classic methods, the Siam tracker and the proposed method (Siam+CM) were
 473 implemented to analyse the video data. The parameter settings are the same as that mentioned
 474 in section 3.2.

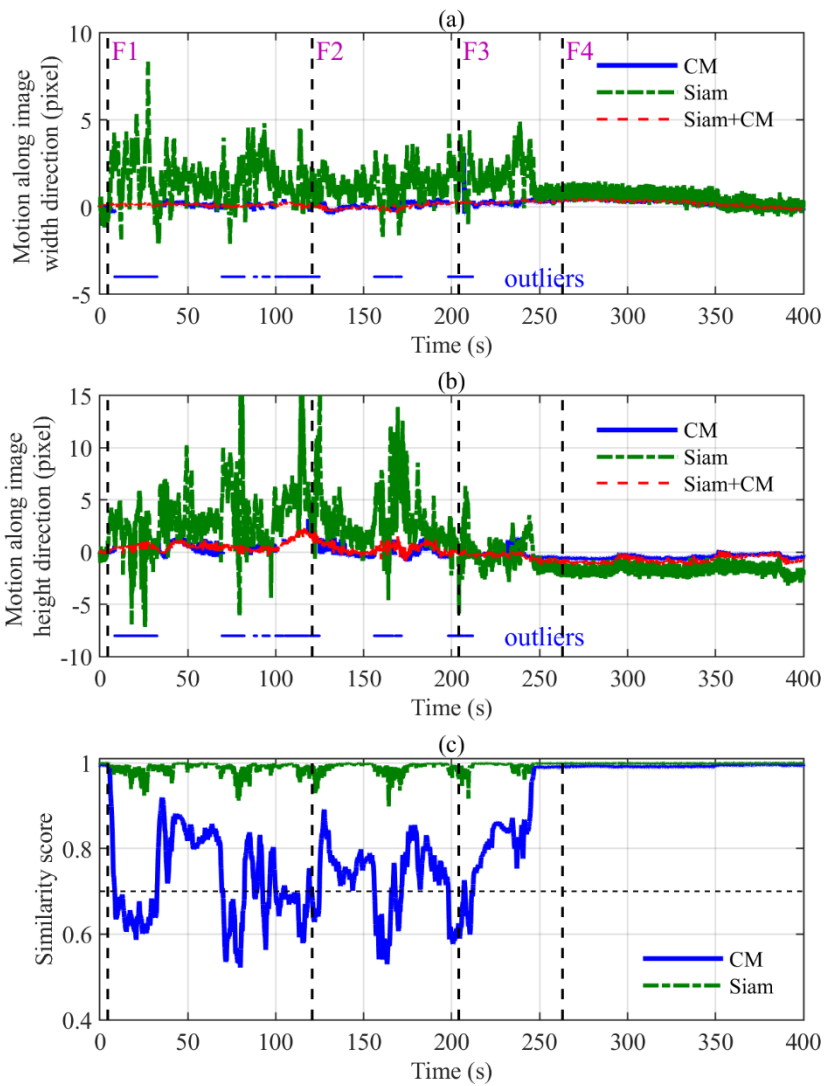
475 Figure 13 (a) and (b) shows the time histories of target motion along image width and height
 476 directions, and the similarity scores are given in Figure 13 (c). It is noted that the similarity
 477 scores in CM and Siam methods are calculated using different functions. In CM, the similarity
 478 score is measured by zero-mean normalised cross correlation coefficient of image intensity

479 value between the template and the proposal region. In Siam, the similarity score is measured
480 by the cross correlation of feature representation embedded by a modified AlexNet between the
481 template and the proposal, subtracting the weighted sum of the cross correlation in embedding
482 space between the template and the distractor set.

483 Four frames at the time steps of 4.6 s, 120.7 s, 204.0 s and 262.9 s are provided in Figure 14 for
484 the demonstration of localisation results by the three methods.

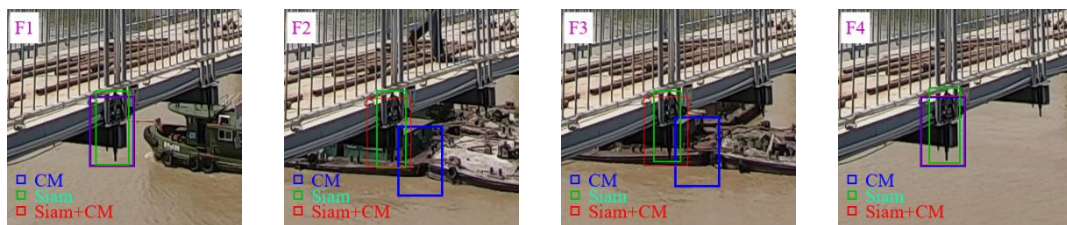
485 For the Siam, the predicted bounding boxes in four sample frames are sensible even in presence
486 of background distractions but the width and height for the bounding box is varied among video
487 frames. Therefore, the predicted target motions which are taken from the centroid coordinates
488 of the predicted bounding boxes are not accurate.

489 The proposed method (Siam+CM) is to refine the Siam output by correcting the centroid drift
490 due to bounding box size changes. The tracked results are stable within 3 pixels in both
491 directions. Compared with the measurements by the SIFT, the correlation coefficient reaches
492 96.0%. The tracked results by Siam+CM is reproduced independently in Figure 13 for better
493 visualisation. Except the deck target, the tracked motions of a stationary target on the west
494 tower is given together. It shows that the tracking resolution is within 0.2 pixel and that there
495 is no apparent drift indicating a stable mounting condition of the camera system during the
496 recording.

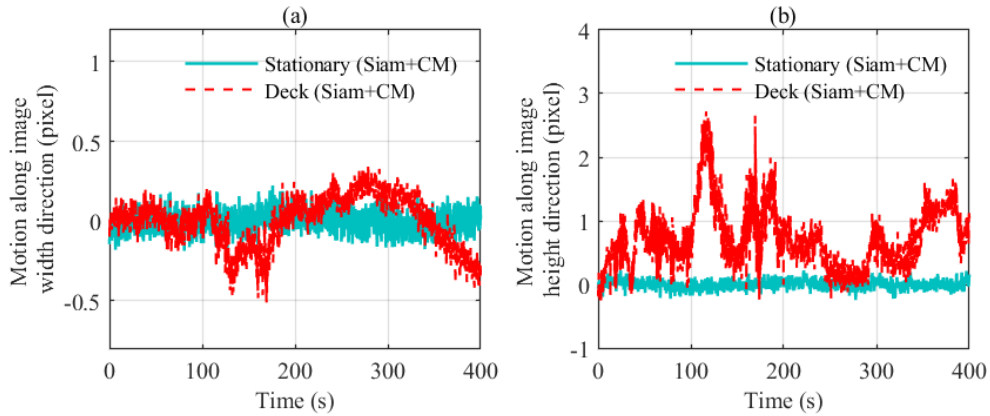


497

498 Figure 13 Image motions along two image directions (a-b) and similarity scores (c) using the
 499 three methods (CM, Siam, Siam+CM) for the deck target.

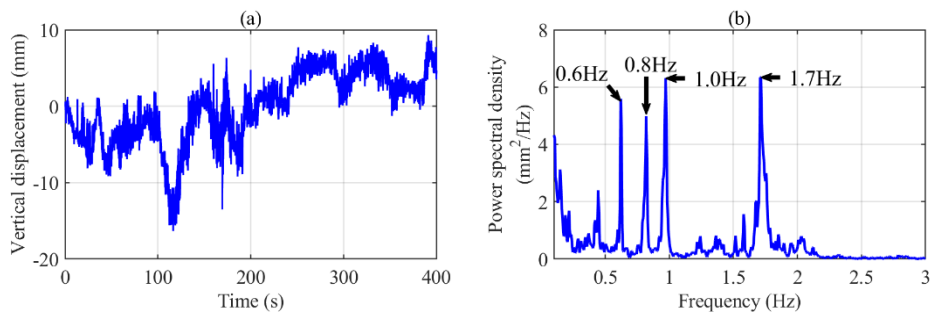


500 Figure 14 Predicted bounding boxes by the three methods in four frames F1~4.



501 Figure 15 Image motions along two image directions (a-b) using the proposed method for the
 502 stationary and deck targets.

504 The 2D structural displacement at the deck target region is estimated by transforming the
 505 image coordinates of predicted target centroids using the pre-determined planar homography
 506 matrix. The results are shown in Figure 16 with the time history and frequency information.
 507 The maximum vertical displacement is approximately 1.6 cm compared with the initial state.
 508 It occurred when a group of maintenance workers passed through. The frequency components
 509 indicate more than four peaks below 2 Hz which means the bridge serviceability could be an
 510 issue of concern.



512 Figure 16 Measured displacement time history for the deck target by the proposed method (a)
 513 and the corresponding power spectral density (c).

514 5 FIELD TEST ON A LONG-SPAN BRIDGE

515 This section describes a case study using the proposed target tracking method for measuring
 516 deformation of a long-span road bridge in normal operation. The recorded video streams were
 517 approximately 26 minutes before the sunset and the selected target region in video frames

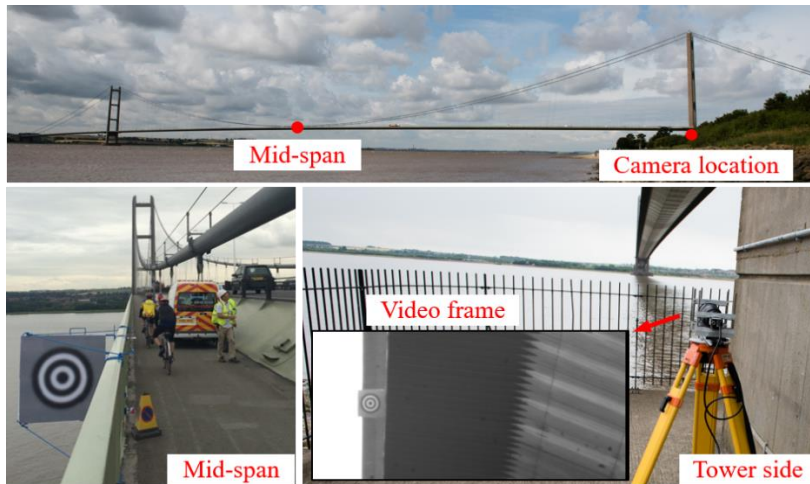
518 experienced severe patterns changes due to moving shaded area. Section 5.1 and 5.2 describe
519 the test setup and the results obtained, respectively.

520 *5.1 Bridge and test information*

521 The tested bridge is a long-span suspension bridge, the Humber Bridge in UK with a main span
522 of 1410 m. A single day of field test using the vision-based monitoring system was performed
523 to measure the displacement at mid-span of the bridge which has been reported in [57]. The
524 study in [57] was about improving GPS measurement by a data fusion method with the vision-
525 based data as data comparison. The focus here is different.

526 The vision-based monitoring system used in the test was a commercial system by Imetrum
527 Limited, UK. The camera was equipped with the lens of 300 mm focal length mounted on a
528 tripod at the base of the north tower as shown in Figure 17 (b). Essentially the camera is zoomed
529 in on the artificial target of concentric rings which has been mounted in a 1 m x 1 m metal frame
530 attached to the parapet on the east side of the bridge with on sample frame shown in Figure 17
531 (c). The acquired images are grayscale. The bridge long-term monitoring system includes two
532 GPS rovers (Leica GMX902) mounted on the main cables at mid-span and a GPS base station
533 at the bridge tower. The measured displacement by the GPS is used as reference in this study.
534 Currently, the GPS is the common choice of displacement sensing in long-span bridges [58].
535 The accuracy level of GPS data is suggested to be up to 15 mm and 35 mm for horizontal and
536 vertical measurements, respectively, at 98.5 percentile level without gross errors such as cycle
537 slip or multipath [59].

538 The acquired videos were output for post-processing using the proposed method. The planar
539 homography based on dimension correspondences in image plane and structural surface plane
540 are adopted to transform the image coordinates (i.e. pixel) to structural coordinates (i.e. mm)
541 and the known dimensions are from the width and height of the mounted metal frame. The
542 output includes the two-dimensional translations along vertical and transverse directions of the
543 bridge. A target region covering the main circular pattern on the artificial target shown in Figure
544 17 (c) is selected for analysis. The target pattern used as the template was from one frame in
545 normal condition. The video streams for analysis were of 26 min recorded approximately one
546 hour before the sunset and there are severe target pattern variations in the first 15 min.



547

548

Figure 17 Test configuration on the Humber bridge.

549

The video records include over 10 hours from a single day and most of them are under ideal

550

environmental conditions (calm day, stable light conditions and stable tripod mounting). The

551

chosen video stream for analysis includes apparent target patterns due to shade effect. A

552

previous study [60] evaluated the working performance of three tracking methods (i.e. the

553

template matching, LK optical flow, and SIFT matching) in this scenario. All of the three

554

methods failed to acquire a robust measurement. In this study, the video records with shade

555

effect is analysed by the proposed method (Siam+CM). In the first 16 min, due to the low sun

556

elevation in the west, the target panel on the east side was partially in the shadow of the bridge

557

railing and the target patterns are varied with time as shown in Figure 18 (a). Also, there was a

558

sharp pattern change in less than one second during the vehicle passage. When one tall vehicle

559

passed the mid-span of the bridge between the sun and the target, sunlight was completely

560

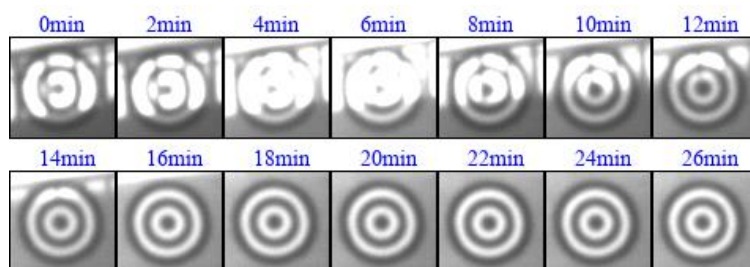
blocked, making the whole target pattern visible in the image within 0.2 seconds as shown in

561

Figure 18 (b). Then the target pattern recovered to the previous situation with parapet shades.

562

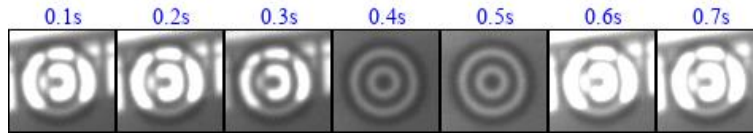
In the end 10 min, there is no apparent target pattern change.



563

564

(a)



(b)

567 Figure 18 Target pattern variations in recorded videos in 26 mins due to shade effect (a) and
 568 sharp pattern changes in 0.7 second due to vehicle passage (b).

569 5.2 Measurement results

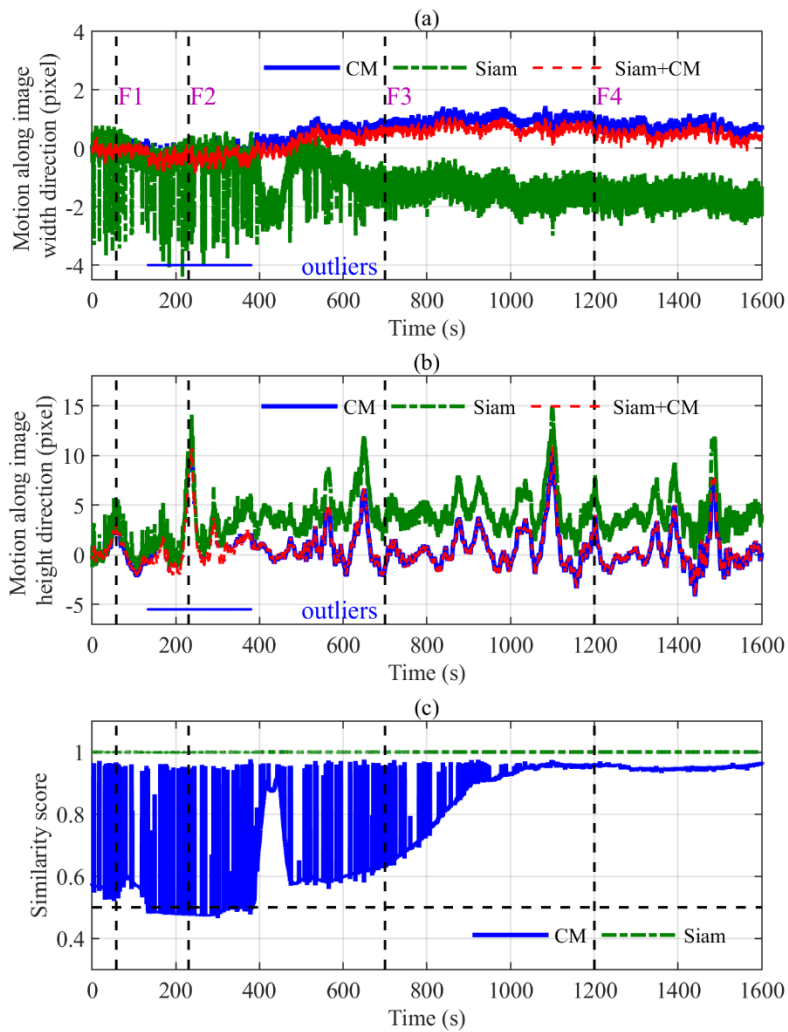
570 Three methods were implemented to analyse the video data including the CM, Siam and
 571 Siam+CM. The results are shown in Figure 19 and Figure 20.

572 The similarity scores by CM in the first 1000 seconds are mostly in a low status below 0.8 but
 573 also involve frequent jumps to a high value (over 0.9) due to vehicle passage. Some
 574 measurement outliers are observed at around 200 s as shown in Figure 20 (c) and it is observed
 575 that the circular patterns are less salient. The threshold of similarity scores for unreliable
 576 measurement evaluation is set as 0.5 through trials. The measurement after 400 s are highly
 577 similar with that by the proposed method. It indicates that from the CM is apparently not robust
 578 in presence of shading effect and apparent partial pattern variations.

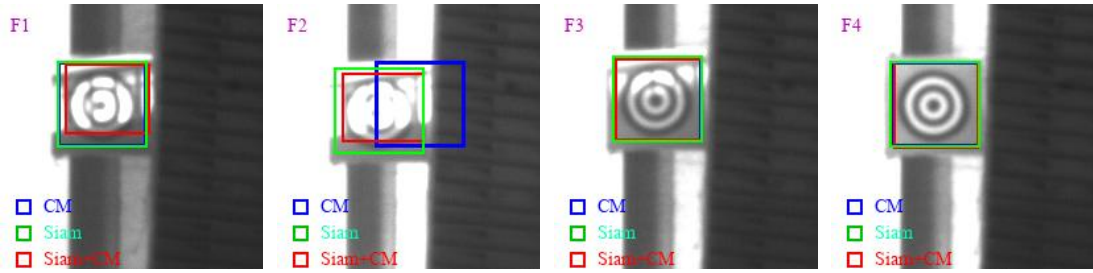
579 For the Siam, the similarity score during the whole video records is always higher than 0.90
 580 which indicates high confidence on measured results. As shown in Figure 20, the predicted
 581 bounding boxes in four sample frames are sensible but the bounding box is of varied dimensions.
 582 The predicted target motions taken from the centroid coordinates of the predicted bounding
 583 boxes are not accurate.

584 The tracked results by the proposed method (Siam+CM) are stable and the transformed vertical
 585 displacement data are presented in Figure 21 together with the GPS measurement. The cross
 586 correlation between two displacement signals is 94.7% and the RMS difference of the two
 587 measurements in 1600 seconds is 0.97 cm. The peak displacement under vehicle passage by
 588 vision-based method is 14.48 cm at 238 s and 14.72 cm at 1099 s, close to the GPS measurement
 589 (15.81 cm at 238 s and 14.60 cm at 1099 s). The previous study [57] compared the GPS and
 590 vision-based data by the CM method acquired at different time of the same day in ideal
 591 condition. The RMS difference of the two measurements is 0.75 cm, slightly smaller than the

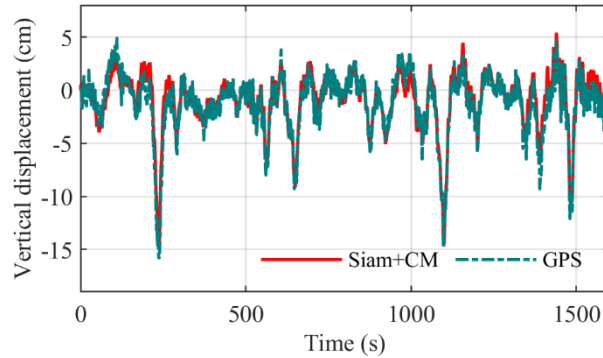
592 value (0.97 cm) in this study. Since the reference GPS displacement has limited accuracy
 593 (centimetre level [59]), the comparison study here could indicate that vision-based
 594 measurement integrating the proposed method could also reach centimetre-level accuracy in
 595 over 715 m camera-to-target distance when apparent shading and lighting changes occur. It
 596 could be a non-contact alternative to the GPS for measuring displacement data in long-span
 597 bridges.



599 Figure 19 Image motions along two image directions (a-b) and similarity scores (c) using the
 600 three methods.



601 Figure 20 Predicted bounding boxes by three methods in four frames F1~4.



602
603 Figure 21 Measured displacement time history at bridge mid-span by the proposed vision-based
604 method and the GPS.

605 6 CONCLUSIONS

606 This study proposes a novel distraction-free vision-based displacement measurement approach
607 and provides indoor and field validation tests in challenging scenarios. The main conclusions
608 are as follows:

- 609 1. The Siamese tracker (the DaSiamRPN combined with the UpdateNet) is an effective tool
610 for coarsely localising the target regions in video frames. It uses deep feature
611 representations and learned similarity measures for matching and also considers adaptive
612 template update with time. It provides convincing tracking results under illumination
613 variations. Also, it could adapt to the presence of severe target pattern changes through
614 size changes of bounding boxes and local shift to cover the adjacent area. Thus, it is
615 suitable for long-time continuous measurement.
- 616 2. The proposed method integrating the Siamese tracker with correlation-based template
617 matching inherits the advantages of the Siamese tracker and also corrects the error in the
618 Siamese tracker's output due to the image size changes of the estimated target regions. The
619 method was validated in short-range and long-range monitoring campaigns considering

620 the scenarios with severe background variations, illumination changes and shade effect,
621 providing stable and accurate displacement measurement results.

622 3. Existing applications of vision-based systems for bridge displacement measurement are
623 usually limited to short-time monitoring tests. The proposed method resolves the problem
624 of measurement robustness to environmental variations and could be potentially tied to a
625 vision-based system stably fixed on a bridge component for long-time measurement.

626 4. The study evaluated the proposed method over three application cases and further study is
627 necessary to evaluate the measurement accuracy and uncertainty for quality assurance of
628 vision-based measurement system.

629 CREDIT AUTHOR STATEMENTS

630 **Yan Xu:** Conceptualization, Methodology, Software, Validation, Investigation, Data curation,
631 Formal analysis, Methodology, Investigation, Visualization, Writing - original draft, Writing -
632 Review & Editing, Funding acquisition. **Jian Zhang:** Conceptualization, Methodology, Data
633 curation, Writing - original draft, Writing - Review & Editing. **James Brownjohn:**
634 Conceptualization, Methodology, Writing- Reviewing and Editing.

635 DECLARATION OF COMPETING INTEREST

636 The authors declare that they have no known competing financial interests or personal
637 relationships that could have appeared to influence the work reported in this paper.

638 ACKNOWLEDGEMENT

639 This research was supported by the National Key R&D Program of China (2019YFC1511102)
640 and the Jiangsu Natural Science Foundation (BK20190372). The authors express their sincere
641 appreciation for their support. Finally, the authors would like to thank the four anonymous
642 reviewers for their constructive comments.

643 REFERENCES

644 [1] H. Yoon, H. Elanwar, H. Choi, M. Golparvar-Fard, B.F. Spencer, Target-free
645 approach for vision-based structural system identification using consumer-grade

- 646 cameras, *Struct. Control Heal. Monit.* 23 (2016) 1405–1416.
647 <https://doi.org/10.1002/stc.1850>.
- 648 [2] E. Caetano, S. Silva, J. Bateira, Application of a vision system to the monitoring of
649 cable structures, in: *Seventh Int. Symp. Cable Dyn., 2007*: pp. 225–236.
- 650 [3] B.K. Oh, J.W. Hwang, Y. Kim, T. Cho, H.S. Park, Vision-based system identification
651 technique for building structures using a motion capture system, *J. Sound Vib.* 356
652 (2015) 72–85. <https://doi.org/10.1016/j.jsv.2015.07.011>.
- 653 [4] D. Feng, M.Q. Feng, Model Updating of Railway Bridge Using In Situ Dynamic
654 Displacement Measurement under Trainloads, *J. Bridg. Eng.* 20 (2015) 1–12.
655 [https://doi.org/10.1061/\(ASCE\)BE.1943-5592.0000765](https://doi.org/10.1061/(ASCE)BE.1943-5592.0000765).
- 656 [5] Y.-J. Cha, J.G. Chen, O. Büyükoztürk, Output-only computer vision based damage
657 detection using phase-based optical flow and unscented Kalman filters, *Eng. Struct.*
658 132 (2017) 300–313. <https://doi.org/10.1016/j.engstruct.2016.11.038>.
- 659 [6] T. Ojio, C.H. Carey, E.J. O'Brien, C. Doherty, S.E. Taylor, Contactless Bridge Weigh-
660 in-Motion, *J. Bridg. Eng.* 21 (2016) 04016032.
661 [https://doi.org/10.1061/\(ASCE\)BE.1943-5592.0000776](https://doi.org/10.1061/(ASCE)BE.1943-5592.0000776).
- 662 [7] G.A. Stephen, J.M.W. Brownjohn, C.A. Taylor, Measurements of static and dynamic
663 displacement from visual monitoring of the Humber Bridge, *Eng. Struct.* 15 (1993)
664 197–208.
- 665 [8] J.H.G. Macdonald, E.L. Dagless, B.T. Thomas, C.A. Taylor, Dynamic measurements
666 of the Second Severn Crossing, *Proc. Inst. Civ. Eng. - Transp.* 123 (1997) 241–248.
667 <https://doi.org/10.1680/itrans.1997.29978>.
- 668 [9] D. Feng, M. Feng, E. Ozer, Y. Fukuda, A Vision-Based Sensor for Noncontact
669 Structural Displacement Measurement, *Sensors.* 15 (2015) 16557–16575.
670 <https://doi.org/10.3390/s150716557>.
- 671 [10] X.W. Ye, Y.Q. Ni, T.T. Wai, K.Y. Wong, X.M. Zhang, F. Xu, A vision-based system
672 for dynamic displacement measurement of long-span bridges: algorithm and
673 verification, *Smart Struct. Syst.* 12 (2013) 363–379.
674 https://doi.org/10.12989/sss.2013.12.3_4.363.

- 675 [11] W.Y. Liao, W.H. Chen, Y.Q. Ni, Y. Xia, Development of a vision-based real-time
676 displacement measurement system for Guangzhou New TV Tower, in: F. Casciati,
677 Giordano M. (Eds.), Proc. 5th Eur. Work. Struct. Heal. Monit., Sorrento, Naples, Italy,
678 2010: pp. 450–455.
- 679 [12] S. Baker, I. Matthews, Lucas-Kanade 20 Years On: A Unifying Framework, *Int. J.*
680 *Comput. Vis.* 56 (2004) 221–255.
681 <https://doi.org/10.1023/B:VISI.0000011205.11775.fd>.
- 682 [13] J. Guo, C. Zhu, Dynamic displacement measurement of large-scale structures based on
683 the Lucas–Kanade template tracking algorithm, *Mech. Syst. Signal Process.* 66–67
684 (2016) 425–436. <https://doi.org/10.1016/j.ymsp.2015.06.004>.
- 685 [14] D. Feng, M.Q. Feng, Computer vision for SHM of civil infrastructure: From dynamic
686 response measurement to damage detection – A review, *Eng. Struct.* 156 (2018) 105–
687 117. <https://doi.org/10.1016/j.engstruct.2017.11.018>.
- 688 [15] S.S. Beauchemin, J.L. Barron, The computation of optical flow, *ACM Comput. Surv.*
689 27 (1995) 433–466. <https://doi.org/10.1145/212094.212141>.
- 690 [16] D. Sun, S. Roth, M.J. Black, Secrets of optical flow estimation and their principles, in:
691 *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2010: pp. 2432–2439.
692 <https://doi.org/10.1109/CVPR.2010.5539939>.
- 693 [17] M. Ehrhart, W. Lienhart, Development and evaluation of a long range image-based
694 monitoring system for civil engineering structures, in: P.J. Shull (Ed.), Proc. SPIE
695 *Struct. Heal. Monit. Insp. Adv. Mater. Aerospace, Civ. Infrastruct.*, San Diego,
696 California, United States, 2015: p. 94370K. <https://doi.org/10.1117/12.2084221>.
- 697 [18] Y.F. Ji, C.C. Chang, Nontarget Image-Based Technique for Small Cable Vibration
698 Measurement, *J. Bridg. Eng.* 13 (2008) 34–42. [https://doi.org/10.1061/\(ASCE\)1084-
699 0702\(2008\)13:1\(34\)](https://doi.org/10.1061/(ASCE)1084-0702(2008)13:1(34)).
- 700 [19] E. Caetano, S. Silva, J. Bateira, A vision system for vibration monitoring of civil
701 engineering structures, *Exp. Tech.* 35 (2011) 74–82. [https://doi.org/10.1111/j.1747-
702 1567.2010.00653.x](https://doi.org/10.1111/j.1747-1567.2010.00653.x).

- 703 [20] D.J. Fleet, A.D. Jepson, Computation of Component Image Velocity from local phase
704 information, *Int. J. Comput. Vis.* 5 (1990) 77–104.
- 705 [21] J.G. Chen, N. Wadhwa, Y.-J. Cha, F. Durand, W.T. Freeman, O. Buyukozturk, Modal
706 identification of simple structures with high-speed video using motion magnification,
707 *J. Sound Vib.* 345 (2015) 58–71. <https://doi.org/10.1016/j.jsv.2015.01.024>.
- 708 [22] Y. Yang, C. Dorn, T. Mancini, Z. Talken, G. Kenyon, C. Farrar, D. Mascareñas, Blind
709 identification of full-field vibration modes from video measurements with phase-based
710 video motion magnification, *Mech. Syst. Signal Process.* 85 (2017) 567–590.
711 <https://doi.org/10.1016/j.ymsp.2016.08.041>.
- 712 [23] J.G. Chen, A. Davis, N. Wadhwa, F. Durand, W.T. Freeman, O. Buyukozturk, Video
713 Camera-Based Vibration Measurement for Civil Infrastructure Applications, *J.*
714 *Infrastruct. Syst.* 23 (2017) 11. [https://doi.org/10.1061/\(asce\)is.1943-555x.0000348](https://doi.org/10.1061/(asce)is.1943-555x.0000348).
- 715 [24] C.Z. Dong, O. Celik, F.N. Catbas, E.J. O’Brien, S. Taylor, Structural displacement
716 monitoring using deep learning-based full field optical flow methods, *Struct.*
717 *Infrastruct. Eng.* 16 (2020) 51–71. <https://doi.org/10.1080/15732479.2019.1650078>.
- 718 [25] D.G. Lowe, Distinctive image features from scale-invariant keypoints, *Int. J. Comput.*
719 *Vis.* 60 (2004) 91–110.
- 720 [26] H. Bay, A. Ess, T. Tuytelaars, L. Van Gool, Speeded-up robust features (SURF),
721 *Comput. Vis. Image Underst.* 110 (2008) 346–359.
- 722 [27] E. Rublee, V. Rabaud, K. Konolige, G. Bradski, ORB: An efficient alternative to SIFT
723 or SURF, in: *Proc. IEEE Int. Conf. Comput. Vis., IEEE*, 2011: pp. 2564–2571.
724 <https://doi.org/10.1109/ICCV.2011.6126544>.
- 725 [28] R. Szeliski, *Computer Vision: Algorithms and Applications*, Springer Science &
726 Business Media, London, 2011. <https://doi.org/10.1007/978-1-84882-935-0>.
- 727 [29] M. Calonder, V. Lepetit, C. Strecha, P. Fua, Brief: Binary robust independent
728 elementary features, in: *Eur. Conf. Comput. Vis., Crete, Greece*, 2010: pp. 778–792.
- 729 [30] T. Khuc, F.N. Catbas, Completely contactless structural health monitoring of real-life
730 structures using cameras and computer vision, *Struct. Control Heal. Monit.* 24 (2017)
731 e1852. <https://doi.org/10.1002/stc.1852>.

- 732 [31] T. Khuc, F.N. Catbas, Computer vision-based displacement and vibration monitoring
733 without using physical target on structures, *Struct. Infrastruct. Eng.* 13 (2017) 505–
734 516. <https://doi.org/10.1080/15732479.2016.1164729>.
- 735 [32] M. Ehrhart, W. Lienhart, Monitoring of Civil Engineering Structures using a State-of-
736 the-art Image Assisted Total Station, *J. Appl. Geod.* 9 (2015) 174–182.
737 <https://doi.org/10.1515/jag-2015-0005>.
- 738 [33] Y. Xu, J.M.W. Brownjohn, Review of machine-vision based methodologies for
739 displacement measurement in civil structures, *J. Civ. Struct. Heal. Monit.* 8 (2018) 91–
740 110. <https://doi.org/10.1007/s13349-017-0261-4>.
- 741 [34] C.Z. Dong, F.N. Catbas, A review of computer vision-based structural health
742 monitoring at local and global levels, *Struct. Heal. Monit.* (2020).
743 <https://doi.org/10.1177/1475921720935585>.
- 744 [35] L. Zheng, Y. Yang, Q. Tian, SIFT Meets CNN: A Decade Survey of Instance
745 Retrieval, *IEEE Trans. Pattern Anal. Mach. Intell.* 40 (2018) 1224–1244.
746 <https://doi.org/10.1109/TPAMI.2017.2709749>.
- 747 [36] Y.J. Cha, W. Choi, O. B ü y ü k ö z t ü r k, Deep Learning-Based Crack Damage Detection
748 Using Convolutional Neural Networks, *Comput. Civ. Infrastruct. Eng.* 32 (2017) 361–
749 378. <https://doi.org/10.1111/mice.12263>.
- 750 [37] F. Hu, J. Zhao, Y. Huang, H. Li, Structure-aware 3D reconstruction for cable-stayed
751 bridges: A learning-based method, *Comput. Civ. Infrastruct. Eng.* 36 (2021) 89–108.
752 <https://doi.org/10.1111/mice.12568>.
- 753 [38] M. Dutton, W.A. Take, N.A. Hoult, Curvature Monitoring of Beams Using Digital
754 Image Correlation, *J. Bridg. Eng.* 19 (2014) 05013001.
755 [https://doi.org/10.1061/\(asce\)be.1943-5592.0000538](https://doi.org/10.1061/(asce)be.1943-5592.0000538).
- 756 [39] B.F. Spencer, V. Hoskere, Y. Narazaki, Advances in Computer Vision-Based Civil
757 Infrastructure Inspection and Monitoring, *Engineering.* 5 (2019) 199–222.
758 <https://doi.org/10.1016/j.eng.2018.11.030>.
- 759 [40] C. Koch, K. Georgieva, V. Kasireddy, B. Akinci, P. Fieguth, A review on computer
760 vision based defect detection and condition assessment of concrete and asphalt civil

- 761 infrastructure, *Adv. Eng. Informatics*. 29 (2015) 196–210.
762 <https://doi.org/10.1016/j.aei.2015.01.008>.
- 763 [41] C.V. Dung, L.D. Anh, Autonomous concrete crack detection using deep fully
764 convolutional neural network, *Autom. Constr.* 99 (2019) 52–58.
765 <https://doi.org/10.1016/j.autcon.2018.11.028>.
- 766 [42] F. Fang, L. Li, Y. Gu, H. Zhu, J.-H. Lim, A novel hybrid approach for crack detection,
767 *Pattern Recognit.* 107 (2020) 107474. <https://doi.org/10.1016/j.patcog.2020.107474>.
- 768 [43] Y. Bao, Z. Tang, H. Li, Y. Zhang, Computer vision and deep learning–based data
769 anomaly detection method for structural health monitoring, *Struct. Heal. Monit.* 18
770 (2019) 401–421. <https://doi.org/10.1177/1475921718757405>.
- 771 [44] R. Yang, S.K. Singh, M. Tavakkoli, N. Amiri, Y. Yang, M.A. Karami, R. Rai, CNN-
772 LSTM deep learning architecture for computer vision-based modal frequency
773 detection, *Mech. Syst. Signal Process.* 144 (2020).
774 <https://doi.org/10.1016/j.ymsp.2020.106885>.
- 775 [45] L. Bertinetto, J. Valmadre, J.F. Henriques, A. Vedaldi, P.H.S. Torr, Fully-
776 convolutional siamese networks for object tracking, *Eur. Conf. Comput. Vis.* (2016)
777 850–865. https://doi.org/10.1007/978-3-319-48881-3_56.
- 778 [46] J. Valmadre, L. Bertinetto, J. Henriques, A. Vedaldi, P.H.S. Torr, End-to-end
779 representation learning for Correlation Filter based tracking, *Proc. - 30th IEEE Conf.*
780 *Comput. Vis. Pattern Recognition, CVPR 2017. 2017-Janua* (2017) 5000–5008.
781 <https://doi.org/10.1109/CVPR.2017.531>.
- 782 [47] B. Li, J. Yan, W. Wu, Z. Zhu, X. Hu, High Performance Visual Tracking with Siamese
783 Region Proposal Network, *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern*
784 *Recognit.* (2018) 8971–8980. <https://doi.org/10.1109/CVPR.2018.00935>.
- 785 [48] Z. Zhu, Q. Wang, B. Li, W. Wu, J. Yan, W. Hu, Distractor-aware siamese networks
786 for visual object tracking, in: *Proc. Eur. Conf. Comput. Vis.*, 2018: pp. 101–117.
787 https://doi.org/10.1007/978-3-030-01240-3_7.
- 788 [49] L. Zhang, A. Gonzalez-Garcia, J. van de Weijer, M. Danelljan, F.S. Khan, Learning
789 the model update for siamese trackers, in: *Proc. IEEE Int. Conf. Comput. Vis.*, 2019.

- 790 [50] A. Krizhevsky, I. Sutskever, G.E. Hinton, ImageNet classification with deep
791 convolutional neural networks, in: *Commun. ACM*, 2017: pp. 84–90.
792 <https://doi.org/10.1145/3065386>.
- 793 [51] <https://github.com/foolwood/DaSiamRPN>, (n.d.).
- 794 [52] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A.
795 Karpathy, A. Khosla, M. Bernstein, A.C. Berg, L. Fei-Fei, ImageNet Large Scale
796 Visual Recognition Challenge, in: *Int. J. Comput. Vis.*, 2015: pp. 211–252.
797 <https://doi.org/10.1007/s11263-015-0816-y>.
- 798 [53] E. Real, J. Shlens, S. Mazzocchi, X. Pan, V. Vanhoucke, YouTube-BoundingBoxes: A
799 large high-precision human-annotated data set for object detection in video, *Proc. -*
800 *30th IEEE Conf. Comput. Vis. Pattern Recognition, CVPR 2017. 2017-Janua (2017)*
801 *7464–7473*. <https://doi.org/10.1109/CVPR.2017.789>.
- 802 [54] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D.
803 Ramanan, C.L. Zitnick, P. Dollár, Microsoft COCO: Common Objects in Context, in:
804 *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2015: pp. 3686–3693.
805 <http://arxiv.org/abs/1405.0312>.
- 806 [55] M. Kristan, A. Leonardis, J. Matas, M. Felsberg, R. Pflugfelder, L.Č. Zajc, T. Vojíř,
807 G. Bhat, A. Lukežič, A. Eldesokey, G. Fernández, Á. García-Martín, Á. Iglesias-Arias,
808 A.A. Alatan, A. González-García, A. Petrosino, A. Memarmoghadam, A. Vedaldi, A.
809 Muhič, A. He, A. Smeulders, A.G. Perera, B. Li, B. Chen, C. Kim, C. Xu, C. Xiong,
810 C. Tian, C. Luo, C. Sun, C. Hao, D. Kim, D. Mishra, D. Chen, D. Wang, D. Wee, E.
811 Gavves, E. Gundogdu, E. Velasco-Salido, F.S. Khan, F. Yang, F. Zhao, F. Li, F.
812 Battistone, G. De Ath, G.R.K.S. Subrahmanyam, G. Bastos, H. Ling, H.K. Galoogahi,
813 H. Lee, H. Li, H. Zhao, H. Fan, H. Zhang, H. Possegger, H. Li, H. Lu, H. Zhi, H. Li,
814 H. Lee, H.J. Chang, I. Drummond, J. Valmadre, J.S. Martin, J. Chahl, J.Y. Choi, J. Li,
815 J. Wang, J. Qi, J. Sung, J. Johnander, J. Henriques, J. Choi, J. van de Weijer, J.R.
816 Herranz, J.M. Martínez, J. Kittler, J. Zhuang, J. Gao, K. Grm, L. Zhang, L. Wang, L.
817 Yang, L. Rout, L. Si, L. Bertinetto, L. Chu, M. Che, M.E. Maresca, M. Danelljan,
818 M.H. Yang, M. Abdelpakey, M. Shehata, M. Kang, N. Lee, N. Wang, O. Miksik, P.

819 Moallem, P. Vicente-Moñivar, P. Senna, P. Li, P. Torr, P.M. Raju, Q. Ruihe, Q.
820 Wang, Q. Zhou, Q. Guo, R. Martín-Nieto, R.K. Gorthi, R. Tao, R. Bowden, R.
821 Everson, R. Wang, S. Yun, S. Choi, S. Vivas, S. Bai, S. Huang, S. Wu, S. Hadfield, S.
822 Wang, S. Golodetz, T. Ming, T. Xu, T. Zhang, T. Fischer, V. Santopietro, V. Štruc, W.
823 Wei, W. Zuo, W. Feng, W. Wu, W. Zou, W. Hu, W. Zhou, W. Zeng, X. Zhang, X.
824 Wu, X.J. Wu, X. Tian, Y. Li, Y. Lu, Y.W. Law, Y. Wu, Y. Demiris, Y. Yang, Y. Jiao,
825 Y. Li, Y. Zhang, Y. Sun, Z. Zhang, Z. Zhu, Z.H. Feng, Z. Wang, Z. He, The sixth
826 visual object tracking VOT2018 challenge results, in: Eur. Conf. Comput. Vis., 2018.
827 https://doi.org/10.1007/978-3-030-11009-3_1.

828 [56] M. Guizar-Sicairos, S.T. Thurman, J.R. Fienup, Efficient subpixel image registration
829 algorithms, *Opt. Lett.* 33 (2008) 156–158. <https://doi.org/10.1364/OL.33.000156>.

830 [57] Y. Xu, J.M.W. Brownjohn, D. Hester, K.Y. Koo, Long-span bridges: Enhanced data
831 fusion of GPS displacement and deck accelerations, *Eng. Struct.* 147 (2017) 639–651.
832 <https://doi.org/10.1016/j.engstruct.2017.06.018>.

833 [58] J.M.W. Brownjohn, K.-Y. Koo, A. Scullion, D. List, Operational deformations in
834 long-span bridges, *Struct. Infrastruct. Eng.* 11 (2015) 556–574.
835 <https://doi.org/10.1080/15732479.2014.951857>.

836 [59] A. Nickitopoulou, K. Protopsalti, S. Stiros, Monitoring dynamic and quasi-static
837 deformations of large flexible engineering structures with GPS: Accuracy, limitations
838 and promises, *Eng. Struct.* 28 (2006) 1471–1482.
839 <https://doi.org/10.1016/j.engstruct.2006.02.001>.

840 [60] Y. Xu, J.M.W. Brownjohn, Vision-based systems for structural deformation
841 measurement: Case studies, *Proc. Inst. Civ. Eng. Struct. Build.* 171 (2018) 917–930.
842 <https://doi.org/10.1680/jstbu.17.00134>.

843